

CCE: Confidence-Consistency Evaluation for Time Series Anomaly Detection

Zhijie Zhong, Zhiwen Yu*, *Senior Member IEEE*, Yiu-ming Cheung, *Fellow IEEE*, Kaixiang Yang, *Member IEEE*

Abstract—Time Series Anomaly Detection metrics serve as crucial tools for model evaluation. However, existing metrics suffer from several limitations: insufficient discriminative power, strong hyperparameter dependency, sensitivity to perturbations, and high computational overhead. This paper introduces Confidence-Consistency Evaluation (CCE), a novel evaluation metric that simultaneously measures prediction confidence and uncertainty consistency. By employing Bayesian estimation to quantify the uncertainty of anomaly scores, we construct both global and event-level confidence and consistency scores for model predictions, resulting in a concise CCE metric. Theoretically and experimentally, we demonstrate that CCE possesses strict boundedness, Lipschitz robustness against score perturbations, and linear time complexity $\mathcal{O}(n)$. Furthermore, we establish RankEval, a benchmark for comparing the ranking capabilities of various metrics. RankEval represents the first standardized and reproducible evaluation pipeline that enables objective comparison of evaluation metrics. Both CCE and RankEval implementations are fully open-source¹.

Index Terms—Time Series Anomaly Detection, Evaluation, Bayesian Estimation, Uncertainty Estimation.

I. INTRODUCTION

Time Series Anomaly Detection (TSAD) has extensive applications in industrial monitoring, cybersecurity, financial fraud detection, and other domains, aiming to identify anomalous patterns in time series data [1, 2, 3, 4, 5, 6]. With the advancement of deep learning technologies, an increasing number of state-of-the-art methods have been proposed for TSAD [7, 8, 9, 10, 11, 12, 13]. However, the current development of TSAD has reached a bottleneck, primarily due to the lag in evaluation metrics development compared to model advancement [14, 15, 16]. In early TSAD evaluations, classical metrics such as Precision, F1-score, and AUC-ROC were commonly used, but they primarily focused on the point-level detection capability of models. However, in TSAD, interval (event) anomalies are often of greater concern, making the evaluation of event-level detection capability particularly

important. To address this, an increasing number of studies have proposed interval evaluation metrics to assess event-level detection capability, including R-based F1, F1 with point adjustment (F1-PA), Reduced-F1, Affiliation F1 (Aff-F1), VUS-ROC, PATE, and Unbiased Aff-F1 (UAff-F1) scores [17, 18, 19, 20, 21, 22, 15]. However, based on our analysis and previous literature, we have identified several limitations in these metrics:

- 1) **Low Discriminability:** Previous studies [14, 15] have indicated that F1-PA [18] exhibits significant flaws when evaluating time series anomaly detection models. When faced with random scores, it fails to accurately reflect the actual performance of models, creating an illusion of progress in TSAD. Aff-F1 often yields scores greater than 0.67 in most cases, creating a false impression of good model performance.
- 2) **Hyperparameter Dependency:** F1-PA, Aff-F1, and UAff-F1 rely on anomaly threshold selection during evaluation, resulting in slow evaluation speed and limiting their use for model selection. VUS-ROC and PATE do not require anomaly thresholds but need event buffer size settings, which affects the stability and reproducibility of evaluation results and are unsuitable for scenarios with varying anomaly event lengths.
- 3) **Low Robustness:** Existing evaluation metrics are sensitive to minor perturbations in anomaly scores, leading to unstable evaluation results. F1, F1-PA, and UAff-F1 all depend on anomaly score thresholds for calculation, making them highly sensitive to minor changes in anomaly scores and resulting in unstable evaluation outcomes.
- 4) **High Computational Cost:** VUS-ROC and PATE metrics have high computational costs, limiting their application and evaluation on large-scale datasets. Typically, excessive complexity is not conducive to rapid validation of new model development.
- 5) **Accuracy-oriented:** Existing evaluation metrics primarily focus on model detection accuracy, considering only the relationship between detection results and labels. Except for threshold-independent metrics such as AUC-ROC and VUS-ROC, other metrics ignore the detection consistency issue of models, often assuming that model anomaly scores are accurate and reliable. However, in practical applications, model predictions may contain uncertainty, especially in scenarios with high data noise or complex anomaly patterns.

Furthermore, we found that current research lacks com-

Zhijie Zhong is with Pengcheng Laboratory, Shenzhen, Guangdong, 518066, China, and also with the School of Future Technology, South China University of Technology, Guangzhou, Guangdong 510650, China.

Zhiwen Yu is with the School of Computer Science and Engineering, South China University of Technology, Guangzhou, Guangdong 510650, China, and also with the Pengcheng Laboratory, Shenzhen, Guangdong 518066, China. Email: zhwyu@scut.edu.cn. Telephone number: 86-20-62893506. Fax number: 86-20-39380288.

Yiu-ming Cheung is with Department of Computer Science, Hong Kong Baptist University, Hong Kong SAR, China. E-mail: ymc@comp.hkbu.edu.hk

Kaixiang Yang is with the School of Computer Science and Engineering, South China University of Technology, Guangzhou, Guangdong 510650, China.

*Corresponding author: Zhiwen Yu.

¹Project site: <https://emorzz1g.github.io/CCE/>

prehensive and reproducible evaluation benchmarks for the metrics themselves, resulting in unclear limitations and applicability of different metrics. To address the above problems and metric limitations, we propose a novel evaluation metric called Confidence-Consistency Evaluation (CCE) and a benchmark for evaluating metrics, termed Rank-based evaluation benchmark (RankEval). To the best of our knowledge, CCE is the first metric that evaluates TSAD model performance by incorporating prediction confidence and uncertainty. We theoretically and experimentally demonstrate its robustness and computational efficiency. The fundamental motivation of CCE is that *metrics should simultaneously consider prediction score confidence and consistency, rather than merely focusing on prediction results*. Specifically, we first model anomaly scores as Beta probability distributions and quantify the uncertainty of anomaly scores based on Bayesian statistical principles. We then convert uncertainty awareness into anomaly score consistency, and finally consider both event-level and global confidence and consistency to obtain the CCE score. The workflow of CCE is illustrated in Figure 1.

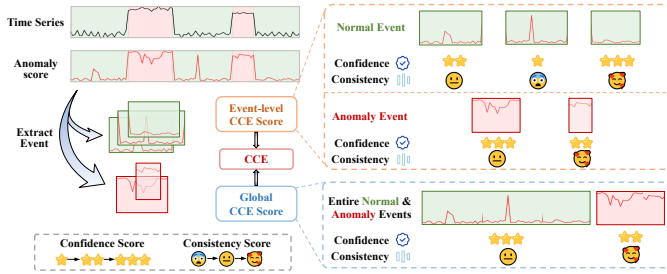


Fig. 1. Workflow of CCE framework. Note that the confidence and consistency score are continuous.

Moreover, to quantitatively and reproducibly evaluate the capabilities of different metrics, we designed the RankEval benchmark, motivated by a simple principle: *the ideal model ranking should be consistent with the metric ranking*. The RankEval benchmark includes various time series anomaly detection datasets, encompassing both synthetic and real-world datasets, while evaluating the ranking capability, robustness, and computational efficiency of metrics.

We summarize our contributions as follows:

- 1) We propose a novel confidence-consistency-based evaluation metric CCE for TSAD and a corresponding benchmark for evaluating metrics, termed RankEval, which contributes to further advancing the design of TSAD metrics and models.
- 2) Through mathematical analysis, we verify numerous excellent properties of CCE, including boundedness, robustness, and low complexity.
- 3) We validate the effectiveness of CCE on the RankEval benchmark and conduct comprehensive evaluation and analysis of multiple existing metrics. All our code is open-source², and we establish an automated standard

evaluation process for updating the RankEval benchmark³.

II. RELATED WORK

A. Point-level Evaluation Metrics

The most widely adopted point-level evaluation metrics are based on threshold-based binary classification approaches, with typical representatives including F1, precision, recall, and accuracy [23, 24, 25]. These metrics treat anomaly detection as a point-by-point binary classification problem, where the model first outputs an anomaly score sequence, then converts the scores to binary labels through thresholding, and finally performs point-by-point comparison with the ground truth labels. This approach is susceptible to noise and anomaly threshold selection. Considering the impact of anomaly threshold selection on evaluation accuracy, threshold-independent point-level metrics such as AUC-ROC and AUC-PR have also been proposed to evaluate the point-level detection capability of models. However, TSAD should focus more on the model's ability to detect entire intervals, and this evaluation approach ignores event-level assessment while overemphasizing point-level detection capability. Additionally, point-level metrics are extremely sensitive to annotation and anomaly score shifts, where slight temporal misalignment can cause dramatic fluctuations in metrics, thereby affecting evaluation reliability [15, 14, 26, 27].

B. Event-level Evaluation Metrics

To overcome the limitations of point-by-point metrics, previous research has treated continuous anomaly segments as events and proposed various event-level evaluation metrics. These metrics can be broadly categorized into three main approaches based on their underlying principles.

The first approach focuses on improving the original F1 score to directly evaluate interval anomalies. F1-PA [18] calculates F1 scores by performing point adjustment on detection results, thereby improving the interval evaluation capability of the original F1 metric. Building upon this concept, R-based F1 [17] evaluates the event-level detection capability of models by calculating the overlap rate between detection results and labels. Similarly, Reduced-F1 [19] computes F1 scores by reducing duplicate detected anomaly points, which mitigates the high-weight impact of long-interval anomalies.

The second approach emphasizes affiliation-based evaluation. Aff-F1 [20] evaluates the event-level detection capability of models by calculating the affiliation between detection results and anomaly events. To address data bias issues, UAff-F1 [15] improves upon Aff-F1 by eliminating such biases, thereby enhancing the discriminative capability of the original metric.

The third approach utilizes buffer-based evaluation strategies. VUS-ROC [21] evaluates the event-level detection capability of models by adding event buffers and calculating the area under ROC curves at different thresholds. Following a similar principle, PATE [22] also adds buffers and computes event-level precision and recall to assess model performance.

²CCE's GitHub: <https://github.com/EmorZz1G/CCE>

³RankEval: <https://emorzz1g.github.io/CCE/>

C. Uncertainty Estimation

Typically, the introduction of uncertainty helps models quantify the reliability of their predictions [28, 29, 8]. For example, AnoFormer [30] proposes using entropy to compute prediction uncertainty to clarify the decision boundary between normal and anomalous patterns. Further, TFAD [8] explores the principle of time-frequency domain uncertainty in temporal representations and proposes an anomaly detection method based on the time-frequency domain. LBAA [31] and COUTA [9] both use calibrated anomalies to enhance the uncertainty estimation capability of models. ImDiffusion [10] directly models anomaly uncertainty using diffusion models. Although uncertainty estimation has been validated in improving model detection capability, its application in the evaluation phase of time series anomaly detection remains unexplored.

III. METHODOLOGY

In this section, we first introduce the problem formulation and Bayesian uncertainty estimation method, then present the Confidence-Consistency Evaluation (CCE) Metric.

A. Problem Formulation

Definition 1 (Problem Formulation). *Given a time series $T = \{t_1, t_2, \dots, t_n\}$ with corresponding anomaly scores $S = \{s_1, s_2, \dots, s_n\}$ and ground truth labels $Y = \{y_1, y_2, \dots, y_n\}$ where $y_i \in \{0, 1\}$, where n represents the number of time points, our goal is to evaluate the quality of anomaly detection models by quantifying both prediction confidence and uncertainty consistency.*

It is worth noting that time series typically contain C channels, but our method is applicable to both univariate and multivariate time series anomaly detection, since CCE only evaluates anomaly scores.

Definition 2 (Event Segmentation). *Let $\mathcal{A} = \{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{N_a}\}$ denote the set of anomaly events and $\mathcal{E} = \{\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_{N_e}\}$ denote the set of normal events, where \mathcal{A}_i and \mathcal{E}_j represent an anomaly event and normal event, respectively. For any event, we can use p and q to denote the start and end times. For brevity, we use $\mathcal{E}_i = (p, q)$ to represent an event, and $|\mathcal{E}_i|$ to represent the duration (length) of the event.*

Theorem 1 (Uncertainty Estimation Function). *For a time series, such as anomaly scores S , the uncertainty estimation function $U : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as:*

$$U(S) = \mathbb{E}[\text{Var}(S)] \quad (1)$$

where $\text{Var}(S)$ represents the variance of the prediction scores.

B. Bayesian Uncertainty Estimation

To construct the subsequent Confidence-Consistency Evaluation Framework, we need to perform uncertainty estimation on anomaly scores. For this purpose, we propose to use Bayesian uncertainty estimation [28, 29] as our primary method due to its theoretical soundness and practical advantages. The Bayesian approach provides a principled way to

quantify uncertainty by modeling the underlying probability distribution of the prediction scores. We recommend the Bayesian approach for uncertainty estimation for three main reasons:

- 1) **Natural Boundedness:** The Beta distribution naturally constrains values to $[0, 1]$, which aligns with normalized anomaly scores. In contrast, other distributions such as Gaussian distribution estimation have no value constraints, which cannot effectively constrain anomaly scores, ultimately leading to experimental errors due to different anomaly score scales.
- 2) **Flexible Shape:** The Beta distribution can capture various shapes (symmetric, skewed, U-shaped) depending on the parameters, making it suitable for different types of anomaly patterns. In most cases, anomalies themselves are rare, which results in anomaly score distributions typically being skewed and dissimilar to Gaussian distributions.
- 3) **Conjugate Properties:** The Beta distribution has desirable mathematical properties that facilitate uncertainty propagation and combination. If we use model output p_o to represent the probability of a certain anomaly, and our prior knowledge of anomaly scores is described by a Beta distribution, then after observing new data, the posterior distribution of p_o remains a Beta distribution.

1) **Beta Distribution Modeling:** For a sequence of anomaly scores $S = \{s_1, s_2, \dots, s_n\}$ where $s_i \in [0, 1]$, we model each event (a continuous segment of time points) as a Beta distribution parameterized by α and β . This approach treats each event as a collective distribution rather than modeling individual time points separately. The following statistical measures are defined: (1) $\bar{s} = \frac{1}{n} \sum_{i=1}^n s_i$ represents the mean of the anomaly scores. (2) $m_2 = \frac{1}{n} \sum_{i=1}^n (s_i - \bar{s})^2$ represents the second central moment (variance) of the anomaly scores.

We model each sequence as a Beta distribution:

$$S \sim \text{Beta}(\alpha, \beta), \quad (2)$$

where the parameters are defined using the method of moments:

$$\begin{aligned} \alpha &= \bar{s} \left(\frac{\bar{s}(1 - \bar{s})}{m_2} - 1 \right) \\ \beta &= (1 - \bar{s}) \left(\frac{\bar{s}(1 - \bar{s})}{m_2} - 1 \right) \end{aligned} \quad (3)$$

2) **Uncertainty Estimation:** Given the Beta distribution parameters for each event, we can compute the uncertainty measure U for that event as:

$$U = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \quad (4)$$

This uncertainty measure represents the variance of the Beta distribution that models the specific event, providing an event-level measure of uncertainty for the anomaly detection model's outputs on that particular event segment.

C. Confidence-Consistency Evaluation Framework

The Confidence-Consistency Evaluation (CCE) framework is designed to assess the quality of anomaly detection models

by evaluating both the confidence of predictions and the consistency of uncertainty estimates. The framework consists of three main components: event-level scoring, global scoring, and confidence-consistency evaluation computation.

Definition 3 (Anomaly Event Confidence). *In TSAD, anomaly scores s_i are typically used to represent the probability that a time point t_i is anomalous. Therefore, for an anomaly event $\mathcal{E}_i = (p, q)$, we define the confidence of this anomaly event as:*

$$\text{Conf}(\mathcal{E}_i) = \max\left(\frac{1}{|\mathcal{E}_i|} \sum_{i=p}^q s_i - \tau, 0\right). \quad (5)$$

Here, τ represents the anomaly confidence threshold, with a default value of 0.5. The introduction of a confidence threshold is necessary because in practical applications, we typically only care about anomaly events with higher confidence levels. Furthermore, if we consider cases where the anomaly accuracy is below τ , the confidence of the anomaly event becomes:

$$\text{Conf}^*(\mathcal{E}_i) = \frac{1}{|\mathcal{E}_i|} \sum_{i=p}^q s_i - \tau. \quad (6)$$

In the subsequent text, unless otherwise specified, we will refer to the confidence with an asterisk as relaxed confidence.

Definition 4 (Normal Event Confidence). *Similarly, for normal events $\mathcal{A}_j = (p, q)$, we define the normal event confidence as:*

$$\text{Conf}(\mathcal{A}_j) = \max\left(1 - \tau - \frac{1}{|\mathcal{A}_j|} \sum_{j=p}^q s_j, 0\right). \quad (7)$$

Analogously, here $1 - \tau$ represents the normal confidence threshold, meaning that the thresholds for normal and anomaly confidence are symmetric about 0.5. If we consider cases where the accuracy is below $1 - \tau$, the confidence of the normal event becomes:

$$\text{Conf}^*(\mathcal{A}_j) = 1 - \tau - \frac{1}{|\mathcal{A}_j|} \sum_{j=p}^q s_j. \quad (8)$$

Definition 5 (Prediction Consistency). *For each event \mathcal{E}_i or \mathcal{A}_j , we define the prediction consistency as:*

$$\begin{aligned} \text{Cons}(\mathcal{E}_i) &= \exp\left(-\frac{1}{|\mathcal{E}_i|} \sum_{k=p}^q U_k\right) \\ \text{Cons}(\mathcal{A}_j) &= \exp\left(-\frac{1}{|\mathcal{A}_j|} \sum_{k=p}^q U_k\right). \end{aligned} \quad (9)$$

Where U_k is the uncertainty of the score at time k . $\text{Cons}(\mathcal{E}_i)$ and $\text{Cons}(\mathcal{A}_j)$ represent the confidence consistency of anomaly events and normal events, respectively. Intuitively, if an event has higher uncertainty, its consistency is lower, indicating that the model's prediction for that event is less reliable.

1) Event-Level CCE Scoring: To evaluate the assessment accuracy of the model for each anomaly event and normal event, we need to compute the CCE score for each event. For anomaly events $\mathcal{E}_i = (p, q)$ and normal events $\mathcal{A}_j = (p, q)$, we define their event-level scores as follows.

For each anomaly event $\mathcal{E}_i = (p, q)$, we compute the anomaly event score:

$$S_{\text{anom}}(\mathcal{E}_i) = \text{Conf}(\mathcal{E}_i) \times \text{Cons}(\mathcal{E}_i) \quad (10)$$

Similarly, for normal events $\mathcal{A}_j = (p, q)$, we compute the normal event score:

$$S_{\text{norm}}(\mathcal{A}_j) = \text{Conf}(\mathcal{A}_j) \times \text{Cons}(\mathcal{A}_j) \quad (11)$$

Definition 6 (Event-Level Score). *The event-level score for anomaly and normal events, without loss of generality, is defined as:*

$$S_{\text{event}} = \alpha \bar{S}_{\text{anom}} + (1 - \alpha) \bar{S}_{\text{norm}}. \quad (12)$$

Where \bar{S}_{anom} and \bar{S}_{norm} are the average scores across all anomaly and normal events respectively, and α is a weight parameter used to balance anomaly and normal events, typically set to 0.5, which can be adjusted based on detection requirements in practical scenarios.

2) Global CCE Scoring: To evaluate the model's performance across the entire time series, we design a global scoring mechanism. We treat all normal events across the entire time series as a comprehensive normal event, and all anomaly events across the entire time series as a comprehensive anomaly event. Specifically, $S_{\text{anom}} = \{s_i \mid i \in \{\mathcal{A}\}\}$, where $\{\mathcal{A}\}$ represents the set of time points for all anomaly events. Similarly, $S_{\text{norm}} = \{s_j \mid j \in \{\mathcal{E}\}\}$, where $\{\mathcal{E}\}$ represents the set of time points for all normal events.

Definition 7 (Global Score). *The global CCE score for anomaly and normal events is defined as:*

$$S_{\text{global}} = \eta S_{\text{anom}}^{\text{global}} + (1 - \eta) S_{\text{norm}}^{\text{global}} \quad (13)$$

where $S_{\text{anom}}^{\text{global}}$ and $S_{\text{norm}}^{\text{global}}$ are the global scores for anomaly and normal events respectively.

Global scores for anomaly and normal events are computed as follows:

$$\begin{aligned} S_{\text{anom}}^{\text{global}} &= \text{Conf}(S_{\text{anom}}) \times \text{Cons}(S_{\text{anom}}) \\ S_{\text{norm}}^{\text{global}} &= \text{Conf}(S_{\text{norm}}) \times \text{Cons}(S_{\text{norm}}) \end{aligned} \quad (14)$$

Definition 8 (Confidence-Consistency Evaluation). *The final confidence-consistency evaluation score is defined as:*

$$S_{\text{CCE}} = S_{\text{event}} + S_{\text{global}} \quad (15)$$

This score combines both event-level and global scores to provide a comprehensive evaluation of the model's performance. Finally, we present the pseudo-code implementation of CCE in Algorithm 1.

Algorithm 1 Confidence-Consistency Evaluation

Require: Ground truth labels \mathbf{Y} , anomaly scores \mathbf{S} , confidence thresholds τ , scale parameter γ

Ensure: CCE score S_{CCE}

- 1: Normalize scores: $\mathbf{S} = \frac{\mathbf{S} - \min(\mathbf{S})}{\max(\mathbf{S}) - \min(\mathbf{S})}$
- 2: Extract events: $\mathcal{A}, \mathcal{E} = \text{extract_events}(\mathbf{Y})$
- 3: Compute Beta parameters and uncertainties for each event using Eq. 3
- 4: Compute event-level CCE scores S_{event} using Definitions III-C, 4, 5, and 12
- 5: Compute global CCE score S_{global} using Definitions 7 and 5
- 6: Compute final CCE score: $S_{\text{CCE}} = S_{\text{event}} + S_{\text{global}}$
- 7: **return** S_{CCE}

IV. THEORY ANALYSIS

In this section, we analyze the theoretical properties and computational complexity of CCE. CCE possesses the following properties, where Properties 1, 2, and 6 are straightforward to establish, while Properties 3, 4, and 5 are proven in this section.

Property 1 (Intuition): CCE quantifies the model's ability to distinguish between normal and anomalous events by evaluating both confidence and consistency. High confidence and high consistency indicate that the model's predictions for normal and anomalous events are more reliable.

Property 2 (Interpretability): CCE provides interpretable scores that can be used to identify model weaknesses and areas for improvement. By analyzing event-level and global scores, we can better understand the model's performance on different types of events, thereby guiding the direction of model improvement.

Property 3 (Boundedness): The final score $S_{\text{CCE}} \in [0, 1]$, where 0 indicates poor performance and 1 indicates excellent performance. If we remove constraints to confidence, the CCE score has a lower bound of -1 and an upper bound of 1.

Property 4 (Robustness): CCE is robust to small perturbations in anomaly scores due to its Lipschitz continuity. This means that when anomaly scores undergo minor changes due to noise or computational errors, the change in the CCE metric will be strictly controlled within L times the magnitude of the perturbation.

Property 5 (Scalability): The computational complexity of CCE is $\mathcal{O}(n)$, where n is the length of the time series data. CCE can be efficiently computed for large datasets, making it suitable for large-scale anomaly detection tasks.

Property 6 (Symmetry): The evaluation is symmetric with respect to anomaly and normal events, ensuring fair assessment. When $\eta = 0.5$, this property guarantees that the model's evaluation of anomaly and normal events is balanced, without bias toward any particular type of event.

A. Proof of CCE Boundedness

Theorem 2. The variance of a Beta distribution has an upper bound of $\frac{1}{4}$.

Proof. We prove that the variance of a Beta distribution has an upper bound, with a maximum value of $\frac{1}{4}$.

The probability density function of a Beta distribution depends on two shape parameters $\alpha > 0$ and $\beta > 0$, and its variance formula is:

$$\text{Var}(\mathbf{S}) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}, \quad (16)$$

where \mathbf{S} is a random variable following a Beta distribution, i.e., $\mathbf{S} \sim \text{Beta}(\alpha, \beta)$.

First, let $n = \alpha + \beta$ and $p = \frac{\alpha}{n}$ (where $\beta = n(1 - p)$ and $0 < p < 1$). The variance formula can be rewritten as:

$$\text{Var}(\mathbf{S}) = \frac{np \cdot n(1 - p)}{n^2(n + 1)} = \frac{p(1 - p)}{n + 1} \quad (17)$$

For a fixed n , the maximum value of $p(1 - p)$ is $\frac{1}{4}$ (achieved when $p = \frac{1}{2}$). At this point, the variance simplifies to $\frac{1}{4(n + 1)}$, which clearly decreases as n increases.

To find the global maximum, we need to consider the case when n approaches its minimum value. Since $\alpha > 0$ and $\beta > 0$, we have $n > 0$. When $n \rightarrow 0^+$ (i.e., both α and β approach 0), the variance approaches $\frac{1}{4}$.

Therefore, regardless of how α and β are chosen (as long as they are positive), the variance of the Beta distribution will never exceed $\frac{1}{4}$. \square

Corollary 3. Both confidence and consistency scores are bounded.

Proof. For the event-level confidence of anomaly events, we have: $\text{Conf}(\mathcal{A}_i) = \max\left(\frac{1}{|\mathcal{A}_i|} \sum_{k=p}^q s_k - \tau, 0\right) \in [0, 1 - \tau]$. If we do not consider applying positive value constraints to confidence, then $\text{Conf}^*(\mathcal{A}_i) \in [-\tau, 1 - \tau]$.

Similarly, for the event-level confidence of normal events, we have: $\text{Conf}(\mathcal{E}_j) = \max\left(1 - \frac{1}{|\mathcal{E}_j|} \sum_{k=p}^q s_k - \tau, 0\right) \in [0, 1 - \tau]$. If we do not consider applying positive value constraints to confidence, then $\text{Conf}^*(\mathcal{E}_j) \in [-\tau, 1 - \tau]$.

Since the uncertainty $U_k \in [0, 0.25]$ (the upper bound of Beta distribution variance is $\frac{1}{4}$), the event-level consistency satisfies:

$$\begin{aligned} \text{Cons}(\mathcal{A}_i) &= \exp\left(-\frac{1}{|\mathcal{A}_i|} \sum_{k=p}^q U_k\right) \in [e^{-0.25}, 1] \\ \text{Cons}(\mathcal{E}_j) &\in [e^{-0.25}, 1] \end{aligned} \quad (18)$$

Similarly, global confidence and consistency also satisfy the above bounds. \square

Theorem 4. S_{event} and S_{global} are bounded.

Proof. If we consider applying positive value constraints to confidence, then:

$$\begin{aligned} S_{\text{anom}}(\mathcal{A}_i) &= \text{Conf}^*(\mathcal{A}_i) \times \text{Cons}(\mathcal{A}_i) \in [0, 1 - \tau] \\ S_{\text{norm}}(\mathcal{E}_j) &= \text{Conf}^*(\mathcal{E}_j) \times \text{Cons}(\mathcal{E}_j) \in [0, 1 - \tau] \end{aligned} \quad (19)$$

If we consider relaxing the constraints on confidence, then:

$$\begin{aligned} S_{\text{anom}}^*(\mathcal{A}_i) &= \text{Conf}(\mathcal{A}_i) \times \text{Cons}(\mathcal{A}_i) \in [-\tau, 1 - \tau] \\ S_{\text{norm}}^*(\mathcal{E}_j) &= \text{Conf}(\mathcal{E}_j) \times \text{Cons}(\mathcal{E}_j) \in [-\tau, 1 - \tau] \end{aligned} \quad (20)$$

Typically, we expect the model's confidence to be greater than 0.5, so we can assume $\tau \in [0.5, 1]$. At this point, we have:

$$S_{\text{anom}}(\mathcal{A}_i), S_{\text{norm}}(\mathcal{E}_j) \in [0, 0.5] \quad (21)$$

For relaxed confidence S_{anom}^* and S_{norm}^* , the maximum and minimum values are related to their corresponding confidence thresholds, but they always satisfy:

$$S_{\text{anom}}^*(\mathcal{A}_i), S_{\text{norm}}^*(\mathcal{E}_j) \in [-0.5, 0.5] \quad (22)$$

Since S_{event} and S_{global} are simple weighted functions, we have:

$$\begin{aligned} S_{\text{event}}, S_{\text{global}} &\in [0, 0.5] \\ S_{\text{event}}^*, S_{\text{global}}^* &\in [-0.5, 0.5] \end{aligned} \quad (23)$$

□

Corollary 5. *CCE is bounded, with $S_{\text{CCE}} \in [0, 1]$ and $S_{\text{CCE}}^* \in [-1, 1]$.*

Since CCE is the sum of S_{event} and S_{global} , this conclusion is straightforward to obtain.

B. Proof of CCE Robustness

Since the entire process of TSAD often involves a large amount of noise, where the most critical factor is the calculation error of anomaly scores, we need to prove that the CCE metric is stable against perturbations in anomaly scores. This ensures that when anomaly scores undergo minor changes due to noise or computational errors, the change in the CCE metric will be strictly controlled within L times the magnitude of the perturbation.

Theorem 6 (CCE satisfies Lipschitz continuity against perturbations in anomaly scores). *According to the assumption of Beta distribution modeling for anomaly scores, anomaly scores are normalized to $[0, 1]$. Let $\mathbf{s} = (s_1, \dots, s_n)$ and $\tilde{\mathbf{s}} = (\tilde{s}_1, \dots, \tilde{s}_n)$ satisfy $\|\mathbf{s} - \tilde{\mathbf{s}}\|_2 \leq \delta$. For S_{CCE} , there exists a constant $L > 0$, and L is close to 1, such that*

$$|S_{\text{CCE}}(\mathbf{s}) - S_{\text{CCE}}(\tilde{\mathbf{s}})| \leq L\delta.$$

Proof. We prove that the CCE metric satisfies Lipschitz continuity against perturbations in anomaly scores. Since CCE is composed of event-level scores and global scores, we need to prove that both event-level scores and global scores satisfy Lipschitz continuity against perturbations in anomaly scores. For any anomaly event \mathcal{E}_i (interval $[p, q]$), we explicitly calculate $|\text{Conf}(\mathcal{E}_i) - \text{Conf}(\tilde{\mathcal{E}}_i)|$ and $|\text{Cons}(\mathcal{E}_i) - \text{Cons}(\tilde{\mathcal{E}}_i)|$ respectively, then combine them to obtain the Lipschitz constant of event-level scores, and finally provide the global L .

First, for the difference in confidence between anomaly events:

$$\begin{aligned} |\text{Conf}(\mathcal{E}_i) - \text{Conf}(\tilde{\mathcal{E}}_i)| &\leq \left| \frac{1}{|\mathcal{E}_i|} \sum_{k=p}^q (s_k - \tilde{s}_k) \right| \\ &\leq \frac{1}{|\mathcal{E}_i|} \sum_{k=p}^q |s_k - \tilde{s}_k| \leq \delta. \end{aligned} \quad (24)$$

Similarly for normal events, the upper bound of the difference is still δ .

For each event \mathcal{E}_i , the uncertainty U is computed from Beta distribution parameters. The sensitivity of uncertainty to individual score changes can be bounded by analyzing the partial derivatives with respect to the event-level statistics.

Now we compute the sensitivity of uncertainty to individual score changes. For any time point $k \in [p, q]$, we have:

$$\begin{aligned} \frac{\partial \bar{s}_{\mathcal{E}_i}}{\partial s_k} &= \frac{1}{|\mathcal{E}_i|} \\ \frac{\partial m_{2, \mathcal{E}_i}}{\partial s_k} &= \frac{2}{|\mathcal{E}_i|} (s_k - \bar{s}_{\mathcal{E}_i}) \end{aligned} \quad (25)$$

Since $|s_k - \bar{s}_{\mathcal{E}_i}| \leq 1$ (as $s_k \in [0, 1]$), we have:

$$\left| \frac{\partial m_{2, \mathcal{E}_i}}{\partial s_k} \right| \leq \frac{2}{|\mathcal{E}_i|} \quad (26)$$

Using the chain rule and the fact that U is a smooth function of α and β , and α, β are smooth functions of $\bar{s}_{\mathcal{E}_i}$ and m_{2, \mathcal{E}_i} , we can bound the partial derivative of U with respect to s_k :

$$\left| \frac{\partial U}{\partial s_k} \right| \leq \frac{C_1}{|\mathcal{E}_i|} \quad (27)$$

where C_1 is a constant that depends on the range of α and β values. Since $\alpha, \beta > 0$ and are bounded by the event size and score statistics, C_1 is finite.

Therefore, for any event \mathcal{E}_i , the change in uncertainty satisfies:

$$|U(\mathbf{s}_{\mathcal{E}_i}) - U(\tilde{\mathbf{s}}_{\mathcal{E}_i})| \leq \frac{C_1}{|\mathcal{E}_i|} \|\mathbf{s}_{\mathcal{E}_i} - \tilde{\mathbf{s}}_{\mathcal{E}_i}\|_2 \leq \frac{C_1}{|\mathcal{E}_i|} \delta \quad (28)$$

Setting $L_U = \frac{C_1}{|\mathcal{E}_i|}$, we have:

$$|U(\mathbf{s}_{\mathcal{E}_i}) - U(\tilde{\mathbf{s}}_{\mathcal{E}_i})| \leq L_U \delta \quad (29)$$

To compute the exact value of C_1 , we analyze the partial derivatives more carefully. Let us define $n = \alpha + \beta$ and $p = \frac{\alpha}{n}$, then:

$$U = \frac{\alpha\beta}{n^2(n+1)} = \frac{p(1-p)}{n+1} \quad (30)$$

The partial derivative with respect to s_k is:

$$\begin{aligned} \frac{\partial U}{\partial s_k} &= \frac{\partial U}{\partial p} \cdot \frac{\partial p}{\partial s_k} + \frac{\partial U}{\partial n} \cdot \frac{\partial n}{\partial s_k} \\ &= \frac{1-2p}{n+1} \cdot \frac{\partial p}{\partial s_k} - \frac{p(1-p)}{(n+1)^2} \cdot \frac{\partial n}{\partial s_k} \end{aligned} \quad (31)$$

Since $|p| \leq 1$ and $|1-2p| \leq 1$, we have:

$$\left| \frac{\partial U}{\partial s_k} \right| \leq \frac{1}{n+1} \cdot \left| \frac{\partial p}{\partial s_k} \right| + \frac{1}{(n+1)^2} \cdot \left| \frac{\partial n}{\partial s_k} \right| \quad (32)$$

From the definitions of α and β , we can show that:

$$\left| \frac{\partial p}{\partial s_k} \right| \leq \frac{2}{|\mathcal{E}_i|} \quad \text{and} \quad \left| \frac{\partial n}{\partial s_k} \right| \leq \frac{2}{|\mathcal{E}_i|} \quad (33)$$

Therefore:

$$\left| \frac{\partial U}{\partial s_k} \right| \leq \frac{2}{|\mathcal{E}_i|(n+1)} + \frac{2}{|\mathcal{E}_i|(n+1)^2} \quad (34)$$

To find a uniform upper bound that works for all possible values of n , we need to analyze the function $f(n) = \frac{2}{n+1} + \frac{2}{(n+1)^2}$ for $n > 0$.

The derivative of $f(n)$ with respect to n is:

$$f'(n) = -\frac{2}{(n+1)^2} - \frac{4}{(n+1)^3} < 0 \quad (35)$$

This shows that $f(n)$ is a decreasing function of n . The maximum value occurs at the minimum possible value of n .

From the Beta distribution properties, we know that $\alpha, \beta > 0$, which implies $n = \alpha + \beta > 0$. However, for meaningful uncertainty estimation in practice, we typically have $n \geq 1$. At $n = 1$, we have:

$$f(1) = \frac{2}{2} + \frac{2}{4} = 1 + 0.5 = 1.5 \stackrel{\text{def}}{=} C_1. \quad (36)$$

Thus,

$$|U_k - \tilde{U}_k| \leq L_U |s_k - \tilde{s}_k| \leq L_U \delta. \quad (37)$$

Using the 1-Lipschitz property of \exp (on $[0, 1]$),

$$\begin{aligned} |\text{Cons}(\mathcal{E}_i) - \text{Cons}(\tilde{\mathcal{E}}_i)| &\leq \left| \frac{1}{\ell_i} \sum_{k=p}^q (U_k - \tilde{U}_k) \right| \\ &\leq \frac{1}{\ell_i} \sum_{k=p}^q L_U \delta = L_U \delta. \end{aligned} \quad (38)$$

Event-level scores are the product of two terms, i.e., $S_i = \text{Conf}(\mathcal{E}_i) \cdot \text{Cons}(\mathcal{E}_i)$. Since $\text{Conf}(\mathcal{E}_i) \leq 1 - \tau$ and $\text{Cons}(\mathcal{E}_i) \leq 1$, we use the product difference inequality $|ab - \tilde{a}\tilde{b}| \leq |a||b - \tilde{b}| + |\tilde{b}||a - \tilde{a}|$, to obtain:

$$|S_i - \tilde{S}_i| \leq (1 - \tau) \cdot L_U \delta + 1 \cdot \delta \leq (0.5L_U + 1)\delta.$$

Subsequently, we can obtain the upper bound of the difference between normal and anomaly event-level scores:

$$\begin{aligned} |\bar{S}_{\text{norm}} - \tilde{\bar{S}}_{\text{norm}}| &\leq (0.5L_U + 1)\delta, \\ |\bar{S}_{\text{anom}} - \tilde{\bar{S}}_{\text{anom}}| &\leq (0.5L_U + 1)\delta. \end{aligned} \quad (39)$$

CCE scores are obtained by summing up event-level scores and global scores, i.e., $S_{\text{CCE}} = \alpha \bar{S}_{\text{anom}} + (1 - \alpha) \bar{S}_{\text{norm}}$. Therefore, the perturbation amount of CCE scores is:

$$\begin{aligned} |S_{\text{CCE}} - \tilde{S}_{\text{CCE}}| &\leq \alpha(0.5L_U + 1)\delta + (1 - \alpha)(0.5L_U + 1)\delta \\ &= (0.5L_U + 1)\delta. \end{aligned} \quad (40)$$

This shows that the CCE metric satisfies Lipschitz continuity against perturbations in anomaly scores, with the Lipschitz constant: $L = 0.5L_U + 1$.

Substituting $L_U = \frac{C_1}{|\mathcal{E}_i|} = \frac{1.5}{|\mathcal{E}_i|}$, we get:

$$L = \frac{0.75}{|\mathcal{E}_i|} + 1 \quad (41)$$

For typical event sizes, we have:

- For events with length $|\mathcal{E}_i| = 1$: $L = \frac{0.75}{1} + 1 = 1.75$
- For events with length $|\mathcal{E}_i| = 2$: $L = \frac{0.75}{2} + 1 = 1.375$
- For events with length $|\mathcal{E}_i| = 5$: $L = \frac{0.75}{5} + 1 = 1.15$
- For events with length $|\mathcal{E}_i| = 10$: $L = \frac{0.75}{10} + 1 = 1.075$
- For events with length $|\mathcal{E}_i| \geq 20$: $L \leq 1.0375$

As can be seen from Table S5, the average length of events in real-world scenarios is generally much greater than 20; therefore, $L \leq 1.0375$.

This ensures the stability of the CCE metric when there are errors in calculating anomaly scores. The Lipschitz constant is always close to 1, indicating that the CCE metric is robust to small perturbations in anomaly scores. \square

C. Computational Complexity

To strictly analyze the computational complexity of CCE, we decompose it into the following six main steps, and provide the time complexity upper bound for each step.

Theorem 7 (CCE Complexity). *Given a time series of length n , $\mathcal{O}_{\text{CCE}} = \mathcal{O}(n)$.*

Proof. We decompose the calculation of CCE into the following steps:

- 1) **Event Extraction:** According to the given labels \mathbf{Y} , we extract all anomaly events and normal events. This step can be completed through a single linear scan: $\mathcal{O}(n)$
- 2) **Beta Parameter Computation:** For each anomaly score s_i , we calculate its corresponding Beta distribution parameters α_i and β_i . Since the calculation of each score is independent, the complexity is: $\mathcal{O}(n)$
- 3) **Uncertainty Estimation:** For each anomaly score, we calculate the variance U_i of its Beta distribution. Similarly, since each variance calculation is independent and requires only constant time: $\mathcal{O}(n)$
- 4) **Event-Level Scoring:** For each event (anomaly or normal), we calculate its confidence and consistency scores. The complexity of each calculation is proportional to the length of the event, but the total length of all events is at most n , therefore: $\mathcal{O}(n)$
- 5) **Global Scoring:** We calculate the global confidence and consistency scores for the entire time series. This step only requires summing up all anomaly scores and normal scores respectively, therefore: $\mathcal{O}(n)$
- 6) **CCE Score Computation:** We perform weighted summation of event-level scores and global scores, which requires only constant time: $\mathcal{O}(1)$

In summary, all steps of CCE can be completed within $\mathcal{O}(n)$ time, therefore: $\mathcal{O}_{\text{CCE}} = \mathcal{O}(n)$. \square

V. EXPERIMENTS

In this section, we first introduce how to use the RankEval benchmark to evaluate different metrics. To provide a clear overview, the experiments are designed to address the following research questions (**RQs**):

- **RQ1:** How efficient are different metrics in terms of computational performance (Sec. V-B)?
- **RQ2:** Does CCE have effectiveness and robustness compared to other metrics (Sec. V-C)?
- **RQ3:** How do different metrics and models perform in terms of visualization (Sec. V-D)?
- **RQ4:** Do the hyperparameters and task settings of CCE affect its evaluation capability (Sec. V-E)?

A. Experimental Setup

In this section, we introduce the datasets, TSAD models, and evaluation methods used in RankEval.

1) **Datasets**: To ensure reproducibility of results, RankEval uses widely adopted and publicly available real-world datasets in the TSAD field, including MSL, SMD, PSM, SWAT, Creditcard, and UCR, among others [15, 32]. Additionally, to verify the performance of different metrics under theoretical conditions, we constructed 30 synthetic datasets. The time series properties of these datasets are elaborated in Table S5 in Appendix.

2) **TSAD Models**: To analyze the performance of different metrics in real-world scenarios, we used two classic machine learning models (LOF [33], IForest [34]) and five deep learning models (LSTMAD, USAD [35], AnomalyTransformer (A.T.) [36], TimesNet [37], Donut). We used the public implementations of these models, and all models used default hyperparameter settings to ensure fairness and consistency of evaluation.

Additionally, to analyze the performance of different metrics under theoretical conditions, we designed three different anomaly score generation models (ASGM) and their corresponding noise perturbation versions. For the default anomaly score generation models, we named them AccQ, LowDisAccQ, and PreQ-NegP. For the noise perturbation versions of anomaly score generation models (ASGM-R), we named them AccQ-R, LowDisAccQ-R, and PreQ-NegP-R. These ASGM-R models can be used to evaluate the robustness of different metrics. The specific ASGM-R models are as follows:

- 1) **AccQ**: The accuracy of this model is fixed at q . If the true label is anomalous, then when correctly predicted (with probability q), the anomaly score is $s^{(1)} = 0.9 + 0.1\mathcal{U}$; when incorrectly predicted (with probability $1-q$), the anomaly score is $s^{(2)} = 0.05\mathcal{U}$. If the true label is normal, then when correctly predicted, the anomaly score is $s^{(1)}$; when incorrectly predicted, the anomaly score is $s^{(2)}$. Here, \mathcal{U} represents the standard uniform distribution.
- 2) **LowDisAccQ**: Similar to AccQ, this model's accuracy is fixed at q , but the discriminative power of anomaly scores is lower: $s^{(1)} = 0.6 + 0.1\mathcal{U}$, $s^{(2)} = 0.4\mathcal{U}$.
- 3) **PreQ-NegP**: This model's anomaly precision is fixed at q , and the false positive rate is fixed at p . The model's default anomaly score is $0.1\mathcal{U}$. For anomalies, with probability q , the anomaly score is $0.9\mathcal{U} + 0.1$; for normal cases, with probability p , the anomaly score is $0.9\mathcal{U} + 0.1$.
- 4) **AccQ-R, LowDisAccQ-R, PreQ-NegP-R**: For the above three anomaly score generation methods, we add Gaussian noise with mean 0 and standard deviation σ to obtain the corresponding noise perturbation versions.

Subsequently, when no ambiguity arises, we use the ASGM names to refer to the corresponding task types.

3) **Evaluation Method**: To evaluate the capability of different metrics in anomaly detection tasks, we adopt a ranking-based evaluation method, which is why our benchmark is called RankEval. This method can effectively measure the

ranking capability of metrics for model performance, i.e., whether metrics can correctly rank models with better performance ahead of those with worse performance.

Given an expected ranking $\mathcal{R}^* = [r_1^*, r_2^*, \dots, r_n^*]$ and an actual ranking $\mathcal{R} = [r_1, r_2, \dots, r_n]$ generated by some metric, we use the following three indicators to evaluate ranking quality:

- 1) **Spearman's Rank Correlation** [38]: Measures the linear correlation between two rankings, defined as:

$$\text{Sp} = 1 - \frac{6 \sum_{i=1}^n (r_i^* - r_i)^2}{n(n^2 - 1)} \quad (42)$$

where r_i^* and r_i respectively represent the ranking of the i -th model in the expected ranking and actual ranking. $\rho_s \in [-1, 1]$, with values closer to 1 indicating more consistent rankings.

- 2) **Kendall's Tau** [39]: Measures the agreement between two rankings, defined as:

$$\text{Kd} = \frac{C - D}{C + D} \quad (43)$$

where C and D respectively represent the number of concordant pairs and discordant pairs. For any two model pairs (i, j) , if i appears before j in the expected ranking and also appears before j in the actual ranking, it is considered a concordant pair; otherwise, it is considered a discordant pair. $\tau \in [-1, 1]$, with values closer to 1 indicating more consistent rankings.

- 3) **Mean Rank Deviation** [40]: Measures the average deviation between the actual ranking and expected ranking, defined as:

$$\text{MD} = \frac{1}{n} \sum_{i=1}^n |r_i^* - r_i| \quad (44)$$

where $|r_i^* - r_i|$ represents the absolute difference in rankings of the i -th model between the two rankings. Smaller MD values indicate more accurate rankings.

This evaluation method is more interpretable, with ranking results that are intuitive and easy to understand, facilitating the analysis of metric effectiveness. Moreover, it avoids the influence of differences in numerical ranges and distributions among different metrics.

B. Computational Latency Analysis (RQ1)

1) **Latency Performance**: To analyze the computational efficiency of different metrics, we evaluated the latency distribution across all metrics on the AccQ task, as illustrated in Figure 2. Among all evaluated metrics, PATE exhibited the highest latency with an average of 36,526.02 ms, followed by VUS-ROC at 12,691.87 ms. In contrast, CCE demonstrated the lowest latency among interval-based metrics, achieving an average of only 37.85 ms. Remarkably, CCE's computational overhead is comparable to point-level metrics such as F1 and AUC-ROC. As a result, CCE provides significant benefits in computing as an interval metric, improving efficiency by up to **965 times**.

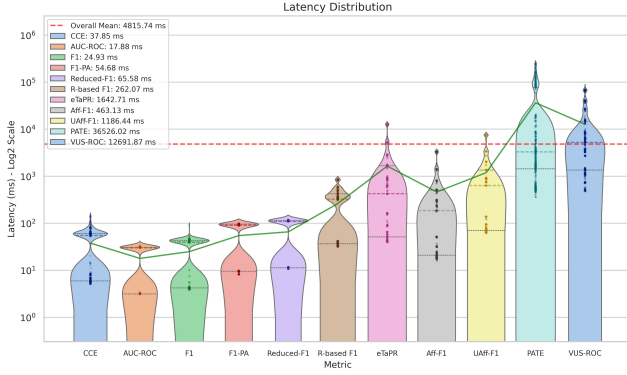


Fig. 2. Latency distribution of different metrics. The violin plot shows the median and interquartile range. The green line represents the average latency of each metric.

2) *CCE Computational Efficiency Analysis*: Figures 3a, 3b, and 3c respectively demonstrate the latency distribution of CCE under different scenarios. From Figure 3a, it is clear that CCE’s computational efficiency is independent of task type, with consistent latency distributions. Figure 3b shows the impact of different time series lengths on CCE latency. We find that CCE’s latency increases approximately linearly, which is consistent with our previous theoretical analysis. Figure 3c demonstrates the impact of the number of anomalies on CCE latency. Overall, the number of event segments has little impact on CCE, which can be approximated as a constant, depending on the computational capabilities of the computer.

C. Metric Effectiveness and Robustness (RQ2)

To verify the effectiveness and robustness of CCE compared to other metrics, we analyzed the performance of different metrics under various parameter settings: $p \in [0.1, 0.2, \dots, 0.9, 1.0]$, $q \in [0.01, 0.05, 0.1, 0.3]$, and $\sigma \in [0.0, 0.05, 0.1]$.

1) *Effectiveness Verification*: For synthetic tasks, anomaly scores are generated by explicitly given parameters (q, p) in ASGM. Therefore, we can directly obtain the expected performance ranking of each model under ideal conditions based on (q, p) and form the expected ranking \mathcal{R}^* . Specifically: AccQ(-R) and LowDisAccQ(-R) are solely determined by accuracy q , with \mathcal{R}^* arranged in ascending order by q ; PreQ-NegP(-R) gives rise to two single-factor tasks: PreQ-NegP-Q only examines anomaly precision q , with \mathcal{R}^* arranged in ascending order by q ; PreQ-NegP-P only examines false positive rate p , with \mathcal{R}^* arranged in descending order by p (lower false positive rate is better), ignoring the other dimension. Subsequently, we compare the rankings generated by each metric for the same set of models with \mathcal{R}^* , calculating Spearman, Kendall, and MD scores to measure metric capability.

The results are shown in Table I. Overall, R-based F1 has the worst ranking consistency capability, with an Sp score of 0.666. The ranking capabilities of F1, F1-PA, Reduced-F1, and eTaPR show small differences, but compared to more advanced metrics like Aff-F1, UAFF-F1, and VUS-ROC, they still lag behind by more than 10%, making these metrics somewhat outdated. Surprisingly, AUC-ROC is not only suitable for point

anomalies but also performs excellently in interval anomaly scenarios. Ultimately, from a holistic perspective, CCE is the best metric across all scenarios, maintaining consistent ranking consistency. VUS-ROC and AUC-ROC are both recommended metrics, while (U)Aff-F1 is only recommended for use in specific scenarios (we will analyze the reasons in the next section).

It should be noted that AUC-ROC and VUS-ROC are not perfect metrics. Further analysis in Appendix A-D (Figs. S11c and S11d) reveals that in noisy scenarios, false positives affect the ranking capability of these metrics, as they pay more attention to the model’s anomaly prediction capability.

2) *Robustness Analysis*: Figs. 4a and 4b respectively demonstrate the performance of different metrics on the AccQ and PreQ-NegP tasks. On the AccQ task, only CCE, AUC-ROC, and VUS-ROC are not affected by noise, while other metrics experience a decline in their ranking capability as noise increases. On the PreQ-NegP task, only CCE, F1, F1-PA, and Reduced-F1 are not affected by noise, while other metrics suffer from incorrect assessment of false positives due to increased noise. Overall, only CCE achieves complete noise robustness, meaning that increasing noise does not affect the metric’s ranking capability.

D. Visualization Analysis (RQ3)

1) *Synthetic Data Visualization*: Figure 5 shows the visualization results of five PreQ-NegP-R models on synthetic data. Their parameters are shown in the first column of Table II, where the numbers after Q and P represent q and p respectively, and the noise size $\sigma = 0.1$. ER represents the theoretical ranking of the model. In this case, only CCE and VUS-ROC can obtain the correct ranking of model performance. Among them, Reduced-F1 has the worst estimation ability, and the score of PreQ0.9-NegP0.1-R, which is expected to rank third, is actually the highest. eTaPR and Aff-F1 also have estimation inaccuracies.

2) *Real Data Visualization*: To better observe and understand anomalies in the dataset, we selected ECG and Power demand from UCR to study the performance of different models in real-world scenarios. The visualization results of model predictions are shown in Figs. 6a and 6b respectively. Table III shows the scores of these models on different metrics.

In ECG, the models with better performance are LSTMAD, USAD, and A.T. In this case, AUC-ROC performs the worst and cannot accurately evaluate USAD’s capability, because AUC-ROC cannot accurately reflect situations where USAD’s anomaly detection capability is not prominent enough, as point-based evaluation is too strict. Additionally, Aff-F1’s score for LSTMAD is 65.0, which is the lowest among all models, which is not reasonable.

In Power, the models with better performance are LOF, LSTMAD, and IForest, while the anomaly scores of USAD and A.T. in the anomaly region are far below the normal interval. However, at this point, VUS-ROC gives USAD a score of 70.4, considering it a better model, but from the visualization, this model is clearly inferior to LOF, yet their VUS-ROC scores are close.

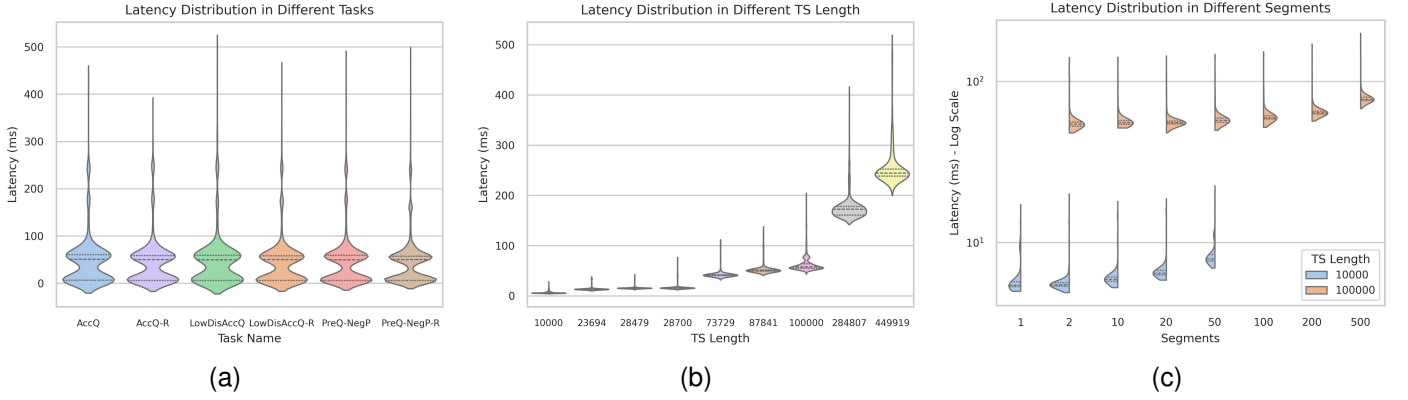


Fig. 3. CCE latency analysis under different scenarios: (a) task type impact, (b) time series length impact, and (c) anomaly segment count impact.

TABLE I

COMPARISON OF RANKING CONSISTENCY ACROSS DIFFERENT INDICATORS (THE LARGER THE SP AND KD, THE BETTER; THE SMALLER THE MD, THE BETTER). **BOLD** INDICATES THE BEST, UNDERLINED INDICATES THE SECOND BEST.

Task	Score	CCE	AUC-ROC	F1	F1-PA	Reduced-F1	R-based F1	eTaPR	Aff-F1	UAff-F1	VUS-ROC
AccQ	Sp	1.000	1.000	0.340	0.340	0.340	0.294	0.780	0.832	0.825	1.000
AccQ	Kd	1.000	1.000	0.248	0.248	0.248	0.193	0.697	0.789	0.779	1.000
AccQ	MD	0.000	0.000	2.094	2.094	2.095	2.344	1.147	0.587	0.622	0.000
LowDisAccQ	Sp	1.000	1.000	0.998	0.998	0.996	0.901	0.846	0.953	0.944	1.000
LowDisAccQ	Kd	1.000	1.000	0.993	0.993	0.990	0.895	0.760	0.927	0.915	1.000
LowDisAccQ	MD	0.000	0.000	0.025	0.025	0.038	0.465	0.933	0.265	0.299	0.000
PreQ-NegP-Q	Sp	1.000	1.000	0.928	0.925	0.895	0.681	0.878	0.883	0.872	1.000
PreQ-NegP-Q	Kd	1.000	<u>0.999</u>	0.931	0.925	0.892	0.674	0.866	0.855	0.839	1.000
PreQ-NegP-Q	MD	0.000	<u>0.003</u>	0.267	0.287	0.371	1.131	0.494	0.551	0.619	<u>0.002</u>
PreQ-NegP-P	Sp	1.000	0.987	1.000	1.000	1.000	0.789	0.876	0.920	0.891	0.990
PreQ-NegP-P	Kd	1.000	0.982	1.000	1.000	1.000	0.776	0.864	0.895	0.863	0.986
PreQ-NegP-P	MD	0.000	0.026	0.000	0.000	0.000	0.267	0.160	0.150	0.192	0.019
Avg.	Sp	1.000	0.997	0.816	0.816	0.808	0.666	0.845	0.897	0.883	<u>0.998</u>
Avg.	Kd	1.000	0.995	0.793	0.792	0.783	0.635	0.797	0.866	0.849	<u>0.996</u>
Avg.	MD	0.000	0.007	0.597	0.601	0.626	1.052	0.684	0.388	0.433	<u>0.005</u>

TABLE II

PERFORMANCE OF DIFFERENT ASGM MODELS ON SYNTHETIC DATASET.

ER	ASGM/Metric	CCE	Reduced-F1	eTaPR	Aff-F1	VUS-ROC
1	PreQ1.0-NegP0.05-R	0.768	0.750	0.712	0.940	0.983
2	PreQ0.9-NegP0.05-R	0.707	0.706	0.706	0.957	0.974
3	PreQ0.9-NegP0.1-R	0.658	0.800	0.732	0.940	0.969
4	PreQ0.9-NegP0.3-R	0.510	0.353	0.459	0.788	0.905
5	PreQ0.8-NegP0.3-R	0.412	0.480	0.516	0.811	0.888

Furthermore, in both ECG and Power cases, we found that eTaPR always considers the model score to be 0, which weakens the discriminative power of this metric. VUS-ROC gives scores higher than 65 to both Random models, but theoretically, VUS-ROC's scoring for random models should be close to 50. At this point, when the model is close to random, the metric error may exceed 30%.

E. Impact of Hyperparameter on CCE (RQ4)

The CCE metric has only one important hyperparameter: the confidence threshold τ . Unlike the buffer size hyperparameters of Aff-F1 and VUS-ROC, this is a dataset-independent hyperparameter. We set the parameter range $\tau \in$

TABLE III

COMPARISON OF MODELS ON ECG AND POWER DATASETS.

ECG	CCE	AUC-ROC	F1	eTaPR	Aff-F1	VUS-ROC
LOF	-0.01	53.72	2.32	0.0	67.2	54.04
IForest	22.3	80.67	4.74	3.47	69.77	82.74
LSTMAD	4.66	83.53	5.33	3.57	65.0	88.19
USAD	2.85	47.26	1.48	1.41	67.03	54.23
A.T.	-0.21	49.27	0.0	0.0	68.46	50.26
Random	2.82	52.85	1.04	0.0	65.42	68.49
Power	CCE	AUC-ROC	F1	eTaPR	Aff-F1	VUS-ROC
LOF	2.63	68.56	1.83	0.0	62.65	75.45
IForest	-5.1	32.21	0.0	0.0	75.18	50.78
LSTMAD	0.5	64.26	2.34	2.45	67.11	79.5
USAD	-1.75	45.31	0.78	0.0	81.8	70.4
A.T.	0.05	51.6	1.97	0.0	71.06	54.21
Random	-4.9	45.08	0.52	0.0	67.03	65.46

[0.1, 0.3, 0.5, 0.7, 0.9] and conducted experiments on the AccQ task. Table IV demonstrates the impact of this parameter on CCE's evaluation capability. Clearly, under any τ setting, CCE maintains consistent ranking capability, indicating that hyperparameter settings do not affect the absolute ranking of models.

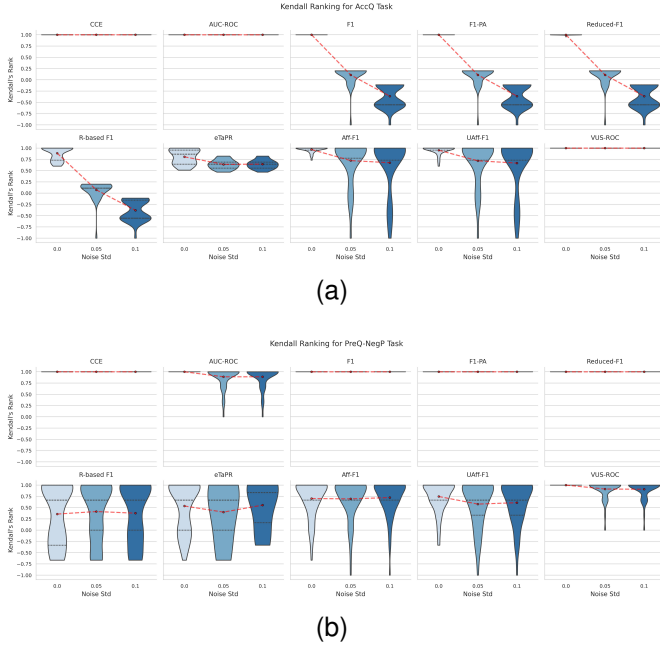


Fig. 4. (a) Performance of different metrics on the AccQ task, (b) Performance of different metrics on the PreQ-NegP task, only considering ranking for p .

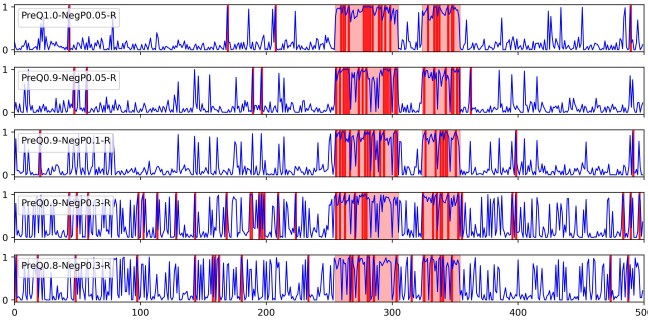


Fig. 5. Visualization of different ASGM models on synthetic dataset.

VI. CONCLUSION

This paper addresses the limitations of previous TSAD evaluation metrics by proposing the CCE metric that characterizes confidence and uncertainty consistency, solving issues such as hyperparameter dependency, low robustness, high computational cost, and lack of prediction consistency evaluation in previous metrics. Additionally, the RankEval benchmark for metric evaluation was constructed, which can assess metric performance across different tasks and robust scenarios. Through theoretical analysis, the boundedness, linear complexity, and robustness of the proposed CCE were proven. Furthermore, extensive numerical experiments validated these properties and hyperparameter independence, while demon-

TABLE IV
SP/KD/MD SCORES OF CCE UNDER DIFFERENT τ SETTINGS.

Threshold	0.1	0.3	0.5	0.7	0.9
τ	1/1/0	1/1/0	1/1/0	1/1/0	1/1/0

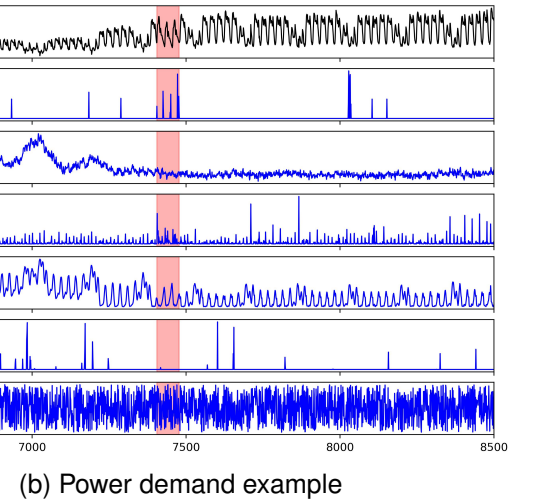
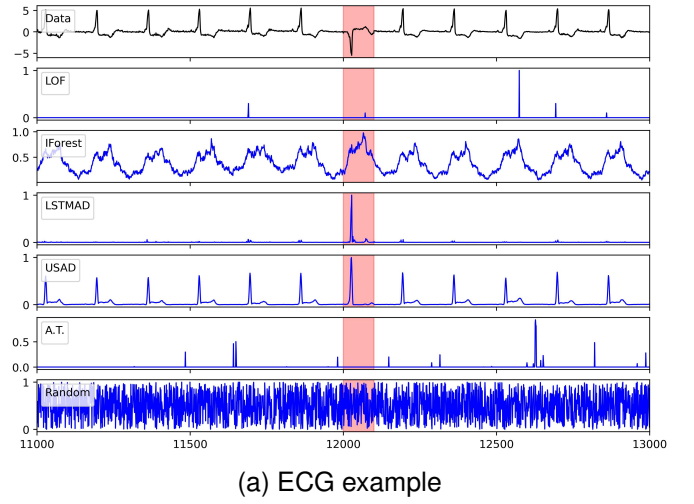


Fig. 6. Visualization of real-world datasets.

strating the robustness limitations of other metrics. Additionally, combined with visualization analysis, the effectiveness and robustness of CCE were further demonstrated. The CCE metric and RankEval benchmark provide a comprehensive TSAD model evaluation solution, advancing TSAD methodologies and facilitating reliable practical model selection.

ACKNOWLEDGMENT

This work was supported in part by National Natural Science Foundation of China No. 92467109, U21A20478, National Key R&D Program of China 2023YFA1011601, and the Major Key Project of PCL, China under Grant PCL2025AS11.

REFERENCES

- [1] T. K. K. Ho, A. Karami, and N. Armanfard, "Graph anomaly detection in time series: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 47, no. 8, pp. 6990–7009, 2025.
- [2] M. Jin, H. Y. Koh, Q. Wen, D. Zambon, C. Alippi, G. I. Webb, I. King, and S. Pan, "A survey on graph neural networks for time series: Forecasting, classification, imputation, and anomaly detection," *IEEE Transactions*

- on *Pattern Analysis and Machine Intelligence*, vol. 46, no. 12, pp. 10466–10485, 2024.
- [3] K. Zhang, Q. Wen, C. Zhang, R. Cai, M. Jin, Y. Liu, J. Y. Zhang, Y. Liang, G. Pang, D. Song, and S. Pan, “Self-supervised learning for time series analysis: Taxonomy, progress, and prospects,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 10, pp. 6775–6794, 2024.
 - [4] R. Laxhammar and G. Falkman, “Online learning and sequential anomaly detection in trajectories,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 6, pp. 1158–1173, 2014.
 - [5] Z. Zhong, Z. Yu, J. Chen, and K. Yang, “Insightful simplicity: Dissimilarity in time series anomaly detection,” in *Proceedings of the ACM Turing Award Celebration Conference - China 2024*, ser. ACM-TURC ’24. New York, NY, USA: Association for Computing Machinery, 2024, p. 242–243.
 - [6] P. Li, Z. Zhong, T. Zhang, Z. Yu, C. Chen, and K. Yang, “A new perspective on time series anomaly detection: Faster patch-based broad learning system,” *arXiv preprint arXiv:2412.05498*, 2024.
 - [7] Z. Zhong, Z. Yu, Y. Yang, W. Wang, K. Yang, and C. L. P. Chen, “Patchad: A lightweight patch-based mlp-mixer for time series anomaly detection,” *IEEE Transactions on Big Data*, pp. 1–15, 2025.
 - [8] C. Zhang, T. Zhou, Q. Wen, and L. Sun, “TFAD: A Decomposition Time Series Anomaly Detection Architecture with Time-Frequency Analysis,” in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, Oct. 2022, pp. 2497–2507, arXiv:2210.09693 [cs].
 - [9] H. Xu, Y. Wang, S. Jian, Q. Liao, Y. Wang, and G. Pang, “Calibrated one-class classification for unsupervised time series anomaly detection,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 11, pp. 5723–5736, 2024.
 - [10] Y. Chen, C. Zhang, M. Ma, Y. Liu, R. Ding, B. Li, S. He, S. Rajmohan, Q. Lin, and D. Zhang, “Imdiffusion: Imputed diffusion models for multivariate time series anomaly detection,” *Proc. VLDB Endow.*, vol. 17, no. 3, p. 359–372, Nov. 2023.
 - [11] T. Zhan, Y. He, Y. Deng, Z. Li, W. Du, and Q. Wen, “Time evidence fusion network: Multi-source view in long-term time series forecasting,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–15, 2025.
 - [12] B. Li, W. Cui, L. Zhang, C. Zhu, W. Wang, I. W. Tsang, and J. T. Zhou, “Diffformer: Multi-resolutional differencing transformer with dynamic ranging for time series analysis,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 11, pp. 13586–13598, 2023.
 - [13] D. Kim, S. Park, and J. Choo, “When model meets new normals: Test-time adaptation for unsupervised time-series anomaly detection,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 12, pp. 13113–13121, Mar. 2024.
 - [14] R. Wu and E. J. Keogh, “Current time series anomaly detection benchmarks are flawed and are creating the illusion of progress,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 3, pp. 2421–2429, 2023.
 - [15] Z. Zhong, Z. Yu, X. Xi, Y. Xu, W. Cao, Y. Yang, K. Yang, and J. You, “Simad: A simple dissimilarity-based approach for time-series anomaly detection,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–12, 2025.
 - [16] S. Sørbrø and M. Ruocco, “Navigating the metric maze: a taxonomy of evaluation metrics for anomaly detection in time series,” *Data Min. Knowl. Discov.*, vol. 38, no. 3, p. 1027–1068, Nov. 2023.
 - [17] N. Tatbul, T. J. Lee, S. Zdonik, M. Alam, and J. Gottschlich, “Precision and recall for time series,” in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3–8, 2018, Montréal, Canada*, S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, pp. 1924–1934.
 - [18] H. Xu, W. Chen, N. Zhao, Z. Li, J. Bu, Z. Li, Y. Liu, Y. Zhao, D. Pei, Y. Feng *et al.*, “Unsupervised anomaly detection via variational auto-encoder for seasonal kpis in web applications,” in *Proceedings of the 2018 world wide web conference*, 2018, pp. 187–196.
 - [19] H. Si, J. Li, C. Pei, H. Cui, J. Yang, Y. Sun, S. Zhang, J. Li, H. Zhang, J. Han *et al.*, “Timeseriesbench: An industrial-grade benchmark for time series anomaly detection models,” *2024 IEEE 35th International Symposium on Software Reliability Engineering (ISSRE)*, pp. 61–72, 2024.
 - [20] A. Huet, J. M. Navarro, and D. Rossi, “Local evaluation of time series anomaly detection algorithms,” in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 635–645.
 - [21] J. Paparrizos, P. Boniol, T. Palpanas, R. S. Tsay, A. Elmore, and M. J. Franklin, “Volume under the surface: a new accuracy evaluation measure for time-series anomaly detection,” *Proceedings of the VLDB Endowment*, vol. 15, no. 11, pp. 2774–2787, 2022.
 - [22] R. Ghorbani, M. J. Reinders, and D. M. Tax, “Pate: Proximity-aware time series anomaly evaluation,” in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, ser. KDD ’24. New York, NY, USA: Association for Computing Machinery, 2024, p. 872–883.
 - [23] Z. Yu, Z. Zhong, K. Yang, W. Cao, and C. L. P. Chen, “Broad learning autoencoder with graph structure for data clustering,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 1, pp. 49–61, 2024.
 - [24] S. Sun, Z. Zhong, N. Yu, X. Gong, and K. Yang, “Humanmod: A multi-rag collaborative llm for inclusive urban public healthcare services,” *Applied Soft Computing*, vol. 184, p. 113684, 2025.
 - [25] Z. Zhong, K. Yang, Z. Yu, Y. Shi, and C. L. Philip Chen,

- “Towards efficient anomaly detection using memory broad learning system,” in *2023 9th International Conference on Control Science and Systems Engineering (ICCSSE)*, 2023, pp. 252–257.
- [26] Q. Liu and J. Paparrizos, “The elephant in the room: Towards a reliable time-series anomaly detection benchmark,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, Eds., vol. 37. Curran Associates, Inc., 2024, pp. 108 231–108 261.
- [27] *TSB-AutoAD: Towards Automated Solutions for Time-Series Anomaly Detection*, ser. PVLDB, vol. 18. VLDB Foundation, 2025.
- [28] Z. Ma and A. Leijon, “Bayesian estimation of beta mixture models with variational inference,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 11, pp. 2160–2173, 2011.
- [29] M. E. Silva, I. Pereira, and B. McCabe, “Bayesian outlier detection in non-gaussian autoregressive time series,” *Journal of Time Series Analysis*, vol. 40, no. 5, pp. 631–648, 2019.
- [30] A.-H. Shin, S. T. Kim, and G.-M. Park, “Time Series Anomaly Detection Using Transformer-Based GAN With Two-Step Masking,” *IEEE Access*, vol. 11, pp. 74 035–74 047, 2023.
- [31] Y. Liu, Y. Tian, Y. Mi, H. Liu, J. Wang, and W. Pedrycz, “Landmark block-embedded aggregation autoencoder for anomaly detection,” *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 55, no. 2, pp. 1004–1019, 2025.
- [32] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana, and E. Keogh, “The ucr time series archive,” *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 6, pp. 1293–1305, 2019.
- [33] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, “Lof: identifying density-based local outliers,” in *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, 2000, pp. 93–104.
- [34] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation forest,” in *2008 eighth ieee international conference on data mining*. IEEE, 2008, pp. 413–422.
- [35] J. Audibert, P. Michiardi, F. Guyard, S. Marti, and M. A. Zuluaga, “USAD: UnSupervised Anomaly Detection on Multivariate Time Series,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Virtual Event CA USA: ACM, 2020, pp. 3395–3404.
- [36] J. Xu, H. Wu, J. Wang, and M. Long, “Anomaly transformer: Time series anomaly detection with association discrepancy,” in *International Conference on Learning Representations*, 2021.
- [37] H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, and M. Long, “Timesnet: Temporal 2d-variation modeling for general time series analysis,” in *The eleventh international conference on learning representations*, 2022.
- [38] C. Spearman, “The proof and measurement of association between two things,” *The American journal of psychology*, vol. 100, no. 3/4, pp. 441–471, 1987.
- [39] M. G. Kendall, “A new measure of rank correlation,” *Biometrika*, vol. 30, no. 1-2, pp. 81–93, 1938.
- [40] A. G. Pacheco and R. A. Krohling, “Ranking of classification algorithms in terms of mean–standard deviation using a-topsis,” *Annals of Data Science*, vol. 5, no. 1, pp. 93–110, 2018.
- [41] A. Abdulaal, Z. Liu, and T. Lancewicki, “Practical approach to asynchronous multivariate time series anomaly detection and localization,” in *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, 2021, pp. 2485–2494.
- [42] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom, “Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding,” in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 387–395.
- [43] J. Goh, S. Adepu, K. N. Junejo, and A. Mathur, “A dataset to support research in the design of secure water treatment systems,” in *Critical Information Infrastructures Security*, G. Havarneanu, R. Setola, H. Nassopoulos, and S. Wolthusen, Eds. Cham: Springer International Publishing, 2017, pp. 88–99.



Zhijie Zhong received the B.S. degree in 2022 from the Harbin Engineering University, Harbin, China and he is currently pursuing the Ph.D. degree in the School of Future Technology, South China University of Technology, Guangzhou, China. His research interests include data mining, machine learning, time series analysis, anomaly detection, and large language model (LLM).

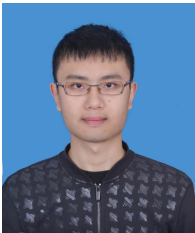


Zhiwen Yu (S'06-M'08-SM'14) is a Professor in School of Computer Science and Engineering, South China University of Technology, China. He received the Ph.D. degree from the City University of Hong Kong, Hong Kong, in 2008. Dr. Yu has authored or coauthored more than 200 refereed journal articles and international conference papers, including more than 70 articles in the journals of IEEE Transactions. His google citation is more than 10000, and h-index is 44. He is an Associate Editor of the IEEE Transactions on systems, man, and cybernetics: systems.

He is a senior member of IEEE and ACM, a Member of the Council of China Computer Federation (CCF).



Yiu-Ming Cheung (Fellow, IEEE) received the PhD degree from the Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, in 2000. He is currently a chair professor with the Department of Computer Science, Hong Kong Baptist University, Hong Kong. His current research interests include machine learning, pattern recognition, and visual computing. He is a fellow of the American Association for the Advancement of Science (AAAS), Institution of Engineering and Technology (IET), British Computer Society (BCS), and Asia-Pacific Artificial Intelligence Association (AAIA). He is the editor-in-chief of IEEE Transactions on Emerging Topics in Computational Intelligence. He also served as an associate editor for IEEE Transactions on Cybernetics, IEEE Transactions on Cognitive and Developmental Systems, IEEE Transactions on Neural Networks and Learning Systems from 2014 to 2020, Pattern Recognition, Knowledge and Information Systems, and Neurocomputing, just to name a few.



Kaixiang Yang (M'21) received the B.S. degree and M.S. degree from the University of Electronic Science and Technology of China and Harbin Institute of Technology, China, in 2012 and 2015, respectively, and the Ph.D. degree from the School of Computer Science and Engineering, South China University of Technology, China, in 2020. He has been a Research Engineer with the 7th Research Institute, China Electronics Technology Group Corporation, Guangzhou, China, from 2015 to 2017, and has been a Postdoctoral Researcher with Zhejiang University from 2020 to 2021. He is now with the School of Computer Science and Engineering, South China University of Technology. His research interests include pattern recognition, machine learning, and industrial data intelligence.

APPENDIX

This is the appendix of *CCE: Confidence-Consistency Evaluation for Time Series Anomaly Detection*.

A. SUPPLEMENTARY EXPERIMENTS

A. Datasets

Table S5 illustrates the characteristics of all datasets employed in RankEval, which encompasses real-world datasets: MSL, SMD, SWAT, Creditcard, and PSM [41, 42, 43].

A total of 30 synthetic datasets were generated, with examples named *100k-20seg-50L* and *100k-20seg-50H*. In such names, *100k* refers to the total length of the time series (corresponding to the TS Length column), *20seg* means that 20 anomaly segments are expected to be generated (i.e., the Segments column), and *50L* indicates that the average length of anomaly segments is 50 with low variance (for instance, lengths ranging from 40 to 60). In contrast, *50H* also represents an average segment length of 50 but with higher variance (such as lengths between 1 and 99). The *Max/Min Seg Length* columns respectively stand for the maximum and minimum lengths of the generated anomaly segments. For time series lengths of 100k and 10k, 15 synthetic datasets were created each. Moreover, 6 real-world datasets were selected according to their characteristics like length, number of anomalies, and anomaly lengths.

TABLE S5
TIME SERIES CHARACTERISTICS OF SYNTHETIC AND REAL DATA.

Dataset Name	TS Length	Segments	Max Seq Length	Min Seq Length
100k-20seg-50L	100000	20	60	40
100k-200seg-50L	100000	200	60	40
100k-20seg-50H	100000	20	99	1
100k-200seg-50H	100000	200	99	1
100k-50seg-20L	100000	50	30	10
100k-500seg-20L	100000	500	30	10
100k-50seg-20H	100000	50	39	1
100k-500seg-20H	100000	500	39	1
100k-10seg-100L	100000	10	110	90
100k-100seg-100L	100000	100	110	110
100k-10seg-100H	100000	10	199	1
100k-100seg-100H	100000	100	199	1
100k-2seg-500L	100000	2	550	450
100k-20seg-500L	100000	20	550	450
100k-2seg-500H	100000	2	999	1
100k-20seg-500H	100000	20	999	1
10k-2seg-50L	10000	2	60	40
10k-20seg-50L	10000	20	60	40
10k-2seg-50H	10000	2	99	1
10k-20seg-50H	10000	20	99	1
10k-5seg-20L	10000	5	30	10
10k-50seg-20L	10000	50	30	10
10k-5seg-20H	10000	5	39	1
10k-50seg-20H	10000	50	39	1
10k-1seg-100L	10000	1	110	90
10k-10seg-100L	10000	10	110	110
10k-1seg-100H	10000	1	199	1
10k-10seg-100H	10000	10	199	1
10k-2seg-500L	10000	2	550	450
10k-2seg-500H	10000	2	999	1
MSL	73729	36	1141	11
Creditcard	284807	465	5	1
SWAT	449919	35	35900	101
SMD-1-1	28479	8	721	2
SMD-2-1	23694	13	452	8
SMD-3-1	28700	4	131	21
PSM	87841	72	8861	1

B. Model Comparison in Real-World Scenario

Table S6 illustrates the average performance of multiple TSAD models across the MSL, SWAT, PSM, SMD, and Creditcard datasets. Current analysis reveals that no single state-of-the-art model achieves superior performance across all evaluation metrics. Additionally, the consistently low CCE scores across these models suggest that current approaches have not sufficiently considered the prediction consistency issue. This observation highlights a promising new direction for advancing time series anomaly detection models in future research.

TABLE S6
COMPARATIVE PERFORMANCE OF DIFFERENT TIME SERIES MODELS IN REAL-WORLD SCENARIO (RESULTS ARE AVERAGED).

Avg.	CCE	AUC-ROC	F1	eTaPR	Aff-F1	VUS-ROC
LOF	0.27	55.25	8.04	7.15	69.69	68.32
IForest	2.1	54.58	9.98	8.83	49.81	67.22
LSTMAD	1.46	79.87	25.5	18.53	61.63	84.71
USAD	1.77	76.84	23.64	16.31	69.51	81.54
A.T.	0.25	50.3	4.66	3.49	62.81	58.79
Donut	0.23	71.07	19.08	14.37	61.26	79.55
TimesNet	0.22	58.52	12.54	12.62	73.68	71.55

C. Robustness Analysis

Figs. S7 and S8 respectively demonstrate the sensitivity of different metrics to noise on the AccQ and LowDisAccQ tasks. Figs. S9 and S10 both show the sensitivity of different metrics to noise on the PreQ-NegP task, separately considering only the effects of p and q respectively. We found that only CCE can maintain robustness to noise across all tasks, while other metrics such as F1, F1-PA, and Reduced-F1 cannot maintain robustness on AccQ and PreQ-NegP tasks. R-based F1, eTaPR, Aff-F1, and UAff-F1 cannot maintain robustness on AccQ, LowDisAccQ, and PreQ-NegP tasks, and even under noise-free conditions, they still cannot achieve perfect ranking, indicating that these metrics may be prone to inflated evaluation. Additionally, AUC-ROC and VUS-ROC show slightly insufficient robustness when considering q , but their performance on other tasks is far better than other comparison metrics.

D. Impact of Noise on Different Metrics

Figs. (S11a-S11d) demonstrate the robustness of different metrics to noise. Under noise-free conditions (i.e., noise std=0), both AUC-ROC and VUS-ROC maintain their ranking capabilities without degradation. However, when noise is present, AUC-ROC exhibits ranking errors at low false positive rates, while VUS-ROC shows ranking errors at high false positive rates. In contrast, Aff-F1 consistently demonstrates ranking inconsistency issues regardless of noise presence, and its error increases progressively with rising false positive rates. This phenomenon occurs because Aff-F1 is based on interval membership, where increased false positives may be interpreted as valid warnings, leading to inflated Aff-F1 scores and consequently introducing experimental bias.

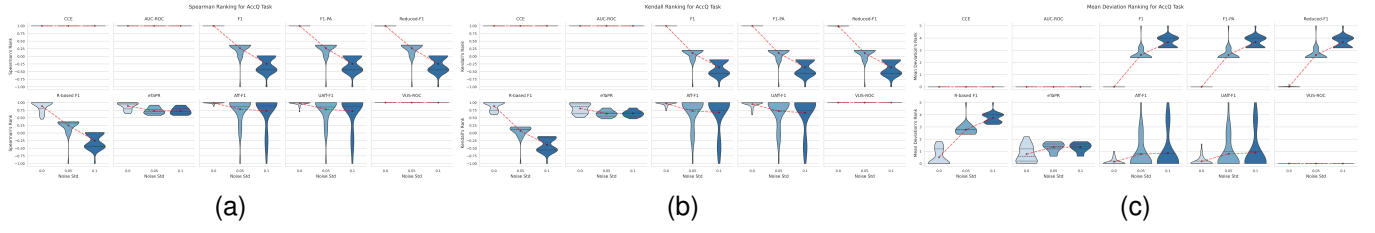


Fig. S7. Evaluation of ranking capability of different metrics on AccQ task. (a) Spearman correlation (b) Kendall correlation (c) Mean Rank Deviation

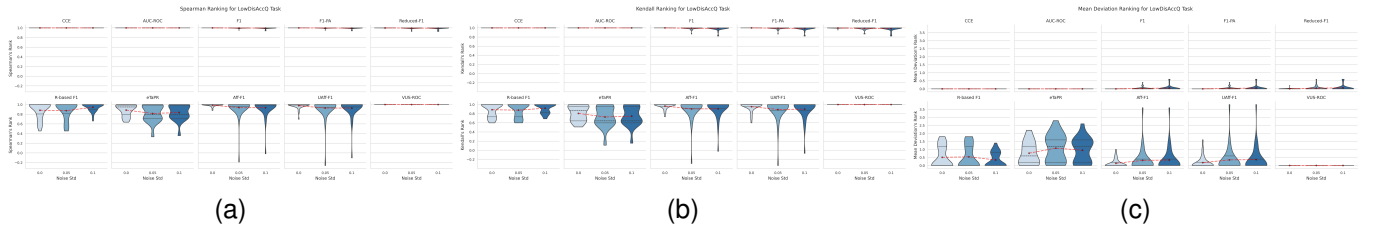


Fig. S8. Evaluation of ranking capability of different metrics on LowDisAccQ task. (a) Spearman correlation (b) Kendall correlation (c) Mean Rank Deviation

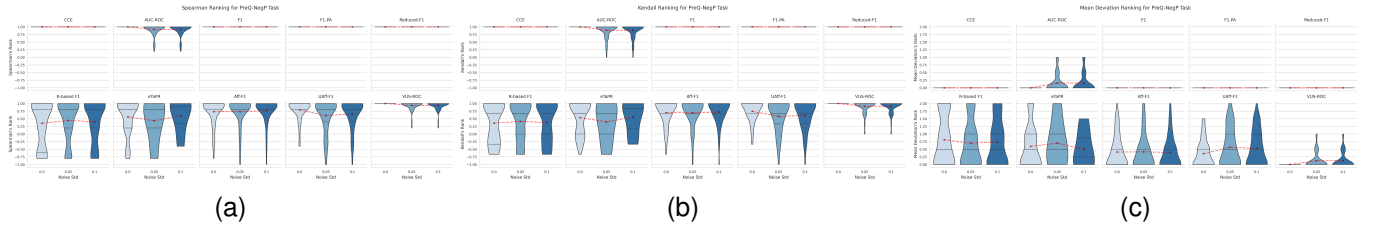


Fig. S9. Evaluation of ranking capability of different metrics on PreQ-NegP task, only considering ranking for p . (a) Spearman correlation (b) Kendall correlation (c) Mean Rank Deviation

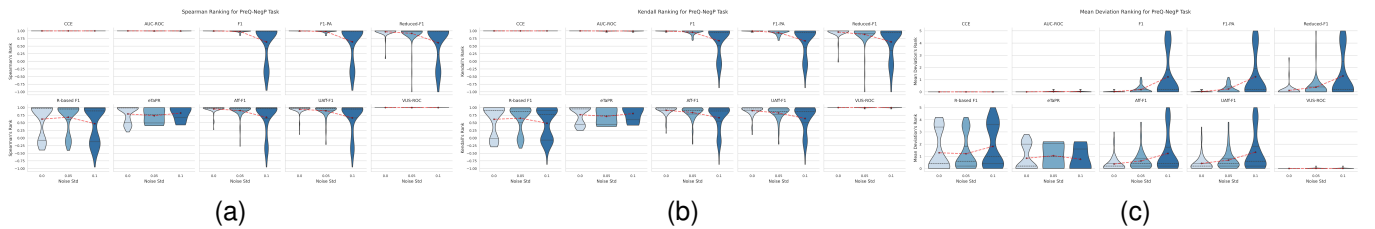


Fig. S10. Evaluation of ranking capability of different metrics on PreQ-NegP task, only considering ranking for q . (a) Spearman correlation (b) Kendall correlation (c) Mean Rank Deviation

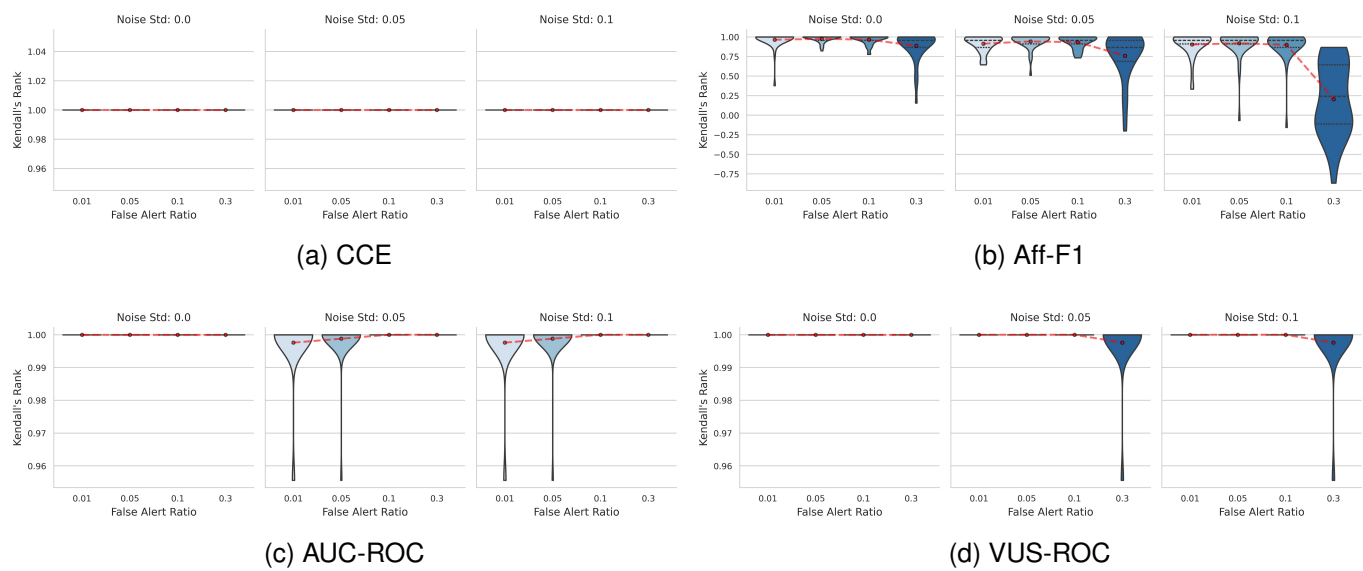


Fig. S11. Robustness of different metrics to noise.