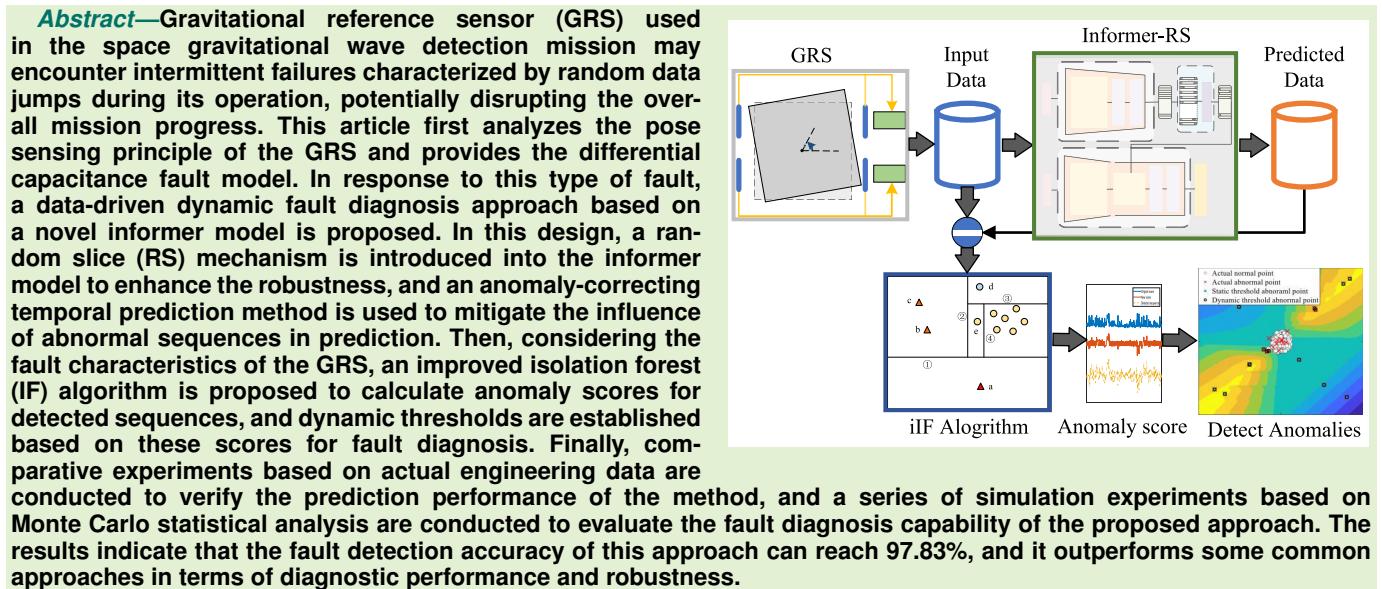


A Dynamic Fault Diagnosis Method for Gravitational Reference Sensor Based on Informer-RS

Cheng Bi[✉], Xiaokui Yue[✉], Zhaohui Dang[✉], Yibo Ding[✉], and Yonghe Zhang[✉]



Index Terms—Dynamic threshold, fault diagnosis, gravitational reference sensor (GRS), informer model.

I. INTRODUCTION

IN THE space gravitational wave detection mission, satellite formations could realize ultralong baseline laser interferometry. Compared with ground-based detectors, this approach is more sensitive to mid-to-low-frequency gravitational waves and has become the focal point of development for various

Received 9 October 2024; revised 26 November 2024; accepted 1 December 2024. Date of publication 11 December 2024; date of current version 14 January 2025. This work was supported in part by the National Key Research and Development Program of China: Gravitational Wave Detection Project under Grant 2021YFC22026 and Grant 2021YFC2202603 and in part by the National Natural Science Foundation of China under Grant 12172288. The associate editor coordinating the review of this article and approving it for publication was Prof. Ling Pei. (*Corresponding author: Zhaohui Dang.*)

Cheng Bi, Xiaokui Yue, Zhaohui Dang, and Yibo Ding are with the National Key Laboratory of Aerospace Flight Dynamics, School of Astronautics, Northwestern Polytechnical University, Xi'an 710072, China (e-mail: verybc@mail.nwpu.edu.cn; xkyue@nwpu.edu.cn; dangzhaohui@nwpu.edu.cn; dingyibo@nwpu.edu.cn).

Yonghe Zhang is with the Institute of Microsatellite Innovation, Chinese Academy of Sciences, Shanghai 201304, China (e-mail: zhangyh@microsate.com).

Digital Object Identifier 10.1109/JSEN.2024.3510739

space-faring nations around the world [1], [2], [3]. As a core component of the mission, gravitational reference sensor (GRS) is capable of sensing the six-degree-of-freedom (DoF) motion of test mass (TM) and provides electrostatic feedback control. However, due to the complexity and variability of the space environment, the GRS is vulnerable to faults [4]. If the faults are left undiagnosed and unaddressed, they may endanger the entire spacecraft system and cause severe economic loss. However, there is currently limited research on the GRS faults, and the fault mechanisms are not yet clear. Therefore, it is imperative to initiate the analysis of the GRS faults by delving into its sensing principles, exploring the propagation mechanisms and manifestations of potential malfunctions. In addition, it is essential to design effective fault diagnosis methods based on the unique fault mechanisms, which will significantly enhance the reliability and safety of the system.

GRS typically consists of two main components: the sensitive structure (SS) and the front-end electronics unit (FEE) [5]. Within the SS, the electrode distribution of the capacitive sensors is located on the electrode housing (EH), while the TM serves as another electrode situated at the center of the

EH. The FEE performs capacitance measurement through the electrodes and calculates the current pose of the TM. Subsequently, processes such as the proportional integral derivative control, the servo amplification feedback, and others are used to achieve acceleration measurement and six-DoF control [6]. Therefore, the GRS faults may occur on the SS and the FEE. The faults on the SS belong to structural faults, which is not considered in this article and is excluded from our study. We specifically focus on investigating faults during the sensing process in the FEE.

Fault diagnosis methods in aerospace missions are typically classified into three categories: model-based methods, signal-based methods, and knowledge-based methods [7]. Model-based methods usually require precise physical models [8], [9], [10], making them challenging to apply to drag-free systems in space gravitational wave detection missions, which has high nonlinearity and strong coupling characteristics. Similarly, knowledge-based methods require a comprehensive understanding of the system and are also not well-suited for the detection missions in the preliminary stage. In recent years, artificial-intelligence-based fault diagnosis methods have garnered significant attention [11], [12], [13], [14]. The methods, along with signal-based methods, fall under the category of data-driven methods. Data-driven methods do not rely on system models but instead use measurable data in the systems as the references for fault diagnosis. By extracting and analyzing data features, they identify abnormal components in the data and achieve fault diagnosis. Consequently, this type of methods is particularly suitable for the fault diagnosis problem of the GRS.

The data-driven methods include dimensionality reduction approaches [15], [16], clustering-based approaches [17], [18], nearest neighbor approaches [19], [20], and other approaches, each with its own advantages and limitations, such as computational complexity, interpretability, and portability [21]. Considering the large volume of telemetry data generated in the current space missions [22], [23], this is suitable for the use of deep neural network (DNN) methods, which require a large amount of data to train the models. DNN models, due to their powerful feature learning capabilities, have been widely applied in various fault diagnosis applications across different domains [24], [25], [26], [27]. Among them, a temporal prediction method based on long short-term memory (LSTM) has demonstrated excellent performance in the aerospace missions [23], [28], [29], [30]. LSTM, a type of recurrent neural network (RNN), is suitable for time-series prediction tasks owing to its ability to capture long-term dependencies by incorporating weighted self-recurrence with the context information. In aerospace missions, the LSTM is often applied to predict telemetry data by leveraging its ability to extract temporal features. Fault diagnosis is then achieved by analyzing the residuals between the measured telemetry data and the prediction data [31], [32]. However, LSTM-based methods are not robust, and they may experience degraded performance when data quality is compromised. Furthermore, some current researches on using LSTM for fault diagnosis are mostly achieved through short sequence times-series forecasting [33], which requires frequent processing of prediction

data that may lead to situations where the processing speed cannot keep up with the prediction speed. However, the LSTM has the problem of gradient explosion and cannot perform parallel computations, leading to low prediction accuracy and significantly increased computation time when dealing with long sequence time-series forecasting (LSTF) problems, which in turn affects the fault diagnosis effectiveness [34].

Transformer is a DNN based on self-attention mechanism [35]. Unlike RNNs, the transformer enables parallel computation, has more efficient computational capabilities, and exhibits stronger learning and expressive abilities. Currently, it has been widely applied in different fields such as natural language processing and computer vision. Some researchers believe that the transformer holds tremendous potential in prediction [36], [37]. However, due to its high time complexity and memory utilization, the transformer cannot be directly applied to time-series prediction [38]. Numerous studies have been conducted to address the issue. The Longformer is proposed by Beltagy et al. [39], which introduces a local attention mechanism with linear complexity in relation to sequence length. Wang et al. [40] introduced fixed linear projection matrices into self-attention, which comprehensively reduces the computational complexity of the entire model. However, this method is only applicable in specific scenarios. Wu et al. [41] incorporated time-series decomposition into the transformer model and used an autocorrelation mechanism instead of self-attention, achieving efficient LSTF. Zhou et al. [34] developed the informer model, which effectively reduces the computational complexity and memory utilization. This model shows excellent performance on four large-scale datasets and has been applied in many engineer applications.

As a new temporal prediction model, the informer has not yet been extensively studied for fault diagnosis in aerospace missions. To provide an efficient fault diagnosis method specifically for GRS in the space gravitational wave detection mission and explore the performance of the informer in this context, this article proposes a dynamic fault diagnosis method based on Informer with random slice (Informer-RS). Comparative experiments are conducted to evaluate the performance of the proposed Informer-RS against other models. The main contributions and innovations of this study are listed below.

- 1) A comprehensive pose sensing principle of the GRS in gravitational wave detection missions is elaborated. Based on this principle, a specific type of fault commonly occurring in the GRS is analyzed, and a novel fault diagnosis strategy tailored to this fault type is proposed.
- 2) The informer model is used to diagnose GRS faults, and its robustness against noise is enhanced by integrating a random slice (RS) mechanism. Considering the impact of abnormal sequences on prediction accuracy, an anomaly-correcting temporal prediction method is proposed, effectively mitigating the influence of such anomalies.
- 3) The isolation forest (IF) algorithm is improved to calculate anomaly scores specifically adapted to the characteristics of the GRS. The enhanced algorithm

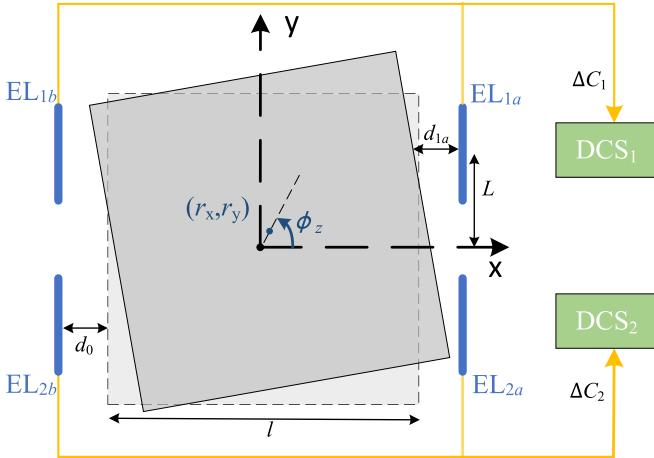


Fig. 1. Pose sensing principle of the TM in GRS.

demonstrates improved stability and the capability to identify anomaly directions. In addition, a dynamic thresholding method is refined to achieve more adaptive and accurate anomaly detection.

The remaining part of this article is organized as follows. Section II elaborates the sensing principle and a fault diagnosis strategy of the GRS. An anomaly-correcting temporal prediction method based on Informer-RS is proposed in Section III. In Section IV, a new anomaly detection method is introduced. Simulation experiments are conducted to verify the effectiveness of the proposed approach in Section V. Finally, the whole article is summarized in Section VI.

II. PROBLEM FORMULATION

GRS is likely to malfunction during space missions, but there is currently little research on the failure mechanism and manifestations of this sensor. This section describes the comprehensive pose sensing principle of the GRS. Based on this principle, a possible fault manifestation of the GRS is proposed, and a corresponding fault diagnosis strategy is developed for this type of fault.

A. Sensing Principle and Fault Model of GRS

Taking the GRS in the laser interferometer space antenna (LISA) mission as an example, its pose sensing principle is based on the variable spacing differential capacitance sensor (DCS) in the FEE [42]. In the EH, there are two DCSs on each axis, and each DCS consists of two electrodes for measuring the differential capacitance, as shown in Fig. 1. In this figure, x -axis is the measured axis, DCS₁ and DCS₂ are the DCSs along the x -axis, EL_{1a} and EL_{1b} are the electrodes on DCS₁, while EL_{2a} and EL_{2b} are the electrodes on DCS₂.

Assuming the initial spacing between the TM and the electrodes is d_0 , the electrode area is S_0 , and the initial capacitance is C_0 , when the TM moves by a distance Δd along the x -axis, according to the formula for the calculation of parallel plate capacitance, the differential capacitance on both sides is given by

$$\Delta C = -2\epsilon_0\epsilon_r S_0 \frac{\Delta d}{d_0^2 - \Delta d^2} \quad (1)$$

where ϵ_0 is the vacuum permittivity and ϵ_r is the relative permittivity. Considering that the TM in the space gravitational wave detection mission only experiences small displacement, i.e., $\Delta d \ll d$, (1) can be simplified as

$$\Delta C \approx -2C_0 \frac{\Delta d}{d_0}. \quad (2)$$

If the differential capacitance value ΔC is measured by the DCS, the displacement Δd can be calculated.

When both the position and orientation of the TM change simultaneously, the parallel relationship between the TM and the electrodes will be no longer maintained. In Fig. 1, the TM undergoes both translational and rotational motions, with the center coordinates being (r_x, r_y) after translation and a counterclockwise rotation angle ϕ_z . The side length of the TM is l , and the distance between the center of the electrode and the y -axis is L . At this point, the distances between the TM and electrodes can be calculated, among which the distance between the TM and EL_{1a} is given by

$$d_{1a} = d_0 - r_x + L \tan \phi_z + 0.5l(1 - \sec \phi_z) - r_y \tan \phi_z. \quad (3)$$

From (3), it can be observed that the displacement and rotation of the TM are coupled, which makes it challenging to independently solve for them. Given that these states are all small, decoupling and elimination of the displacement r_y along the nonmeasured axis can be performed accordingly. Based on the formula for nonparallel plate capacitance, the decoupled differential capacitance values between the electrodes EL_{1a} and EL_{1b}, induced by the translation and rotation of the TM, can be calculated as $2r_x C_0/d_0$ and $-2\phi_z L C_0/d_0$, respectively. Similarly, the differential capacitance values between EL_{2a} and EL_{2b} can be obtained. Therefore, the measurement values for DCS₁ and DCS₂ can be expressed as

$$\begin{cases} \Delta C_1 = 2C_0 \frac{r_x - \phi_z L}{d_0} \\ \Delta C_2 = -2C_0 \frac{r_x + \phi_z L}{d_0} \end{cases} \quad (4)$$

where ΔC_1 represents the measurement value of DCS₁, and ΔC_2 represents the measurement value of the DCS₂. Based on ΔC_1 and ΔC_2 , the pose measurement values of the TM along the x -axis can be computed

$$\begin{cases} r_x = \frac{d_0}{4C_0} (\Delta C_1 - \Delta C_2) \\ \phi_z = -\frac{d_0}{4C_0 L} (\Delta C_1 + \Delta C_2). \end{cases} \quad (5)$$

According to the sensing principle of the GRS, one possible type of fault is that the failure may occur in the two DCSs on each axis. When this type of fault occurs, the sensing displacement and rotation angle of the TM corresponding to the faulty DCS will simultaneously exhibit abnormalities. Furthermore, (5) represents that different DCS faults have different effects on the pose measurement values. When DCS₁ fails, causing an anomaly in the values of ΔC_1 , the measurement values of r_x and ϕ_z will exhibit opposite directional anomalies since their coefficients are different. On the other hand, if DCS₂ fails, the values will exhibit anomalies in the same direction.

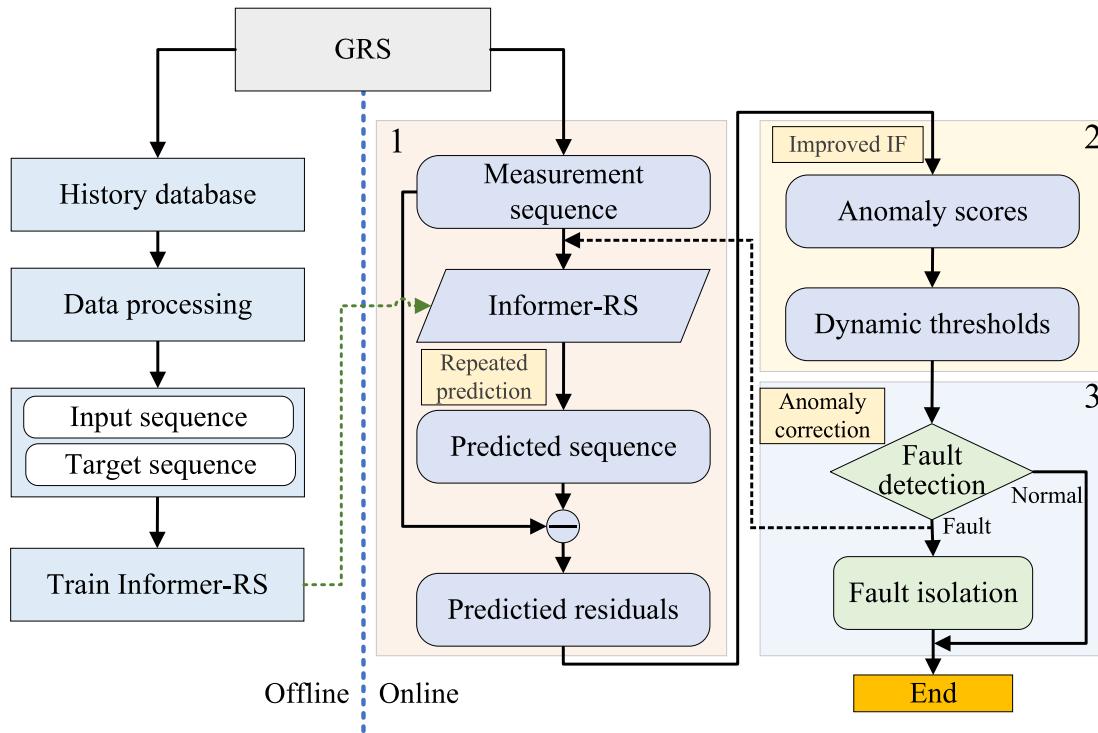


Fig. 2. Flowchart of the proposed fault diagnosis approach.

B. Fault Diagnosis Strategy for GRS

Considering GRS faults caused by DCS issues, there are three possible scenarios.

- 1) Faults in DCSs across different axes simultaneously. Since the measurements from DCSs on different axes are decoupled during computation, they do not interfere with each other. Therefore, these faults can be treated as independent single-axis DCS fault problems to be resolved individually.
- 2) Faults in two DCSs on the same axis simultaneously. In this scenario, either the displacement or angular readouts of the TM are likely to exhibit anomalies. However, these anomalies can be diverse and challenging to predict. As the proposed method aims to diagnose GRS faults based solely on the pose measurements of the TM, it can only detect faults but cannot determine whether the faults are simultaneous. Further determination would require examining the raw DCS measurement data, which lies outside the scope of this study.
- 3) Fault in a single DCS on the same axis. This scenario can be addressed by detecting and isolating faults based on the pose measurements of the TM.

In summary, the goal of fault diagnosis for the GRS, as described in this article, is to detect the fault and identify the faulty DCS when one of the two DCSs on the same axis fails. First, fault detection can be achieved by analyzing the prediction values of the GRS sensing data. Specifically, the situations are reflected in the prediction residuals of r_x and ϕ_z , denoted as $e_{r_x}^t = [e_{r_x}^{(t+1)}, e_{r_x}^{(t+2)}, \dots, e_{r_x}^{(t+m)}]^T$ and $e_{\phi_z}^t = [e_{\phi_z}^{(t+1)}, e_{\phi_z}^{(t+2)}, \dots, e_{\phi_z}^{(t+m)}]^T$, respectively. In normal situations, the residuals of the data are typically distributed

TABLE I
RESIDUAL PERFORMANCE AND REASONS

Case	Residual manifestation	Reason
1	$e_{r_x}^t$ and $e_{\phi_z}^t$ both are normal	GRS is normal
2	$e_{r_x}^t$ and $e_{\phi_z}^t$ deviate in the same direction	DCS ₂ malfunctions
3	$e_{r_x}^t$ and $e_{\phi_z}^t$ deviate in the opposite direction	DCS ₁ malfunctions
4	One of $e_{r_x}^t$ and $e_{\phi_z}^t$ deviates	GRS malfunctions, but DCSs are normal

around zero. However, when anomalies occur, the residuals at abnormal points will deviate from zero. Based on the characteristic of the DCS faults, fault isolation can be realized by examining the different behaviors of the measurement data in the two channels, and the behaviors are reflected in the directions of the deviations in $e_{r_x}^t$ and $e_{\phi_z}^t$. When the DCS malfunctions, both the residuals will simultaneously deviate significantly, and the deviation in the same or opposite directions, respectively, represents different DCS fails. The different manifestations and reasons of residuals $e_{r_x}^t$ and $e_{\phi_z}^t$ are summarized in Table I, and this article only focuses on cases 1–3.

To achieve the fault diagnosis goal, we adopt a fault diagnosis strategy based on time-series prediction and anomaly detection, as shown in Fig. 2. The strategy includes two parts: offline and online. In the offline part, a time-series prediction model is trained based on the history database of the GRS. The online part realizes the fault diagnosis for the GRS, mainly including three parts:

- 1) *Time-Series Prediction*: Input the measurement data of the GRS into the time-series prediction model to obtain the predicted data and residuals.
- 2) *Anomaly Detection*: Calculate the anomaly scores of predicted residuals through relevant algorithms, and set thresholds based on the scores to achieve anomaly detection.
- 3) *Fault Diagnosis*: Based on the anomaly detection results, fault detection and isolation are achieved by referring to Table I, and fault diagnosis is completed.

Subsequently, we will introduce the time-series prediction method and anomaly detection method used in this article through Sections III and IV.

III. ANOMALY-CORRECTING PREDICTION BASED ON INFORMER-RS

Informer-RS model is proposed in this article based on the informer model and an RS mechanism, so the informer and the mechanism are first introduced. Then, an anomaly-correcting prediction method based on Informer-RS is presented to address the issue of abnormal sequences affecting prediction performance.

A. Informer-RS Model

Informer is a deep learning model based on the transformer model, which has an encoder-decoder structure that incorporates self-attention and multihead attention mechanisms. When the transformer is used for time-series prediction, there are still two main challenges to address. One challenge is the high time complexity and memory usage, which are both $O(L^2)$, where L represents the length of the sequence. This leads to a large computational load. The other challenge is the dynamic decoding making the step-by-step inference as slow as RNN-based models. To handle these challenges, the informer mainly introduces the following innovations based on the transformer.

- 1) *Probsparse Self-Attention*: In standard self-attention, the key-value-query mechanism is commonly used, where the similarity between keys and queries is computed to assign corresponding values to all the keys. However, in probsparse self-attention, only the top u most important queries are considered for attention calculation. The probsparse self-attention is defined as

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\overline{\mathbf{Q}}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V} \quad (6)$$

where the number of the dominant query $u = c \ln L_k$, L_k represents the input sequence length for the k th layer, \mathbf{Q} represents the sparse matrix of query, \mathbf{K} represents the matrix of key, \mathbf{V} represents the matrix of value, and d is the input dimension. The sparse method reduces the computational complexity of attention in each layer from $O(L^2)$ to $O(L_k \log L_k)$

- 2) *Distilling*: In the informer, after each probsparse self-attention layer in the encoder and each attention layer in the decoder, there is a convolution layer followed by a max-pooling layer, which halves the computational complexity $O(L_k \log L_k)$ and settles the

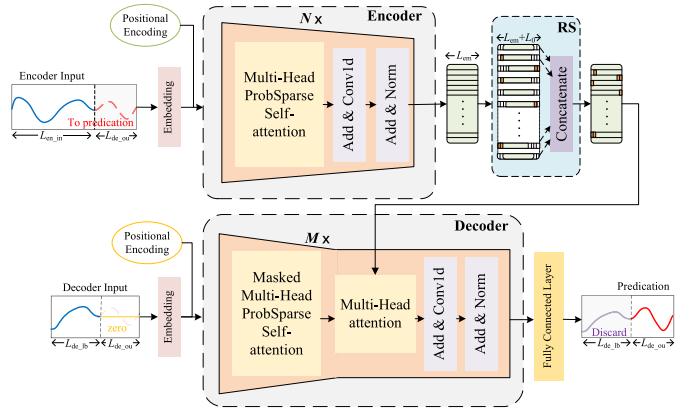


Fig. 3. Structure of the Informer-RS model.

obstacle that the input is too long to stack [38]. The distilling process from layer j to layer $j+1$ is as follows:

$$X_{j+1}^t = \text{MaxPool}\left(\text{ELU}\left(\text{Conv1d}\left([X_j^t]_{\text{att}}\right)\right)\right) \quad (7)$$

where $[\cdot]_{\text{att}}$ represents the attention module, X_j^t is the tensor of layer j , $\text{Conv1d}(\cdot)$ represents a 1-D convolution operation, and $\text{ELU}(\cdot)$ represents the activation function. The step size of $\text{MaxPool}(\cdot)$ is 2, which reduces the computational cost by half.

- 3) *Generative-Style Decoder*: Unlike dynamic decoding, the decoding process is performed through one-forward-step prediction in the informer. This is achieved by concatenating a zero sequence x_0 with a portion x_{token} of the encoder input sequence, which serves as the decoder input

$$x_{\text{de_in}} = \text{Concat}(x_{\text{token}}, x_0). \quad (8)$$

To further enhance the generalization of the model, an RS mechanism is incorporated into the original informer model to obtain the Informer-RS, as shown in Fig. 3. In the RS part of the Informer-RS, the output of the encoder undergoes zero-padding on both sides, resulting in new vectors

$$\bar{x}_{\text{em}} = [\mathbf{0}_{1 \times L_0/2}, x_{\text{em}}, \mathbf{0}_{1 \times L_0/2}] \quad (9)$$

where the length of \bar{x} is $L_{\text{em}} + L_0$. Subsequently, a random section with a length of L_{em} is extracted from the vectors as the new embedding vectors

$$x_{\text{em_new}} = \text{RS}(\bar{x}_{\text{em}}) \quad (10)$$

where $\text{RS}(\cdot)$ represents the RS operation.

Using the RS operation, the issue of overfitting in the model can be effectively mitigated, thereby enhancing the robustness of the model [11]. The main process of Informer-RS is summarized in Algorithm 1. $\text{Embed}(\cdot)$ represents the embedding layer, $\text{PE}(\cdot)$ represents the position encoding, $\text{Encode}(\cdot)$ and $\text{Decode}(\cdot)$ represent the encoder and decoder part, and $\text{FC}(\cdot)$ represents the fully connected layer.

Algorithm 1 Informer

```

Input: The input sequence  $x_{in}$ ;
Output: The prediction sequence  $x_{pred}$ ;
1 function  $f(x_{in})$ ;
2 Encoder input sequence  $x_{em\_in} = x_{in}$ , decoder input
   sequence  $x_{de\_in} = \text{Concat}(x_{token}, x_0)$ 
3 Embed and add position embedding to encoder input
   sequence  $X_{en\_in} = \text{PE}(\text{Embed}(x_{en\_in}))$ 
4 Encoder tensor  $X_{em\_in} = \text{Encode}(X_{en\_in})$ 
5 for  $i = l$  to  $L_{en\_in}$  do
6   Take the  $i$ th row of the encoding tensor  $x_{em} =$ 
      $X_{em\_in}(i)$ 
7   Zero-padding  $\bar{x}_{em} = [\mathbf{0}_{1 \times L_0/2}, x_{em}, \mathbf{0}_{1 \times L_0/2}]$ 
8   RS operation  $x_{em\_new} = \text{RS}(\bar{x}_{em})$ 
9   Fill in the new encoding tensor  $X_{em\_new}(i) =$ 
      $x_{em\_new}$ 
10 end
11 Embed and add position embedding to encoder input
    sequence  $X_{de\_in} = \text{PE}(\text{Embed}(x_{de\_in}))$ 
12 Decode tensor  $X_{em\_out} = \text{Decode}(X_{de\_in}, X_{em\_new})$ 
13 Fully connection  $x_{pred} = \text{FC}(X_{em\_out})$ 
14 return The prediction sequence  $x_{pred}$ .
15 end function

```

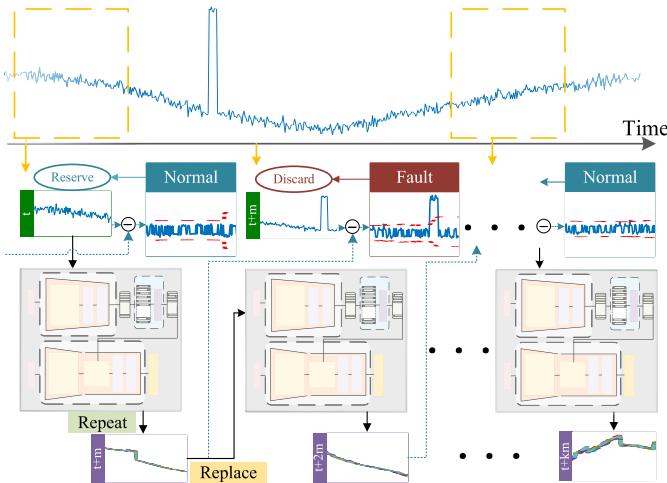


Fig. 4. Flowchart of anomaly-correcting prediction based on the Informer-RS.

B. Anomaly-Correcting Prediction Method

Due to the goal of time-series prediction for fault diagnosis in this study, it is necessary to consider the impact of abnormal data on subsequent predictions. An anomaly-correcting method is used for time-series prediction, as illustrated in Fig. 4. In this figure, Informer-RS is used for predicting the sequence $x = [x_1, x_2, \dots, x_L]^T$, with the input length of L_{en_in} and the prediction length of m ($1 \leq m \leq L_{en_in}$). Because of the RS mechanism in the model, which adds uncertainty during the prediction process, the model can perform multiple repeated predictions, and the average of these predictions is used as the final prediction result to enhance the stability. The anomaly-correcting process in Fig. 4 is as follows.

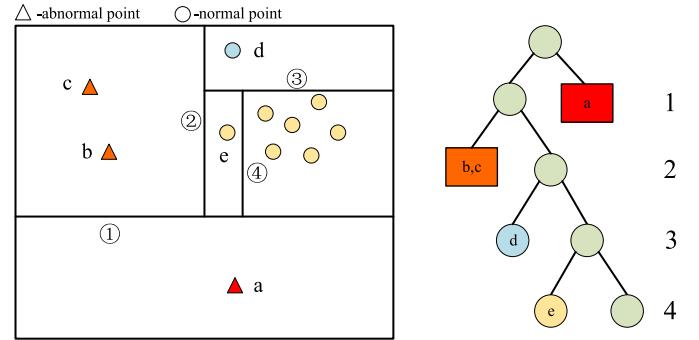


Fig. 5. Principle of the IF algorithm.

Step 1: The predicted mean values $x_{mean}^t = [x_{mean}^{(t+1)}, x_{mean}^{(t+2)}, \dots, x_{mean}^{(t+m)}]^T$ at time t are subtracted from the measurement values $x_{true}^t = [x_{true}^{(t+1)}, x_{true}^{(t+2)}, \dots, x_{true}^{(t+m)}]^T$ at that time to obtain the residuals $e_x^t = [e_x^{(t+1)}, e_x^{(t+2)}, \dots, e_x^{(t+m)}]^T$. The residuals are then analyzed to determine whether there are any anomalies in the measurement sequence.

Step 2: Determine whether the measurement sequence x_{true}^t at time t is abnormal. Upon determining that the sequence is normal, the measurement values are retained and used as part of the input for the next time in the model, and the input at the next time is $x_{en_in}^{t+m} = [x_{true}^{(t-L_{en_in}+m+1)}, \dots, x_{true}^{(t+1)}, \dots, x_{true}^{(t+m)}]^T$.

Step 3: The predicted mean values at time t are subtracted from the measurement values at that time to obtain the residuals.

Step 4: In the measurement sequence at time $t + m$, an abnormality is detected. Therefore, the abnormal part in the measurement values at that time is discarded, and it is replaced with the prediction mean values for anomaly correction. These replaced values are then used as part of the input for the next time, and the input at time $t + 2m$ is $x_{en_in}^{t+2m} = [x_{true}^{(t-L_{en_in}+2m+1)}, \dots, x_{true}^{(t+m+k-1)}, x_{mean}^{(t+m+k)}, \dots, x_{true}^{(t+2m)}]$.

IV. ANOMALY DETECTION METHOD

To avoid the complications of separately analyzing prediction residuals for the displacement and attitude of the TM, an improved IF algorithm is designed to calculate the anomaly scores. In addition, we improve a design method for dynamic thresholds to achieve fault detection, which is adaptive to data changes compared with traditional thresholds.

A. Anomaly Score

The IF algorithm is an ensemble-based, nonparametric, and unsupervised machine learning method widely used for rapid anomaly detection [43]. Unlike some common statistical or clustering-based anomaly detection algorithms, the IF algorithm does not calculate the distance or density of detected points. Instead, it focuses on measuring the ease with which these points can be isolated, effectively reducing the computational complexity of the algorithm.

In the IF algorithm, the forest refers to an ensemble of isolated trees, and each of these trees follows a binary tree structure. In this tree, a random feature is selected to create a binary split on a subsampled subset of the data, and this process is repeated to create two regions at each split. These regions contain fewer and fewer data points, until each region contains only one point or a predefined splitting threshold is reached. The principle of the IF algorithm is shown as Fig. 5. In the figure, labels 1–4 represent binary trees, and points a–e represent different degrees of outliers. Typically, anomaly points are more easily isolated and are closer to the root of the isolation tree, while normal points are harder to be isolated due to their denser distribution, resulting in longer isolation path distances. This distance is used to compute the anomaly score for each data point, and the final score is obtained as the average score from all the isolation trees

$$S = 2^{-\frac{E(h(x))}{c(n)}} \quad (11)$$

where $E(h(x))$ represents the expected path length of node x across all the isolation trees, and $c(n)$ represents the average length of the isolation trees.

Although the low computational cost and suitability for multidimensional data make the IF algorithm suitable for anomaly detection in aerospace sensors, it is challenging to apply it to the fault diagnosis of the GRS. The main reason is that the anomaly score S obtained from the IF algorithm is normalized value, only indicating the degree of the anomaly at the detected point and not providing the information about the direction of the anomaly. Based on the proposed isolation strategy, fault isolation relies on the anomaly direction, so solely depending on S cannot achieve fault isolation. In addition, the IF algorithm is sensitive to outliers, and when faced with some normal points with significant deviation, it may also produce high anomaly scores, leading to false alarms. To address the above challenges, we propose an improved IF algorithm, which computes the new anomaly scores for sequences $X_t = [X^{(t+1)}, X^{(t+2)}, \dots, X^{(t+m)}]^T$ and $Y_t = [Y^{(t+1)}, Y^{(t+2)}, \dots, Y^{(t+m)}]^T$ as follows:

$$Z_t = \sigma_t \cdot \left[\frac{\pi}{10} \arctan(20(S_t - \xi_t - 0.4)) + 0.5 \right] \quad (12)$$

where $\sigma_t = \text{sgn}(X_t) \cdot \text{sgn}(Y_t)$, $\xi_t = \min(S_t)$, S_t represents the anomaly scores generated by the IF algorithm, and σ_t represents the directions of the anomalies. Subtracting ξ_t from the original anomaly scores allows shifting the scores closer to zero. Subsequently, applying the arctangent function can reduce the rate of change of the scores for normal data and increase that for abnormal data. This stabilizes the scores for normal data, making the scores for abnormal data more prominent. The comparison of the anomaly scores in Fig. 6 demonstrates that the proposed improved IF algorithm is more stable and capable of indicating the directions of anomalies.

B. Dynamic Threshold

The design of the thresholds greatly impacts the accuracy of fault detection. If the thresholds are set too low, it may lead to a higher false alarm rate, causing inconvenience for technical personnel. On the other hand, too high thresholds may result

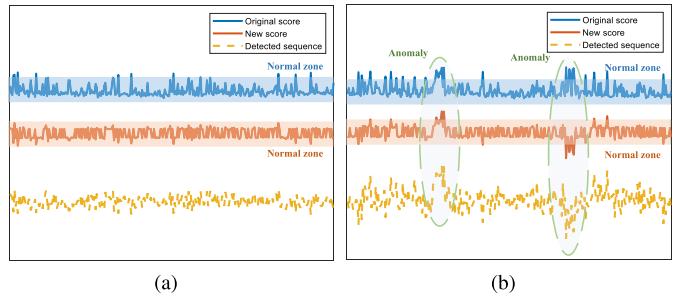


Fig. 6. Comparison of the anomaly scores. (a) Normal data. (b) Abnormal data.

Algorithm 2 Anomaly Scores Calculation and Dynamic Thresholds Design

Input: The augmented residuals \bar{e}_{rx}^t , $\bar{e}_{\phi z}^t$ and the historical anomaly scores Z_B^t ;
Output: The current anomaly scores Z_t and the dynamic thresholds δ_t ;

```

1 function  $g(\bar{e}_{rx}^t, \bar{e}_{\phi z}^t, Z_B^t)$ ;
2 Calculate the augmented anomaly scores  $\bar{Z}^t$  through  $\bar{e}_{rx}^t$  and  $\bar{e}_{\phi z}^t$ 
3 Select the last  $m$  numbers of  $\bar{Z}^t$  as  $Z_t$ 
4 for  $i = l$  to  $m$  do
5 Set  $\bar{Z}_i^t = [Z_B^{(t-h+i)}, \dots, Z_B^{(t-1)}, Z^{(t)}, \dots, Z^{(t+i-1)}]^T$ 
6 Set  $\lambda'_i = \arg \max(F_i(\lambda))$ 
7 Set  $\bar{\delta}_i = \mu(\bar{Z}_i^t) + \lambda'_i \sigma(\bar{Z}_i^t)$  and
    $\underline{\delta}_i = \mu(\bar{Z}_i^t) - \lambda'_i \sigma(\bar{Z}_i^t)$ 
8 Set  $\delta_i = [\bar{\delta}_i, \underline{\delta}_i]^T$ 
9 end
10  $\delta_t = [\delta_1, \delta_2, \dots, \delta_m]^T$ 
11 return The current anomaly scores  $Z_t$  and the dynamic thresholds  $\delta_t$ .
12 end function
```

in a higher rate of missed detections, potentially delaying the discovery of faults and affecting the entire aerospace mission.

The conventional methods assume that the data follow a normal distribution and set the thresholds for anomaly detection based on the 3σ ruler and Z-score method [44]. However, real-world engineering data often do not strictly adhere to a normal distribution. To address this issue, Hundman et al. [23] proposed an unsupervised dynamic threshold design method, which abandons the aforementioned assumption and uses an adversarial method to encourage smaller thresholds while penalizing excessive greedy behavior. As a result, the thresholds are dynamically adjusted based on real-time changes in the data. However, this method requires pruning to reduce the high false alarms caused by the thresholds. Therefore, we make appropriate modifications to this method, eliminating the complex pruning process. The new threshold design method is as follows:

$$\begin{cases} \bar{\delta} = \mu(Z_t) + \lambda' \sigma(Z_t) \\ \underline{\delta} = \mu(Z_t) - \lambda' \sigma(Z_t) \end{cases} \quad (13)$$

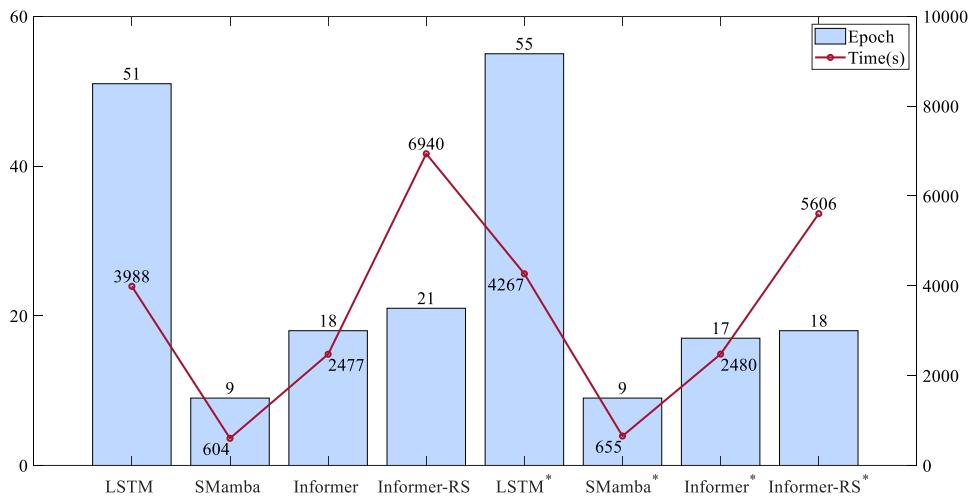


Fig. 7. Train epochs and cost time on each model.

TABLE II
PREDICTION RESULTS OF DIFFERENT MODELS UNDER DIFFERENT DATASETS

Data set	Taiji_data			Taiji_data*			
	Metric	MAE	MSE	Time cost(s)	MAE	MSE	Time cost(s)
LSTM		0.564	0.601	13.46	0.514	0.540	18.46
SMamba		0.482	0.512	3.41	0.405	0.443	3.26
Informer		0.470	0.425	4.03	0.374	0.303	4.69
Informer-RS		0.460	0.409	<u>5.38</u>	0.369	0.301	<u>5.18</u>

¹ The ‘_’ represents the time of a single prediction in repeated predictions.

where $\bar{\delta}$ and $\underline{\delta}$ are the thresholds of sequence \mathbf{Z}_t , and $\mu(\cdot)$ and $\sigma(\cdot)$ represent the mean and standard deviation, respectively. λ' is determined by

$$\lambda' = \arg \max (F(\lambda)) \quad (14)$$

$$F(\lambda) = \frac{(\mu(\mathbf{Z}_t) - \mu(\mathbf{L}_t))/\mu(\mathbf{Z}_t) + (\sigma(\mathbf{Z}_t) - \sigma(\mathbf{L}_t))/\sigma(\mathbf{Z}_t)}{k\|\mathbf{U}_t\|_2 + \varepsilon} \quad (15)$$

where $\mathbf{L}_t = \{\mathbf{Z}_t \in \mathbf{Z}_t \mid \mathbf{Z}_t < \mu(\mathbf{Z}_t) - \lambda\sigma(\mathbf{Z}_t)\}$, $\mathbf{U}_t = \{\mathbf{Z}_t \in \mathbf{Z}_t \mid \mathbf{Z}_t > \mu(\mathbf{Z}_t) + \lambda\sigma(\mathbf{Z}_t)\}$, $k > 0$ and $\varepsilon > 0$ determine the strength of penalizing greedy behavior, and $\lambda \in [2, 5]$ represents the range of dynamic thresholds.

The process of calculating anomaly scores and designing dynamic thresholds is presented in Algorithm 2. To highlight the distinct characteristics of abnormal data, the current sequence is combined with a longer historical sequence during the computation of anomaly scores. Moreover, the current anomaly scores are also considered together with the past anomaly scores during dynamic thresholds design to ensure relatively smooth changes in the thresholds. Here is an explanation for each line of the algorithm.

- 1) In the inputs of the function, the augmented residual sequences $\bar{\mathbf{e}}_{rx}^t = [\mathbf{e}_{rx}^{(t-h+1)}, \dots, \mathbf{e}_{rx}^{(t)}, \mathbf{e}_{rx}^{(t+1)}, \dots, \mathbf{e}_{rx}^{(t+m)}]^T$ and $\bar{\mathbf{e}}_{\phi z}^t = [\mathbf{e}_{\phi z}^{(t-h+1)}, \dots, \mathbf{e}_{\phi z}^{(t)}, \mathbf{e}_{\phi z}^{(t+1)}, \dots, \mathbf{e}_{\phi z}^{(t+m)}]^T$ are formed by combining the current residual sequences \mathbf{e}_{rx}^t and $\mathbf{e}_{\phi z}^t$ with the historical residual sequences, and the historical anomaly scores are denoted by $\mathbf{Z}_B^t = [\mathbf{Z}_B^{(t-h+1)}, \mathbf{Z}_B^{(t-h+2)}, \dots, \mathbf{Z}_B^{(t-1)}]^T$. In the outputs

of the function, the anomaly scores and the dynamic thresholds correspond to the current residual sequences.

- 2) Lines 2 and 3 first calculate the anomaly scores $\bar{\mathbf{Z}}^t$ for the augmented residual sequences and then select the last m values from $\bar{\mathbf{Z}}^t$ to obtain the anomaly scores \mathbf{Z}_t .
- 3) Lines 4–10 design the dynamic thresholds δ_t . Line 5 combines each element in \mathbf{Z}_t with the past $h-1$ anomaly scores to form $\bar{\mathbf{Z}}_i^t$. Lines 6–9 obtain each threshold δ_i by computing λ_i that maximizes $F_i(\lambda)$ related to $\bar{\mathbf{Z}}_i^t$, where $F_i(\lambda)$ is obtained according to (15). Line 10 gathers all δ_i to obtain δ_t .

V. SIMULATION RESULTS AND DISCUSSION

A. Prediction Effect Validation

Considering that the prediction performance of the model directly affects the effectiveness of fault diagnosis, comparative experiments are conducted using the engineering data of class 1C products of “taiji-1” satellite¹ to validate the temporal prediction effectiveness of the proposed model. Because the LSTM has become a relatively mature model in the current space missions, this section selects LSTM, informer, and Informer-RS for comparison. In addition, Mamba, as a new neural network architecture, has received widespread attention in recent years. To demonstrate the superiority of the proposed method, we also chose SMamba [45] as the comparative method. Due to the difficulty in obtaining actual in-orbit data from the GRS of a drag-free satellite, the satellite data triaxial angular data measured by gyroscope A from June 8, 2020 to

¹<https://doi.org/10.57760/sciedb.o00009.00053>

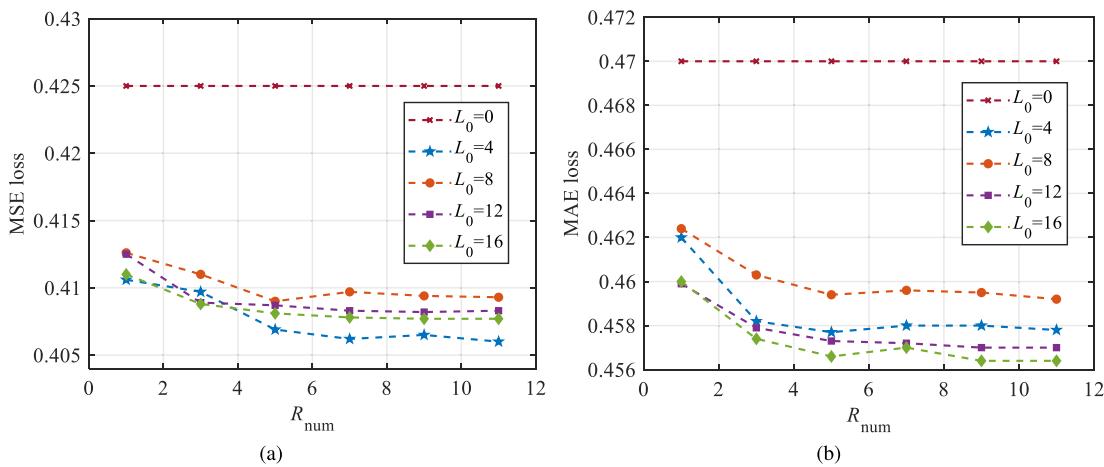


Fig. 8. Prediction results of the Informer-RS in different situations. (a) MSE loss. (b) MAE loss.

TABLE III
PARAMETERS, FLOPS, AND GPU MEMORY OF DIFFERENT MODELS

Model	LSTM	SMamba	Informer	Informer-RS
Parameters	2.00M	17.82M	19.45M	19.45M
FLOPs	1.33G	0.036G	9.43G	9.43G
GPU Memory	362M	840M	520M	628M

June 17, 2020 in the engineering data are selected as the dataset Taiji_data, and the sampling frequency is 1 Hz. To further analyze the influence of noise on different models, the data of Taiji_data, which already contain noticeable noise, are subjected to mean filtering to obtain a new dataset, Taiji_data*. Each dataset comprises 795 082 sets of data, divided into training, validation, and test sets in a ratio of 6:2:2.

All the experiments in this article are carried out on a personal laptop with an i9-12950HX CPU, RTX A4500 GPU, 128 GB of RAM, and 16G of video memory. All the models have an input sequence length of 720, a prediction sequence length of 300, and a lookback window length of 300, and both the input and output data dimensions were set to 3. The hyperparameters for the informer and LSTM are set according to the optimal settings of Zhou et al. [34]. The hyperparameters for the SMamba as set as given in Wang et al. [45]. For the Informer-RS, the hyperparameters are the same as the informer. In addition, the zero-padding length \$L_0\$ and the number of repeated predictions \$R_{\text{num}}\$ in the RS part of the Informer-RS will affect the prediction performance. The default settings are \$L_0 = 8\$ and \$R_{\text{num}} = 5\$, respectively. During the training process, Adam is set as the optimizer, mean squared error (mse) is chosen as the loss function, the batch size is 32, and the number of epochs is 100, with a learning rate of 0.0001. Adaptive decay of the learning rate based on the validation set loss is implemented, and the training will be terminated early if the validation set loss does not decrease for five consecutive epochs.

With the specified settings, the four models are trained on two datasets, and the number of train epochs and cost time for each model are illustrated in Fig. 7. The * indicates the models are trained by Taiji_data*. From the training process, it is evident that the Informer-RS requires the longest training

TABLE IV
PREDICTION TIME UNDER DIFFERENT \$R_{\text{num}}

\$R_{\text{num}}	1	3	5	7	9	11
Total time(s)	5.56	16.32	26.91	37.42	48.35	59.13
Single time(s)	5.56	5.44	5.38	5.35	5.37	5.38

time, with the increase primarily attributed to the time spent on the RS part. After completing the training, 500 sets of data from test set are used to analyze the prediction accuracy and time cost for each model. The prediction accuracy is evaluated using both mse and mean absolute error (MAE), and the total time cost of predicting 500 time series is counted, as presented in Table II. In the table, the Informer-RS exhibits the highest prediction accuracy, and SMamba has the shortest prediction time. Due to the RS part, the Informer-RS has a slightly longer prediction time than the informer, but still significantly less time than the LSTM. Meanwhile, the results also indicate that noise will affect the prediction performance of the models.

Considering the critical importance of computational complexity and memory usage in space missions, we compare the parameters, floating point operations (FLOPs), and GPU memory of four models. The results are presented in Table III. The LSTM, being relatively simple in structure, has the smallest parameters and GPU memory. The SMamba benefits from the Mamba architecture and linearized design, achieving extremely low FLOPs, but its parameters and GPU memory remain relatively high. The informer and the Informer-RS exhibit the highest parameters and FLOPs, while their GPU memory is moderate. In fact, FLOPs do not directly equate to computational cost. The recursive prediction method of the LSTM results in it having the highest computational cost among the models. Although the informer and the Informer-RS have the largest FLOPs, their computation time still competes with the SMamba according to the results in Table II.

To further analyze the impact of \$L_0\$ and \$R_{\text{num}}\$ on the prediction performance of the Informer-RS, comparative experiments are conducted using the test data of the Taiji_data dataset, and the results are presented in Fig. 8. \$L_0 = 0\$ indicates that the model is the informer.

In Fig. 8, the results indicate that the prediction accuracy of the Informer-RS in different situations is much higher than that

TABLE V
PREDICTION RESULTS UNDER DIFFERENT SAMPLING RATES

Sampling rate(Hz)	0.1	0.5	1	5	10
MSE	0.488	0.432	0.409	0.346	0.175
MAE	0.508	0.479	0.460	0.385	0.250
Prediction duration(s)	3000	600	300	60	30

TABLE VI
HYPERPARAMETERS' SETTINGS OF INFORMER-RS

Parameters	Description	Value
L_{en_in}	encoder input length	720
L_{de_ou}	decoder output length	300
L_{de_lb}	lookback window length	300
S_{en_in}	encoder input size	2
S_{de_in}	decoder input size	2
S_{de_out}	decoder output size	2
N	layers of encoder	3
M	layers of decoder	2
N_{head}	number of heads	8
L_{em}	embedding length	512
dropout	dropout	0.1

of the informer. Furthermore, the prediction accuracy of the Informer-RS shows an initial improvement with the increase in R_{num} , followed by stabilization, while there is no apparent linear relationship between the zero-padding length and the prediction accuracy. In addition, under different metrics, the prediction accuracy varies with different L_0 . This is attributed to mse being more sensitive to outliers compared to MAE, and shorter zero-padding lengths may result in smaller potential deviations of predicted values, reducing the likelihood of outliers. Although a larger R_{num} may lead to higher prediction accuracy, the increased time cost associated with repeated predictions needs to be considered, particularly in the context of fault diagnosis. Table IV shows the relationship between R_{num} and the time cost of predicting 500 series, and it is evident that as R_{num} increases, the cost time will significantly increase.

In the proposed approach, the selection of the sampling rate is equally critical. With a fixed prediction sequence length, a higher sampling rate corresponds to a longer prediction duration, which reduces computational cost. However, it also results in more frequent variations in the predicted sequence, potentially degrading prediction accuracy and fault diagnosis performance. Table V presents the prediction results under different sampling rates, which validate the above conclusions. Therefore, a sampling rate of 1 Hz is chosen in this study as a balanced and well-considered option.

B. Fault Diagnosis Simulation

Some simulation experiments are conducted to validate the effectiveness of the proposed fault diagnosis approach in this study. Currently, GRS technology is still not sufficiently mature worldwide, with a large number of relevant experiments only conducted in the Lisa pathfinder (LPF) task. However, due to some objective reasons, we are unable to access these experimental data. Therefore, to maximize the credibility of our simulation experiments in terms of practical applicability, we have rigorously constructed a MATLAB sim-

ulation model² based on the GRS sensing principles described in Section II, without any approximations or linearizations. In addition, all the parameters in the model are consistent with those used in the GRS of the LPF task [46], [47]. Further considering noise has a significant impact on the performance of time-series prediction, which will directly affect the effectiveness of fault diagnosis, the colored noise model for the sensor measurement process is also developed based on relevant research on colored noise in the context of the LPF task [5]. Besides this colored noise, all the subsequent noises mentioned are zero-mean Gaussian white noise.

In the experiments, except for the input and output data dimensions, the other hyperparameter settings of the Informer-RS are the same as those in Section V-A, as shown in Table VI. Similarly, the training process parameters are also not adjusted. Based on the results shown in Fig. 8, we selected $L_0 = 4$ in the Informer-RS, as this configuration achieves the lowest mse. This choice aligns with the need for fewer outliers in predictions to enhance fault diagnosis performance. In addition, we set $R_{num} = 7$ because, beyond this value, the improvement in prediction accuracy becomes negligible, and the computational cost remains reasonably low. The parameters of the improved IF algorithm follow the default settings of the IF algorithm, with the subsample size of 256 and the number of trees of 100. The parameters for setting the dynamic threshold are based on empirical experience. The values of k and ϵ are set to 10 and 1×10^{-7} , respectively, and the length after splicing the historical sequence is $h = 300$.

Considering that Perosanz et al. [4] mentioned that the GRS has experienced frequent data jumping anomalies, we subject the GRS to intermittent faults to analyze the diagnosis effectiveness of the proposed method. Intermittent faults are a type of fault characterized by a short duration and unpredictable patterns. The type of fault is difficult to detect, and if left unaddressed for an extended period, it can easily lead to more serious malfunctions. The specific form of faults follows the mathematical model proposed by Cui et al. [19], and two types of intermittent faults are considered in this study, namely, the bias fault and oscillation fault. These faults are applied to the measurement data of different DCSs through simulation injection.

First, verify the effectiveness of the proposed fault diagnosis approach for the GRS. Based on the approach, fault diagnosis is conducted on sequences of length 1500 s. Only one type of intermittent fault is injected into the sequences at a time, and the specific fault settings are shown in Table VII. The diagnosis results for two types of intermittent faults are presented in Figs. 9 and 10. To clearly demonstrate the effectiveness of fault diagnosis, we have set up larger faults here.

The results of the bias faults are shown in Fig. 9. From Fig. 9(a) and (b), the Informer-RS accurately predicts the measurement data, and the use of the anomaly-correcting method prevents the abnormal sequences from affecting future sequence predictions. In Fig. 9(c), the anomaly scores and the dynamic thresholds are shown, and the residual sequences have undergone translation and scaling processes to simul-

²https://github.com/VERYBC/GRS_model

TABLE VII
FAULT SETTING

Type	Occurrence time (s)	Duration (s)	Size ($10^{-17}C$)	Location	Mathematical model
Bias fault	114	12	40	DCS1	$F(t) = \chi(t) + \xi(t)$
	769	6	30	DCS2	$\chi(t)$ is a normal value, $\xi(t)$ is a constant.
	1324	9	20	DCS1	
Oscillation fault	164	6	30	DCS1	$F(t) = \chi(t) + \delta(t)$
	870	11	24	DCS2	$\chi(t)$ is a normal value, $\delta(t)$ follows $N(0, \sigma^2)$.
	1424	8	26	DCS2	

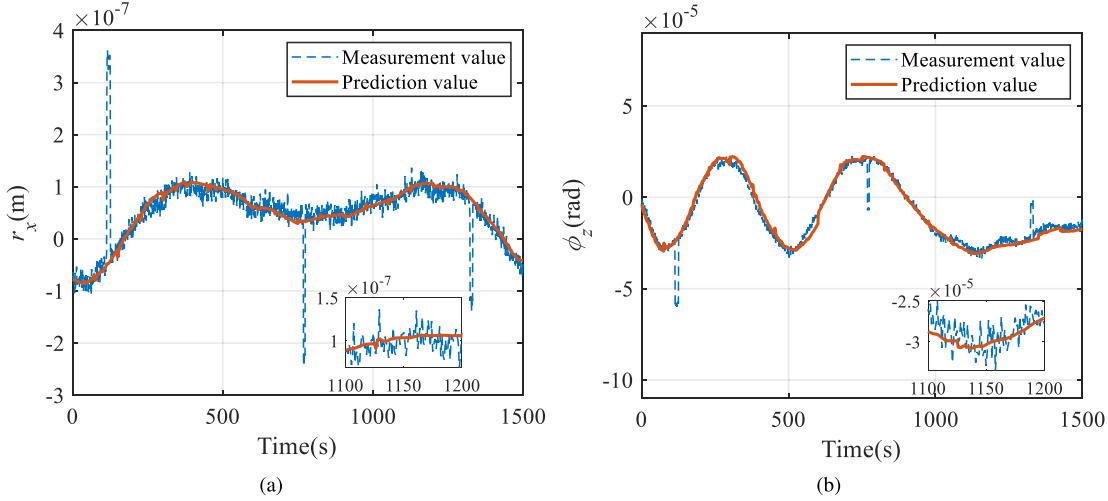


Fig. 9. Fault detection results for bias faults. (a) Measurement and prediction curves of r_x . (b) Measurement and prediction curves of ϕ_z . (c) Anomaly scores and dynamic thresholds. (d) Contour map of anomaly scores.

taneously demonstrate the variations of their curves. Three faults accurately are detected in this figure, and the fault locations can be determined by the variations in the scores. Fig. 9(d) further shows that the distribution of the data and the ability of the proposed anomaly scores to characterize different anomalies. All the abnormal points of the bias fault can be detected by both dynamic and static thresholds, with the static thresholds set using the 3σ rule.

The results of the oscillation faults are shown in Fig. 10. The difference between oscillation faults and bias faults lies in the fact that some abnormal points of the oscillation fault may be distributed within the range of normal points, making them difficult to be detected. As shown in Fig. 10(d),

the results of the static thresholds show substantial missed detections. Compared with the static thresholds, the dynamic thresholds are able to detect the majority of the outliers. However, several abnormal points within the normal range still cannot be detected by the dynamic thresholds.

To further analyze the fault diagnosis capability of the proposed approach, Monte Carlo experiments are conducted. In the experiments, the randomness of some fault characteristics is taken into account. The experimental settings are shown in Table VIII. In the table, F_b represents the bias fault, and F_o represents the oscillation fault. The values of -1 and 1 for fault direction represent subtractive and additive fault, respectively. A fault detection criterion is applied: if at least

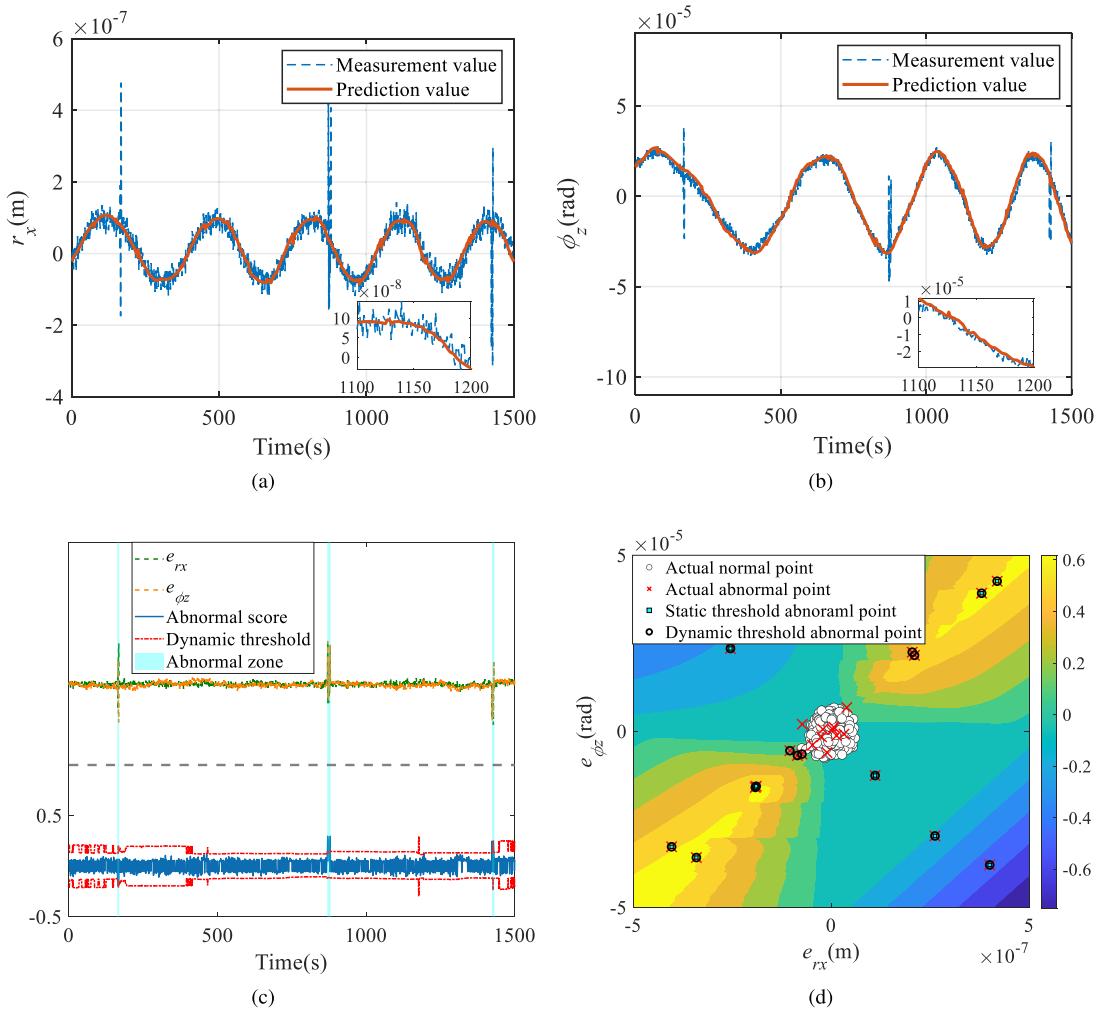


Fig. 10. Fault detection results for oscillation faults. (a) Measurement and prediction curves of r_x . (b) Measurement and prediction curves of ϕ_z . (c) Anomaly scores and dynamic thresholds. (d) Contour map of anomaly scores.

TABLE VIII
MONTE CARLO EXPERIMENT SETTINGS

Parameters	Value
Number of experiments	600
Sequence length (s)	300
Fault type	$\{F_b, F_o\}$
Fault location	$\{\text{DCS}_1, \text{DCS}_2\}$
Fault occurrence time (s)	[1, 281]
Fault duration (s)	{5, 6, ..., 20}
Fault size (10^{-17}C)	[10, 40]
Fault direction	{-1, 1}

four abnormal points are detected within a window of length 5, the current time series is deemed to have anomalies, indicating the occurrence of faults.

The evaluation metrics of fault detection include common classification metrics such as the accuracy, precision, and recall, which are defined as follows:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (16)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (17)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}. \quad (18)$$

In the fault diagnosis problem of this article, TP represents the number of samples detected as faults among the actual fault samples, TN represents the number of samples detected as normal among the actual normal samples, FP represents the number of samples detected as normal among the actual fault samples, and FN represents the number of samples detected as faults among the actual normal samples. Accuracy reflects the accuracy of fault detection, while precision and recall are related to missed detections and false alarms, respectively, and higher score indicate better detection performance.

In typical classification problems, the $F1$ score is commonly used to balance precision and recall, but it treats precision and recall with equal weight. This is not suitable for fault diagnosis tasks, as the severity of missed detection is greatly higher than that of false alarm in such cases. In other words, the importance of precision outweighs that of recall. To address this issue, we propose a new metric called the $B1$ score to instead of $F1$ score. The $B1$ score is formulated as follows:

$$B1 = \frac{\text{Precision} \times \text{Recall}}{(1 - \alpha) \times \text{Precision} + \alpha \times \text{Recall}} \quad (19)$$

TABLE IX
COMPARISON RESULTS OF DIFFERENT METHODS

Methods	Accuracy	Precision	Recall	B1	S1	Cost time(s)
Informer-RS+iIF+DT	0.9783	0.9677	1.0000	0.9772	0.9000	2.1258/ <u>1.6054</u>
Informer-RS+iIF+ST	0.8883	0.8337	1.0000	0.8775	0.6845	2.1230/ <u>1.6026</u>
Informer-RS+IF+DT	0.7567	1.0000	0.7341	0.9020	0.7322	2.1148/ <u>1.5944</u>
Informer-RS+IF+ST	0.9750	0.9628	1.0000	0.9737	0.8660	2.1124/ <u>1.5920</u>
Informer+iIF+DT	0.9550	0.9454	0.9870	0.9575	0.8083	1.6048
Informer+iIF+ST	0.8567	0.7866	1.0000	0.8404	0.6498	1.6020
Informer+IF+DT	0.7783	0.9950	0.7538	0.9079	0.7143	1.5938
Informer+IF+ST	0.8650	0.7990	1.0000	0.8503	0.7453	1.5914
SMamba+iIF+DT	0.9617	0.9429	1.0000	0.9594	0.8105	1.6041
SMamba+iIF+ST	0.8300	0.7469	1.0000	0.8083	0.5914	1.6013
SMamba+IF+DT	0.9367	0.9876	0.9234	0.9674	0.8979	1.5931
SMamba+IF+ST	0.9383	0.9082	1.0000	0.9339	0.7896	1.5907
LSTM+iIF+DT	0.7583	0.6402	1.0000	0.7177	0.1705	1.6268
LSTM+iIF+ST	0.4767	0.2208	1.0000	0.2882	0.0562	1.6240
LSTM+IF+DT	0.9200	0.8809	1.0000	0.9135	0.5775	1.6158
LSTM+IF+ST	0.8633	0.7965	1.0000	0.8484	0.3614	1.6134

¹The ‘_’ represents the fault diagnosis time with only one prediction.

TABLE X
ACCURACY AND *B1* OF DIFFERENT MODELS IN DIFFERENT NOISE

Noise (10^{-17} C)	Informer-RS		Informer		SMamba		LSTM	
	Index	Accuracy	B1	Accuracy	B1	Accuracy	B1	Accuracy
0	0.9700	0.9700	0.9900	0.9901	0.9600	0.9568	0.7600	0.7367
1	0.9800	0.9798	0.9600	0.9696	0.9400	0.9418	0.6600	0.6098
2	0.9200	0.9207	0.9100	0.9104	0.8500	0.8235	0.5700	0.4992
3	0.8200	0.7850	0.8100	0.7719	0.8200	0.7649	0.4900	0.2669
4	0.7400	0.6961	0.7200	0.6766	0.7200	0.6475	0.4500	0.2662
5	0.6600	0.5810	0.6600	0.5856	0.6100	0.5396	0.4200	0.1815
Count	9		3		0		0	

where α represents the proportion of precision, $\alpha \in (0.5, 1)$.

Furthermore, due to the set fault detection criterion, the above metrics cannot represent the accuracy of the approach in detecting abnormal points. Therefore, a new metric called *S1* score is introduced in this study to represent the accuracy of outlier detection

$$S1 = \frac{K}{TP + FP} \quad (20)$$

where K represents the number of samples with accurately detected outliers, which are defined as follows: sequences with bias faults having over 90% of outliers detected, or sequences with oscillation faults having over 60% of outliers detected.

Similarly, the informer, the SMamba, and the LSTM are selected as the baseline models, with the hyperparameters of the informer set the same as those in Table VI. Most of the hyperparameters of the SMamba refer to Wang et al. [45], but the number of encoder layers is set to 5. The hyperparameters of the LSTM still refer to the settings of Zhou et al. [34], where the number of layers is set to 4 and the hidden size is set to 256. Based on the aforementioned evaluation criterion, where $\alpha = 0.7$ is set, the results of the Monte Carlo experiments for different methods are presented in Table IX. In this table, iIF refers to the improved IF algorithm, and DT and ST represent dynamic threshold and static threshold, respectively. The results indicate that the proposed method exhibits the best fault detection performance, attributed to the outstanding

prediction accuracy of the Informer-RS and the stable anomaly scores provided by the iIF method and dynamic threshold setting method. However, the combination of iIF and ST methods leads to a lower precision, indicating an increase in false negatives, mainly due to the fact that ST method often has a larger threshold range than DT method when facing stable data. The combination of IF and DT methods yields high precision but lower recall, suggesting that this outcome is a result of a substantial number of false positives. The drastic fluctuation in anomaly scores generated by IF method makes the dynamic threshold less stable. On the other hand, the combination of IF and ST methods for the informer and Informer-RS shows better fault detection performance compared with the aforementioned two cases, as the unstable anomaly scores paired with a stable threshold make the method relatively more effective. Furthermore, although the SMamba performs less effectively than the informer in temporal prediction, its superior short-term tracking capability leads to smaller prediction variance, resulting in better fault diagnosis performance compared with the informer. It is noteworthy that due to the poor prediction performance of the LSTM, the use of iIF method results in smaller differences in anomaly scores between normal and abnormal data compared with the IF method. The smaller differences can lead to more true negatives and poorer outlier detection performance, which is strongly reflected by the *S1* scores. Moreover, the time spent

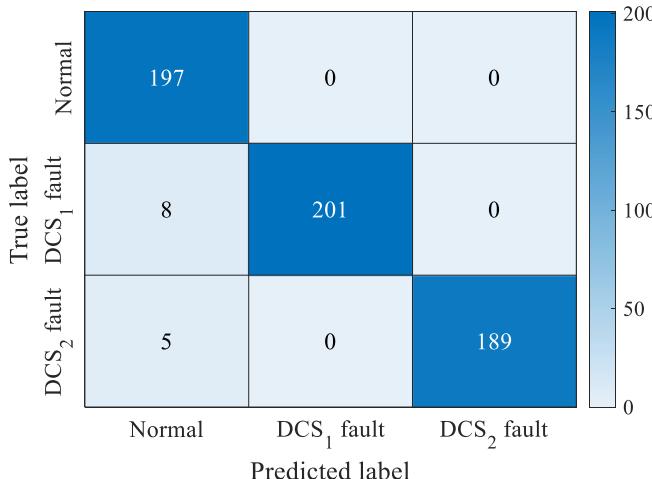


Fig. 11. Confusion matrix of fault isolation.

by all four models is considerably less than the length of the detected sequence (300 s), indicating that the proposed method can meet online prediction requirements. In addition, the increased computational load introduced by iIF and DT methods is minimal compared with IF and ST methods. Notably, the SMamba has the shortest time, while the Informer-RS has the longest time, primarily due to repeated predictions. Reducing this number could bring its computational load below that of the LSTM. The specific number of repeated predictions should be adjusted based on engineering requirements.

The isolation results of the proposed approach in the above experiments are presented in Fig. 11. According to the confusion matrix for fault isolation, only 13 sequences encounter missed detections out of the 600 sequences. Among the detected 390 abnormal sequences, the isolation accuracy reaches 100%, indicating that all faulty DCSs can be successfully isolated.

Finally, to validate the improvement of model stability by introducing the RS part, we conduct 100 Monte Carlo experiments by adding noise of various sizes to the differential capacitance measurement data. The method of combining iIF and DT is used in anomaly detection, and other parameters are kept the same as listed in Table VIII. The robustness of the four models is compared by analyzing their fault diagnosis effectiveness in the face of different noises, and the results are shown in Fig. 12 and Table X. The results in Fig. 12 show that at lower noise levels, the accuracy of outlier detection in the Informer, the SMamba, and the Informer-RS can be compared. However, as the noise increases, the accuracy of the Informer and the SMamba diminishes compared with the Informer-RS, which shows that the RS part can improve the robustness of the Informer model in the presence of noise, and the SMamba demonstrates relatively poor robustness. The accuracy of outlier detection in the LSTM is significantly inferior to the other three models.

The accuracy and $B1$ scores of the four models under different noise levels are summarized in Table X, indicating that the Informer-RS achieves better performance than the other models in fault detection. The Informer maintains a good fault detection performance when the noise is low, even

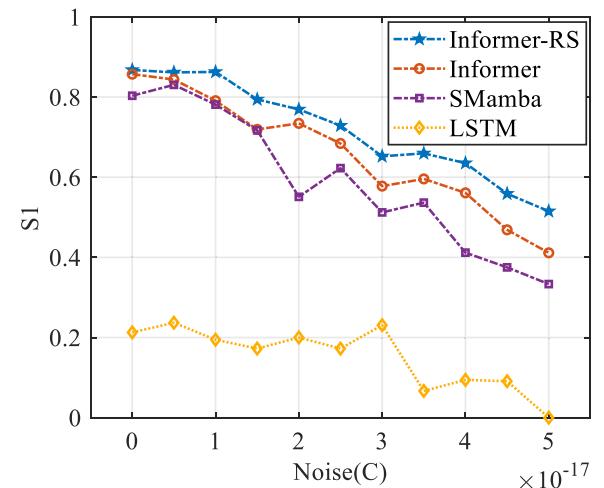


Fig. 12. $S1$ of different models in different noise.

outperforming the Informer-RS in individual cases. With the increasing level of noise, the fault detection performance of the Informer starts to lag behind the Informer-RS, until the noise reaches a sufficiently high level, at which point both the models achieve a comparable level of performance. However, the performance of SMamba is consistently relatively inferior. Considering the results in Fig. 12, despite the higher $S1$ of the Informer-RS, it exhibits a fault detection performance similar to the Informer. This could be attributed to the increased false positives of the Informer after the introduction of higher noise, resulting in a decent fault detection performance in terms of data. This aspect can potentially be improved by further adjusting the fault detection criteria. In summary, the RS mechanism can enhance the fault diagnosis performance of the Informer model to some extent in the presence of noise, and the Informer-RS model combined with repeated predictions possesses stronger robustness.

VI. CONCLUSION

In conclusion, this article presents a dynamic fault diagnosis method, namely, Informer-RS, tailored for GRS in the space gravitational wave detection mission. The proposed method offers several key contributions and enhancements. First, the pose sensing principle of the GRS is elaborated, and based on this principle, a possible type of fault and its diagnosis strategy are proposed. Then, the Informer model is enriched by incorporating an RS mechanism to improve its robustness. In addition, an anomaly-correcting prediction method is introduced to mitigate the impact of faults in prediction of time series. To facilitate the fault diagnosis, an improved IF algorithm and a dynamic threshold design method are applied. Through comparative experiments based on actual in-orbit satellite data, the prediction performance of the Informer-RS has been validated. Furthermore, the proposed approach is evaluated through simulations of two intermittent faults and compared with several common approaches using Monte Carlo statistical analysis. The results demonstrate that the proposed approach can accurately diagnose faults and exhibit superior fault diagnosis performance.

Nevertheless, it is important to acknowledge that while the proposed approach exhibits improved robustness, its performance may still be adversely affected when the measurement data are heavily influenced by noise. Thus, our future work will focus on further enhancing the model's ability to handle unexpected noise and interference, aiming to achieve even higher fault diagnosis accuracy and reliability. In addition, the proposed method has difficulty identifying simultaneous faults in DCSs on the same axis, which will be one of the key directions for our future research. We will continue to monitor the development of GRS technology worldwide, so that we can conduct more comprehensive testing of our proposed method as soon as actual in-orbit data becomes available.

REFERENCES

- [1] H. Liu, X. Niu, M. Zeng, S. Wang, K. Cui, and D. Yu, "Review of micro propulsion technology for space gravitational waves detection," *Acta Astronautica*, vol. 193, pp. 496–510, Apr. 2022.
- [2] S. Vidano, C. Novara, L. Colangelo, and J. Grzymisch, "The LISA DFACS: A nonlinear model for the spacecraft dynamics," *Aerospace Sci. Technol.*, vol. 107, Dec. 2020, Art. no. 106313.
- [3] G. Wang, W.-T. Ni, W.-B. Han, P. Xu, and Z. Luo, "Alternative LISA-TAIJI networks," *Phys. Rev. D, Part. Fields*, vol. 104, no. 2, Jul. 2021, Art. no. 024012.
- [4] F. Perosanz et al., "On board evaluation of the STAR accelerometer," in *First CHAMP Mission Results for Gravity, Magnetic and Atmospheric Studies*. Berlin, Germany: Springer, 2003, pp. 11–18.
- [5] M. Armano et al., "Capacitive sensing of test mass motion with nanometer precision over millimeter-wide sensing gaps for space-borne gravitational reference sensors," *Phys. Rev. D, Part. Fields*, vol. 96, no. 6, Sep. 2017, Art. no. 062004.
- [6] S. Wang, L. Chen, Y. Wang, Z. Zhou, K. Qi, and Z. Wang, "A space inertial sensor ground evaluation system for non-sensitive axis based on torsion pendulum," *Appl. Sci.*, vol. 10, no. 9, p. 3090, Apr. 2020.
- [7] P. M. Frank, "Fault diagnosis in dynamic systems using analytical and knowledge-based redundancy: A survey and some new results," *Automatica*, vol. 26, no. 3, pp. 459–474, May 1990.
- [8] X. Chen, R. Sun, M. Liu, and D. Song, "Two-stage exogenous Kalman filter for time-varying fault estimation of satellite attitude control system," *J. Franklin Inst.*, vol. 357, no. 4, pp. 2354–2370, Mar. 2020.
- [9] N. Tantouris and K. Dellios, "Euphoria-filter: A robust H_2/H_∞ technique amenable to the LISA-Pathfinder thruster faults," *Int. J. Space. Sci. Eng.*, vol. 4, no. 1, pp. 45–63, Aug. 2016.
- [10] S. Gao, W. Zhang, and X. He, "Observer-based multiple faults diagnosis scheme for satellite attitude control system," *Asian J. Control*, vol. 22, no. 1, pp. 307–322, Jan. 2020.
- [11] T. Luo, M. Liu, H. Zhao, G. Duan, and X. Cao, "Data-driven fault monitoring for spacecraft control moment gyro with slice residual attention network," *J. Franklin Inst.*, vol. 359, no. 16, pp. 9313–9333, Nov. 2022.
- [12] T. Zhang et al., "Intelligent fault diagnosis of machines with small imbalanced data: A state-of-the-art review and possible extensions," *ISA Trans.*, vol. 119, pp. 152–171, Jan. 2022.
- [13] H. Nizam, S. Zafar, Z. Lv, F. Wang, and X. Hu, "Real-time deep anomaly detection framework for multivariate time-series data in industrial IoT," *IEEE Sensors J.*, vol. 22, no. 23, pp. 22836–22849, Dec. 2022.
- [14] B. Yang, S. Xu, Y. Lei, C.-G. Lee, E. Stewart, and C. Roberts, "Multi-source transfer learning network to complement knowledge for intelligent diagnosis of machines with unseen faults," *Mech. Syst. Signal Process.*, vol. 162, Jan. 2022, Art. no. 108095.
- [15] R. Fujimaki, T. Yairi, and K. Machida, "An approach to spacecraft anomaly detection problem using kernel feature space," in *Proc. 11th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, Aug. 2005, pp. 401–410.
- [16] Y. K. Takehisa Yairi, "Telemetry-mining: A machine learning approach to anomaly detection and fault diagnosis for space systems," in *Proc. 2nd IEEE Int. Conf. Space Mission Challenges Inf. Technol. (SMC-IT)*, Aug. 2006, pp. 466–476.
- [17] K. Li, Y. Wu, S. Song, Y. Sun, J. Wang, and Y. Li, "A novel method for spacecraft electrical fault detection based on FCM clustering and WPSVM classification with PCA feature extraction," *Proc. Inst. Mech. Eng., G, J. Aerosp. Eng.*, vol. 231, no. 1, pp. 98–108, Jan. 2017.
- [18] Y. Gao, T. Yang, M. Xu, and N. Xing, "An unsupervised anomaly detection approach for spacecraft based on normal behavior clustering," in *Proc. 5th Int. Conf. Intell. Comput. Technol. Autom.*, Jan. 2012, pp. 478–481.
- [19] L. Cui et al., "A method for satellite time series anomaly detection based on fast-DTW and improved-KNN," *Chin. J. Aeronaut.*, vol. 36, no. 2, pp. 149–159, Feb. 2023.
- [20] J. Qin, L. Wang, and R. Huang, "Research on fault diagnosis method of spacecraft solar array based on f-KNN algorithm," in *Proc. Prognostics Syst. Health Manage. Conf. (PHM-Harbin)*, Jul. 2017, pp. 1–4.
- [21] M. Goldstein and S. Uchida, "A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data," *PLoS ONE*, vol. 11, no. 4, Apr. 2016, Art. no. e0152173.
- [22] Z. Zeng, G. Jin, C. Xu, S. Chen, Z. Zeng, and L. Zhang, "Satellite telemetry data anomaly detection using causal network and feature-attention-based LSTM," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–21, 2022.
- [23] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Söderström, "Detecting spacecraft anomalies using LSTMs and non-parametric dynamic thresholding," in *Proc. 24th ACM SIGKDD Int. Conf. Knowl. Disc. Data Min.*, Jul. 2018, pp. 387–395.
- [24] B. Chen et al., "Continual learning fault diagnosis: A dual-branch adaptive aggregation residual network for fault diagnosis with machine increments," *Chin. J. Aeronaut.*, vol. 36, no. 6, pp. 361–377, Jun. 2023.
- [25] S. Ai, J. Song, and G. Cai, "A real-time fault diagnosis method for hypersonic air vehicle with sensor fault based on the auto temporal convolutional network," *Aerospace Sci. Technol.*, vol. 119, Dec. 2021, Art. no. 107220.
- [26] X.-G. Guo, M.-E. Tian, Q. Li, C. K. Ahn, and Y.-H. Yang, "Multiple-fault diagnosis for spacecraft attitude control systems using RBFNN-based observers," *Aerospace Sci. Technol.*, vol. 106, Nov. 2020, Art. no. 106195.
- [27] H.-J. Jin, Y.-P. Zhao, and Z.-Q. Wang, "A rotating stall warning method for aero-engine compressor based on DeepESVDD-CNN," *Aerospace Sci. Technol.*, vol. 139, Aug. 2023, Art. no. 108411.
- [28] M. ElDali and K. D. Kumar, "Fault diagnosis and prognosis of aerospace systems using growing recurrent neural networks and LSTM," in *Proc. IEEE Aerosp. Conf.*, Jun. 2021, pp. 1–20.
- [29] J. Chen, D. Pi, Z. Wu, X. Zhao, Y. Pan, and Q. Zhang, "Imbalanced satellite telemetry data anomaly detection model based on Bayesian LSTM," *Acta Astronautica*, vol. 180, pp. 232–242, Mar. 2021.
- [30] S.-T. Yun and S.-H. Kong, "Data-driven in-orbit current and voltage prediction using bi-LSTM for LEO satellite lithium-ion battery SOC estimation," *IEEE Trans. Aerospace Electron. Syst.*, vol. 58, no. 6, pp. 5292–5306, Dec. 2022.
- [31] P. Malhotra, A. Ramakrishnan, G. Anand, L. Vig, P. Agarwal, and G. Shroff, "LSTM-based encoder-decoder for multi-sensor anomaly detection," 2016, *arXiv:1607.00148*.
- [32] P. Malhotra, L. Vig, G. Shroff, and P. Agarwal, "Long short term memory networks for anomaly detection in time series," in *Proc. Eur. Symp. Artif. Neural Netw. (ESANN)*, 2015, p. 89.
- [33] M. Sirajul Islam and A. Rahimi, "Fault prognosis of satellite reaction wheels using a two-step LSTM network," in *Proc. IEEE Int. Conf. Prognostics Health Manage. (ICPHM)*, Jun. 2021, pp. 1–7.
- [34] H. Zhou et al., "Informer: Beyond efficient transformer for long sequence time-series forecasting," in *Proc. AAAI Conf. Artif. Intell.*, vol. 35, no. 12, May 2021, pp. 11106–11115.
- [35] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, Jun. 2017, pp. 5998–6008.
- [36] N. Wu, B. Green, X. Ben, and S. O'Banion, "Deep transformer models for time series forecasting: The influenza prevalence case," 2020, *arXiv:2001.08317*.
- [37] R. M. Farsani and E. Pazouki, "A transformer self-attention model for time series forecasting," *J. Electr. Comput. Eng. Innov. (JECEI)*, vol. 9, no. 1, pp. 1–10, 2020.
- [38] M. Gong, Y. Zhao, J. Sun, C. Han, G. Sun, and B. Yan, "Load forecasting of district heating system based on informer," *Energy*, vol. 253, Aug. 2022, Art. no. 124179.
- [39] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," 2020, *arXiv:2004.05150*.

- [40] S. Wang, B. Z. Li, M. Khabsa, H. Fang, and H. Ma, "LinFormer: Self-attention with linear complexity," 2020, *arXiv:2006.04768*.
- [41] H. Wu, J. Xu, J. Wang, and M. Long, "Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting," in *Proc. NIPS*, vol. 34, Dec. 2021, pp. 22419–22430.
- [42] H. Inchauspé et al., "Numerical modeling and experimental demonstration of pulsed charge control for the space inertial sensor used in LISA," *Phys. Rev. D, Part. Fields*, vol. 102, no. 4, Aug. 2020, Art. no. 042002.
- [43] J. Jiang, T. Li, C. Chang, C. Yang, and L. Liao, "Fault diagnosis method for lithium-ion batteries in electric vehicles based on isolated forest algorithm," *J. Energy Storage*, vol. 50, Jun. 2022, Art. no. 104177.
- [44] Q. Yao, D. Song, and X. Xu, "Robust finger-vein ROI localization based on the 3σ criterion dynamic threshold strategy," *Sensors*, vol. 20, no. 14, p. 3997, Jul. 2020.
- [45] Z. Wang et al., "Is mamba effective for time series forecasting?" 2024, *arXiv:2403.11144*.
- [46] M. Bassan et al., "Actuation crosstalk in free-falling systems: Torsion pendulum results for the engineering model of the LISA pathfinder gravitational reference sensor," *Astroparticle Phys.*, vol. 97, pp. 19–26, Jan. 2018.
- [47] M. Armano et al., "LISA pathfinder: The experiment and the route to LISA," *Classical Quantum Grav.*, vol. 26, no. 9, Apr. 2009, Art. no. 094001.



Cheng Bi received the bachelor's degree from Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2021, and the master's degree from Northwestern Polytechnical University, Xi'an, China, in 2024, where he is currently pursuing the Ph.D. degree with the School of Artificial Intelligence, Optics and Electronics (iOPEN).

His main research interests include pattern recognition and image processing.



Xiaokui Yue received the Ph.D. degree in flight vehicle design from Northwestern Polytechnical University, Xi'an, China, in 2002.

He is currently a Professor with the School of Astronautics, Northwestern Polytechnical University. He was named as the Distinguished Professor of Chang Jiang Scholar, in 2016. His current research interests include spacecraft dynamics and control, space maneuver and intelligent operation, and relative navigation.



Zhaohui Dang received the B.E. and Ph.D. degrees from the National University of Defense Technology, Changsha, China, in 2009 and 2015, respectively.

He is an Associate Professor and a Doctoral Supervisor at the School of Astronautics, Northwestern Polytechnical University, Xi'an, China. He has served as a Principal Investigator for 16 research projects. His main research interests include space orbital game technology, space artificial intelligence technology, and spacecraft formation flight technology.



Yibo Ding received the Ph.D. degree in aeronautical and astronautical science and technology from Harbin Institute of Technology, Harbin, China, in 2020.

He is currently an Associate Professor with the School of Astronautics, Northwestern Polytechnical University, Xi'an, China. His research interests include aircraft game confrontation and collaborative guidance technology.



Yonghe Zhang received the Ph.D. degree from the University of Chinese Academy of Sciences, Beijing, China, in 2016.

He is currently the Vice President of the Institute of Microsatellite Innovation, Chinese Academy of Sciences, Shanghai, China. His research interests include satellite overall design and simulation, and deep space navigation.