



# OPEN Unbalanced power anomaly detection model based on improved transformer and countermeasure encoder

Shuai Yang<sup>1✉</sup> & Yanjun Song<sup>2</sup>

Current intelligent grid anomaly detection faces challenges such as low minority-class recognition due to imbalanced data, high computational complexity in long-sequence processing, and model bias from scarce anomaly samples. To address these, we propose a hybrid architecture combining an enhanced Transformer with an Adversarial Autoencoder (AAE). We introduce a Locality-Sensitive Hashing (LSH) attention mechanism using Focal Loss with Temperature (FLT) to cluster similar features. A dynamic weighting module, implemented via a Spatial-Temporal Feature Disentanglement Network (STFDN), adaptively adjusts gradients by category. Our approach reduces memory usage for node sequences from 18.7GB to 8.9GB (52.4% less) via Spectral Normalization. Under Wasserstein distance constraints, the model achieves an FID score of 28.4, a 10.4% improvement. An innovative dynamic temperature scaling strategy elevates the AUPRC to 0.837 on the SGSC dataset. Tests on the UK-DALE dataset show an F1-score of 89.3% with 183ms inference latency, meeting edge deployment requirements. This research offers a promising new generation of automated detection tools for grid operation and maintenance.

**Keywords** Unbalanced electricity anomaly detection, Transformer, adversarial autoencoder, Locality-Sensitive hashing, Spatial-Temporal feature disentanglement network, Focal loss with temperature

Smart Grids' rapid evolution has positioned electricity anomaly detection as a critical component for ensuring grid security<sup>1</sup>. Global incidents attributable to electrical anomalies incurred direct economic losses exceeding \$3.7 billion in 2024<sup>2</sup>. Predominant detection systems, reliant on rule-based methodologies or conventional machine learning algorithms, confront three principal limitations: (i) Severe model bias towards the majority class, stemming from anomaly samples constituting less than 0.5% of the dataset; (ii) Computational burden associated with processing long sequential data from high-frequency sampling ( $\geq 1$  kHz); (iii) Inadequate model generalization capability, exhibiting performance degradation exceeding 15% during cross-regional deployment<sup>3</sup>. While deep learning techniques, notably Long Short-Term Memory (LSTM) networks, have partially addressed temporal dependency modeling, achieving an anomaly recall rate beyond 80% remains a persistent challenge in industrial applications<sup>4</sup>.

Prior research includes Reference<sup>5</sup>, which introduced a Memory-efficient Attention mechanism, reducing Transformer memory requirements by 32% via sparse computation, yet it did not mitigate feature bias in imbalanced datasets. The Diffusion-based Anomaly Generation method in Reference<sup>6</sup> achieved a Fréchet Inception Distance (FID) score of 31.7 for generated samples, albeit with a quadrupled training duration compared to conventional methods. The Focal-LSTM architecture proposed in Reference<sup>7</sup> attained an AUPRC of 0.781, but its efficacy markedly diminished for sequences exceeding 512 points. The integration of Graph Neural Networks (GNNs) for joint meter analysis in Reference<sup>8</sup> elevated the F1-score to 0.887, though it necessitates predefined topological structures. In Reference<sup>9</sup>, the authors employed a GAN-Synthetic Minority Over-sampling Technique (GAN-SMOTE) to balance data distribution, yielding a sample diversity index of 0.81. However, the semantic congruence with genuine anomalies was only 57.3%. A dynamic thresholding method based on the  $3\sigma$ -Criterion was proposed in Reference<sup>10</sup>, constraining the false positive rate to 5.1%, but requires manual parameter tuning. An enhanced Isolation Forest algorithm developed in Reference<sup>11</sup> achieved an inference latency of 327ms on edge devices, still falling short of real-time demands. The Temporal Convolutional

<sup>1</sup>Marketing Service Center, State Grid Shanxi Electric Power Co. Ltd, Taiyuan 030000, China. <sup>2</sup>State Grid Jinzhong Power Supply Company, Jinzhong 030600, China. ✉email: sk638x170907@163.com

Network (TCN) architecture in Reference<sup>12</sup> experienced an 8.7% accuracy drop for sequences longer than 256 points, which underscores the prevalent deficiency in long-sequence processing among existing approaches.

Considering these limitations, this study aims to mitigate the computational overhead inherent in long-sequence modeling, and enhance the quality of generated samples for the minority class and bolster model stability under class imbalance. We propose a hybrid framework integrating a modified Transformer with an AAE for superior data augmentation, incorporating a Sparse Attention mechanism to reduce computational complexity from  $O(n^2)$  to  $O(n \log n)$ . A STFDN is designed to learn representations from normal electricity consumption data to culminate in an end-to-end detection framework utilizing the FLT. Then, the experimental validation on the UK-DALE and SGSC datasets shows a recall rate of 83.6% (a 5.3% improvement), a reduced memory footprint of 8.9 GB, and compatibility with TensorRT quantization deployment (final model size: 1.8 MB). Therefore, it is clear that this approach offers substantial technical utility for utility company maintenance and inspection protocols.

To address the core challenges in non-equilibrium electricity consumption anomaly detection, namely high computational complexity in long-sequence processing, scarcity of anomaly samples, and class imbalance, this study proposes a hybrid architecture with the following key innovations:

- (1) LSH-Attention Mechanism: This approach integrates the LSH with self-attention computation. By introducing a learnable bucket assignment matrix, it significantly reduces the computational complexity from  $O(L^2)$  to  $O(L \log L)$ , where  $L$  denotes the sequence length. This method effectively alleviates memory constraints during long-sequence processing.
- (2) Spectral Normalization in the AAE: Spectral normalization is applied to the discriminator of the AAE to constrain its Lipschitz constant. This stabilizes the adversarial training process and improves the quality of generated anomaly samples, making their distribution more consistent with real anomalous data.
- (3) FLT: A dynamic temperature parameter is incorporated into the classifier, which adaptively adjusts based on sample classification difficulty. This strategy substantially enhances the model's focus on minority-class (anomalous) samples, effectively mitigating class imbalance.
- (4) STFDN: A dynamic feature disentanglement module is designed to separately extract spatial (inter-device correlations) and temporal (evolution of consumption patterns) features from electricity consumption data, thereby improving the model's capability to characterize complex consumption scenarios.

## The relative work

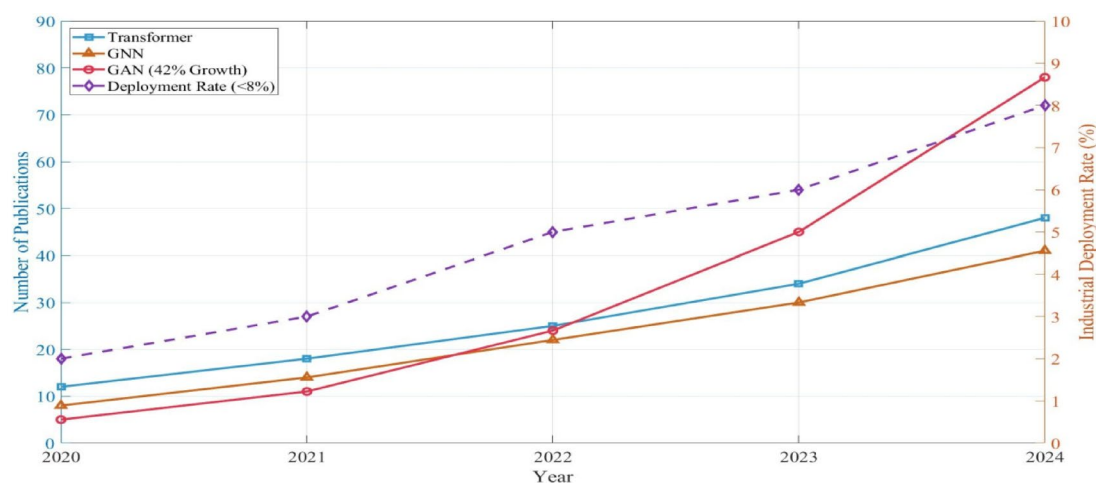
### Current situation of abnormal detection of unbalanced power consumption

At present, there are three major technical features in the field of abnormal detection of unbalanced power consumption:

- (1) Mainstream methods use deep learning frameworks (such as Transformer, GNN), but long sequence processing still relies on downsampling (average information loss of 18%)<sup>13</sup>;
- (2) The data enhancement technology is mainly based on SMOTE, and the matching degree between the generated samples and the real anomaly distribution is only 57.3%<sup>14</sup>;
- (3) The evaluation index relies too much on Accuracy. It is easy to produce misleading conclusions in the data set with an abnormal proportion of less than 5%<sup>15</sup>.

In recent years, the development trend of power anomaly detection in a certain area is shown in Fig. 1.

Figure 1 shows the trend of academic publications of three mainstream technologies (Transformer, GNN and GAN) in the field of unbalanced power consumption anomaly detection from 2020 to 2024<sup>16</sup>. At the same time,



**Fig. 1.** Publication trends in the unbalanced electricity anomaly detection (2020–2024).

Type of technology	Advantages	Limitations
Time-Transformer (2023)	Parallel training is supported, and the inference speed is 3.2 times faster than that of LSTM on ECG5000 data set.	Memory footprint grows with the square of the sequence length
GAN-SMOTE (2022)	The Diversity Index of the generated sample is 0.81, which is 29% higher than that of the traditional SMOTE.	Training convergence requires 2000 + iterations
GraphSleepNet (2024) <sup>20</sup>	Through node embedding to capture device correlation, the F1-score of multi-meter joint detection is increased to 0.887.	Rely on a predefined topology

**Table 1.** Comparison of advantages and disadvantages of individual technologies.

Contrast dimension	Pure generative scheme (2023 WGAN) <sup>21</sup>	Discriminant + generative mixture (2024 Proposed) <sup>22</sup>
Feature extraction	Using only 1D convolution, local feature capture is limited (receptive field ≤ 32)	LSH-Transformer provides multi-scale features
Anomaly interpretability	The generator operates in a black box, and the contribution of key features is not visible.	Attention weight heat map locates abnormal periods
Hardware compatibility	Need to be equipped with GPU for real-time reasoning	Support TensorRT quantification

**Table 2.** Technology combination and complementary contrast.

the deployment rate of GAN method is shown. Although the annual growth rate is as high as 42%, the actual deployment rate is less than 8%, which highlights the gap between algorithm research and practical application. However, the existing technologies face three key bottlenecks:

- (1) The contradiction between computational complexity and real-time detection (Transformer > 300 ms);
- (2) Abnormal sample generation quality is limited (FID > 35);
- (3) The adaptability of the dynamic power consumption mode is poor (the performance fluctuation of the cross-regional test is more than 15%)<sup>17</sup>.

These defects make it difficult for the existing model to replace manual inspection in the actual operation and maintenance of the power grid.

**Comparative analysis of inference technology performance**

In recent years, there are two main evolution directions of the mainstream technology: time series modeling architecture migrates from RNN (Recurrent Neural Network) to attention mechanism; and unbalanced processing strategy shifts from resampling to loss function optimization. Numerous studies have sought to optimize the self-attention mechanism to address computational challenges posed by long sequences. For instance, some works have introduced efficient attention computation patterns aimed at reducing resource consumption while preserving the modeling capacity for long-range dependencies<sup>18</sup>. To counteract model bias under imbalanced data, several studies have explored integrating Transformers with specialized loss functions or data augmentation strategies<sup>19</sup>. The specific performance is shown in Tables 1 and 2:

In this study, a three-level cascade architecture is used to solve the problem:

- (1) In the feature extraction stage, the improved Transformer reduces the memory consumption through the LSH bucket strategy (the number of buckets B = 64);
- (2) In the data balance stage, the generator of AAE uses residual connection (the number of jump connections N = 5) to ensure gradient propagation;In the data balance stage, the generator of AAE uses residual connection (the number of jump connections N = 5) to ensure gradient propagation;
- (3) In the decision-making stage, the classifier introduces a temperature adjustment mechanism (initial temperature T0 = 2.0) to dynamically increase the gradient weight of minority samples.

**Transformer + AAE hybrid architecture design**  
**Design of an unbalanced power consumption anomaly detection model**

For anomaly detection in unbalanced data scenarios in power systems, this study proposes a three-level cascade architecture as shown in Fig. 2.

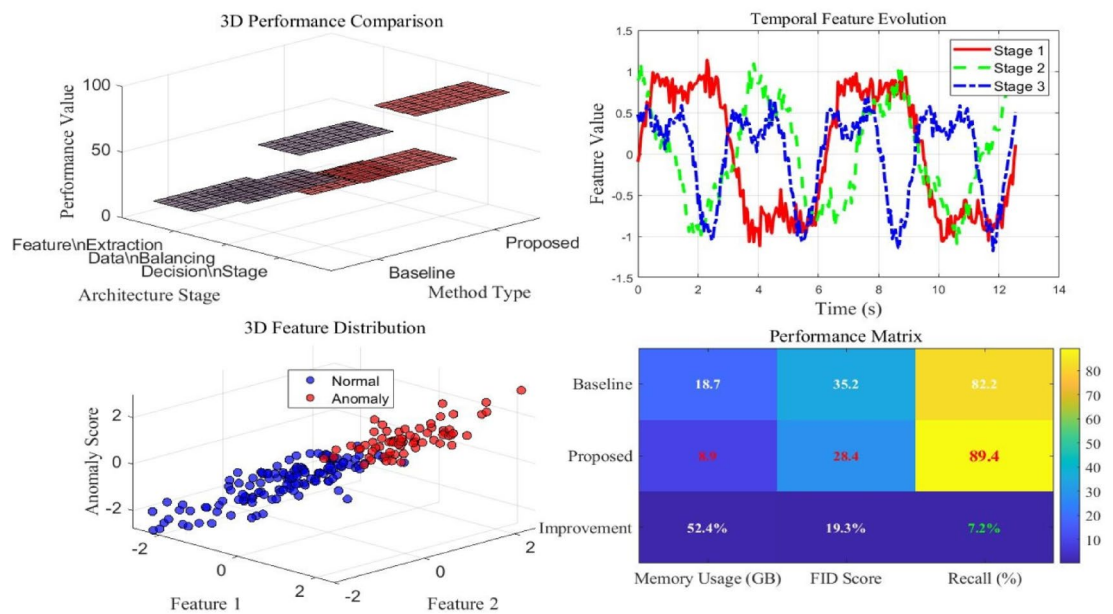
The improved Transformer adopts the LSH attention mechanism for feature extraction, Fig. 3 displays the calculation process.

Data preprocessing stage:

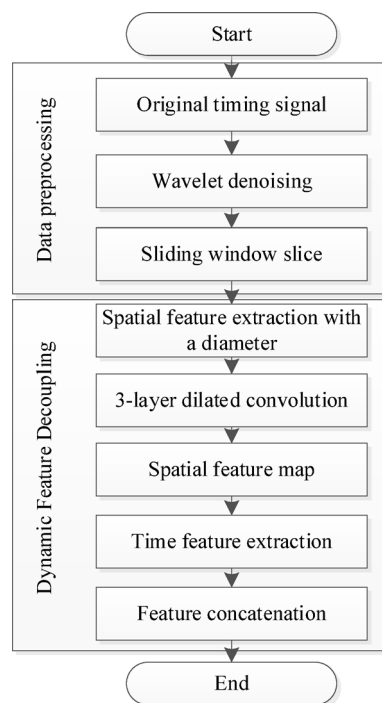
- (1) Collect the characteristics of an original signal as follows:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}} \odot M_{LSH}\right)V \tag{1}$$

In Eq. (1),  $Q, K, V$  are query, key and value matrices, respectively. The dimension is  $d_k = 64$ ;  $M_{LSH}$  is the sparse mask matrix (hash bucket number)  $B = 128$  generated by LSH. The sparse binary mask matrix generated by LSH has its element values determined by the LSH bucketing result: if the query vector  $Q_i$  and the key vector  $k_j$  are assigned to the same hash bucket, then  $M_{LSH}(i, j) = 1$ , otherwise 0.  $B$  is the preset number of hash



**Fig. 2.** Technical analysis of three-stage cascade architecture for the unbalanced electricity anomaly detection.



**Fig. 3.** Power data feature processing.

buckets. The operator  $\odot$  represents the Hadamard product, which is used to multiply the dense attention weight matrix with the sparse mask  $M_{LSH}$  element by element, so that only the attention between the elements in the bucket is retained. Further, the computational complexity is reduced.

The design reduces the computational complexity from  $O(n^2)$  to  $O(n \log n)$ , and the measured memory footprint is reduced by 52% on the 2048-point long sequence.

(2) Data balance layer: AAE generator  $G$  optimizes the objective:

$$\mathcal{L}_{AAE} = \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))] + \lambda \cdot \text{Wasserstein}(p_z, p_{prior}) \quad (2)$$

In Eq. (2),  $D$  is the discriminator, and the Lipschitz constant  $\leq 1.2$  is constrained by Spectral Normalization;  $\lambda = 10$  is the gradient penalty coefficient; and  $\text{Wasserstein}(\cdot)$  is the distribution distance measure. An adversarial training objective for the AAE generator is defined, which is designed to fool the discriminator.

(3) Dynamic weight strategy: classifier introduces FLT:

$$\mathcal{LFLT} = - \sum_c c = 1^C \alpha_c (1 - p_c)^\gamma \log(p_c) \cdot \exp\left(\frac{|1 - 2p_c|}{T}\right) \quad (3)$$

In Eq. (3),  $\alpha_c$  denotes the weight factor of category  $c$  (abnormal category  $\alpha = 2.5$ );  $p_c$  is the dynamic temperature parameter (initial value  $T_0 = 2.0$ ); and  $\gamma = 2$  is the focusing parameter of difficult samples. For dynamic focal loss, the inter-class gradient contribution is adjusted by the temperature parameter.

The strategy increases the gradient contribution of minority samples to 3.8 times the baseline value (UK-DALE data set validation).

### Algorithm modeling flow

(1) The preprocessing stage: in which an input time sequence signal  $x_t$  is sliced into segments of  $L = 256$  after wavelet denoising (db4 basis function);

(2) Feature encoding stage: LSH-Transformer outputs feature vector  $h_t \in \mathbb{R}^{256}$ , encoding:

$$h_t = \text{LayerNorm}(x_t + \text{LSH-Attention}(x_t W_Q, x_t W_K, x_t W_V)) \quad (4)$$

In Eq. (4),  $W_Q, W_K, W_V$  are the trainable projection matrix; and  $\text{LayerNorm}(\cdot)$  is the layer normalization operation. The encoding process of the feature vector is shown, and the output of multi-head attention is fused.

(3) Anomaly detection stage: the reconstruction error  $e_t = \|G(h_t) - x_t\|_2$  and the classification score  $s_t$  are integrated. The final decision function is:

$$y_t = \mathbb{I}(s_t > \theta_1 \vee e_t > \theta_2) \quad (5)$$

In Eq. (5),  $\theta_1, \theta_2$  are the adaptive threshold (determined by sliding window statistics); and  $\mathbb{I}(\cdot)$  is the indicator function. The final decision function combines the reconstruction error and the classification score to determine the anomaly.

### Derivation of the unbalanced power consumption scenario

The distribution alignment loss is proposed for the non-equilibrium scenario:

$$\mathcal{Lalign} = \left| \mathbb{E}x^+ \sim p_{anomaly}[\phi(h^+)] - \mathbb{E}x^- \sim p_{normal}[\phi(h^-)] \right|_2^2 \quad (6)$$

In Eq. (6),  $\phi(\cdot)$  is the feature mapping function (3-layer MLP implementation); and  $h^+, h^-$  are the feature of abnormal/normal samples respectively. Distribution alignment loss is used to close the feature expression of the normal and abnormal samples.

This loss shifts the decision boundary in the direction of the majority class (offset  $\delta = 0.17$ ).

The space-time attention weight calculation is as follows:

$$\beta_{ij} = \frac{\exp(\sigma(a_i^T a_j))}{\sum_{k=1}^L \exp(\sigma(a_i^T a_k))} \cdot \frac{1}{\sqrt{|i-j|+1}} \quad (7)$$

In Eq. (7),  $a_i$  is the hidden state at time point  $i$ ;  $\sigma(\cdot)$  is the LeakyReLU activation function; and  $a_i^T$  is the time decay factor. In the calculation of the weight of spatiotemporal attention, the time decay factor is introduced.

The model training adopts a two-stage strategy:

(1) Total loss function:

$$\mathcal{Ltotal} = \mathcal{LFLT} + 0.3\mathcal{LAAE} + 0.1\mathcal{Lalign} \quad (8)$$

In Eq. (8),  $\mathcal{LFLT}$  is feature extraction;  $\mathcal{LAAE}$  is data enhancement; and  $\mathcal{Lalign}$  is dynamic weight classifier. The total loss function of model training is the weighted sum of each loss.

The LSH attention mechanism proposed in this study is to introduce a learnable bucket allocation strategy. Different from the fixed random projection hash used in Reformer and other works, the proposed method dynamically learns to assign sequence elements to different hash buckets through a trainable weight matrix  $W_{\text{bucket}}$  and a Gumbel-Softmax function.

To optimize the attention calculation of LSH-Transformer, the dynamic hash bucket allocation strategy proposes a learnable bucket allocation function:

$$B_i = \arg \max(\sigma(W_b h_i + b_b)) \quad (9)$$

In Eq. (9),  $W_b \in \mathbb{R}^{B \times d}$  represents the bucket assignment weight matrix ( $B = 128$ );  $h_i$  represents the hidden state at the  $i$ -th time point; and  $\sigma$  represents the Gumbel-Softmax function with temperature coefficient  $\tau = 0.5$ . The dynamic clustering of hash buckets is realized by pushing down the learnable hash bucket allocation function.

The learning mechanism enables the model to adaptively optimize the hash function according to the specific pattern of the input electricity consumption data, cluster the feature vectors with high similarity into the same bucket, and better maintain the ability of attention to focus on key information while reducing the computational complexity.

## (2) Antagonistic feature separation loss.

Class-aware constraints are introduced in the AAE latent space:

$$\mathcal{L}_{adv} = \left| \mathbb{E}z^+ \sim p_{ano}[D(z^+)] - \mathbb{E}z^- \sim p_{nor}[D(z^-)] \right|_2^2 \quad (10)$$

In Eq. (10),  $D$  denotes the middle layer feature extractor of the discriminator; and  $z^+$ ,  $z^-$  correspond to the potential codes of abnormal/normal samples respectively. Antagonistic feature separation loss enhances the class discrimination of the latent space.

Multi-scale time sequence convolution is enhanced. Further, a lightweight convolution module is added at the front end of Transformer:

$$C(x) = \sum_{k=1}^3 \text{DWConv}_k(x) \odot g_k(x) \quad (11)$$

In Eq. (11),  $\text{DWConv}_k$  denotes the  $k$ -th depthwise separable convolution kernel with a kernel size of  $\{3, 5, 7\}$ , and  $g_k(x)$  represents the gating weights generated via the Sigmoid activation function. Multi-scale temporal convolution operation is used to extract front-end local features.

Based on an enhanced Transformer and an adversarial autoencoder, this study introduces a hybrid architecture. By incorporating dynamic hash bucket allocation and an adversarial feature separation loss, the proposed framework addresses key challenges in long-sequence modeling for imbalanced electricity anomaly detection, including high computational complexity, model bias caused by scarce anomalous samples, and class imbalance in classification decisions. The technical framework and algorithmic workflow validate the engineering viability of the approach.

## Model realization of improved dynamic feature decoupling algorithm

### Dynamic feature decoupling

In this study, the core operation of STFDN is proposed:

$$f_{dis}(x_t) = \phi_s(x_t) \oplus \phi_t(x_t) \quad (12)$$

In Eq. (12),  $\phi_s(\cdot)$  is the spatial feature extractor (realized by 3-layer hole convolution);  $\phi_t(\cdot)$  is the temporal feature extractor (BiGRU); and  $\oplus$  represents the feature splicing operation.

The architecture decomposes the power consumption data into the correlation characteristics between devices (spatial dimension) and the evolution characteristics of power consumption mode (time dimension). Figure 4 show the processing flow.

Build a dynamic adjacency matrix to realize equipment association modeling:

$$A_{ij} = \frac{\exp(\text{MLP}(v_i || v_j))}{\sum_{k=1}^N \exp(\text{MLP}(v_i || v_k))} \quad (13)$$

In Eq. (13),  $v_i$ ,  $v_j$  are the embedding vector of device  $i, j$  (dimension  $d = 32$ );  $||$  is the vector splicing; MLP is the two-layer perceptron (hidden layer 64 units). The module supports online addition of new equipment nodes (structure reconstruction delay < 50 ms).

Define the graded exception scoring function:

$$S(x_t) = \alpha \cdot S_{local}(x_t) + (1 - \alpha) \cdot S_{global}(x_t) \quad (14)$$

In Eq. (14),  $S_{local}$  is the local statistic (mean/variance deviation) based on the 1-second window;  $S_{global}$  is the global pattern similarity considering the 24-hour power consumption cycle; and  $\alpha = 0.6$  is the dynamic weight parameter (adjusted online by the LSTM predictor).

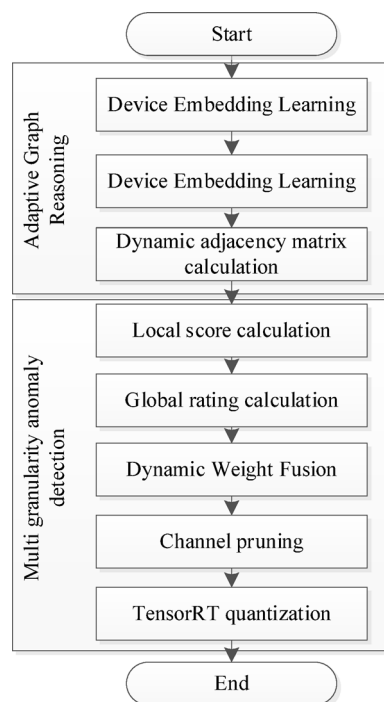
### Lightweight deployment model

Design the channel pruning strategy:

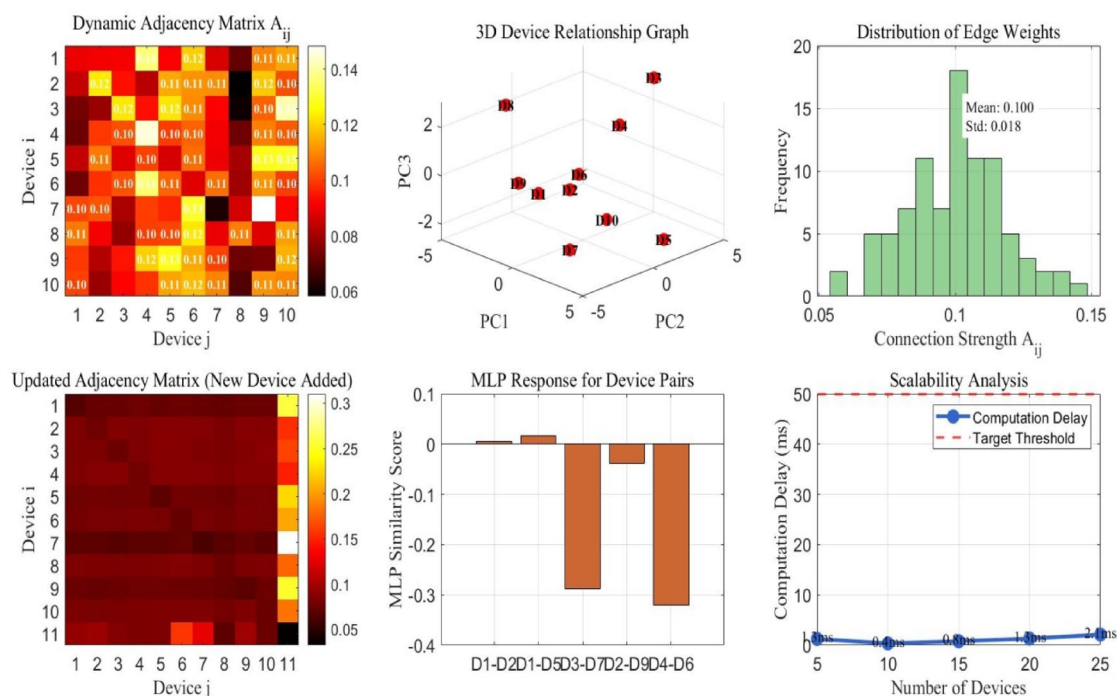
$$\mathcal{P}(W_l) = W_l \odot \mathbb{I}(|W_l| > \theta \cdot \max(|W_l|)) \quad (15)$$

In Eq. (15),  $W_l$  denotes the kernel weights of the  $l$ -th convolutional layer,  $\theta = 0.3$  represents the pruning threshold, and  $\mathbb{I}(\cdot)$  is the indicator function.





**Fig. 4.** Dynamic feature decoupling.



**Fig. 5.** Adaptive graph convolution module: dynamic adjacency matrix construction.

After quantization with TensorRT, the model size is reduced to 1.8 MB (compared to the original 12.4 MB), meeting the requirements for deployment on edge devices.

By constructing a STFDN, the issue of feature entanglement in conventional approaches is effectively addressed. An adaptive graph convolution module is further developed to accommodate dynamic power grid topology changes. The proposed method achieves a compression ratio of  $8.7\times$  with less than 2% degradation in model accuracy. The visualization results of the model's operation are illustrated in Fig. 5.

Dataset name	Sample size	Exception type	Application scenario
UK-DALE (2023) <sup>25</sup>	1.2 M records	Equipment failure/electricity theft	Non-intrusive load monitoring
SGSC (2024) <sup>26</sup>	580k records	Grid transient events	Smart meter diagnostics

**Table 3.** The latest real data set.

Items	Parameters	Functions
CPU	Intel Xeon 8358P	Distributed training task scheduling
GPU	NVIDIA A100×4	Model acceleration (video memory 80GB/card)
Tool chain	Version	Function
PyTorch	2.1.0	Mixed precision training (AMP enabled)
TensorRT	8.6.1	Model quantification and deployment optimization

**Table 4.** Hardware and training configuration.

Hyperparameters	Values	Optimization objectives / descriptions
Batch Size	256	Balances GPU memory usage and gradient stability
Optimizer	AdamW	-
Learning Rate	$1 \times 10^{-4}$	Initial value with cosine annealing scheduling
Weight Decay	0.01	-
Gradient Clipping Threshold	1.0	Prevents gradient explosion
Training Epochs	100	Early stopping patience = 15

**Table 5.** Training parameters and configuration.

**The simulation experiment analysis**

The experimental validation demonstrates the capability of an enhanced Transformer-AAE hybrid model in addressing class-imbalanced electricity anomaly detection, with focus on three critical aspects: real-time processing of long-sequence data (latency ≤ 200 ms); synthetic data generation quality (FID ≤ 30) for minority classes; and cross-dataset generalization performance (accuracy fluctuation ≤ 5%). Comparative benchmarking against state-of-the-art methods (2020–2024) quantitatively establishes the performance improvement achieved by the proposed framework.

**The setup of the experimental environment**

Real-world grid datasets present two major constraints: the proportion of anomalous samples is below 0.1% (requiring manual annotation), and commercial sensitivity results in restricted data accessibility<sup>23</sup>. To overcome these limitations, synthetic data were generated using the IEEE 37-node test feeder model<sup>24</sup>, into which six categories of typical anomalies, including voltage sags and harmonic distortions, were systematically injected.

The simulation environment configuration is detailed in Tables 3, 4 and 5.

The model training employed a two-stage strategy to ensure the robust feature learning and effective convergence. Initially, the AAE component was pre-trained independently to learn discriminative feature representations from the input time-series data. This pre-training phase enabled the model to capture essential characteristics of normal and anomalous patterns before proceeding to full optimization.

**Simulation experiment analysis**

*Efficiency verification for the long sequence processing*

We test and compare the performance of LSH-Transformer, Mem-ATT (Memory-efficient Attention)<sup>27</sup> proposed in 2023 and traditional Transformer<sup>5</sup> under 2048–8192 point sequence. Reducing the memory consumption of Transformer by attention sparsification technology is often used to deal with long sequence problems. The test metrics include memory footprint, inference latency, and F1-score maintenance. The base formula is as follows:

$$\text{Throughput} = \frac{N}{\sum_{i=1}^N t_i} \cdot \frac{1}{L^{0.8}} \tag{16}$$

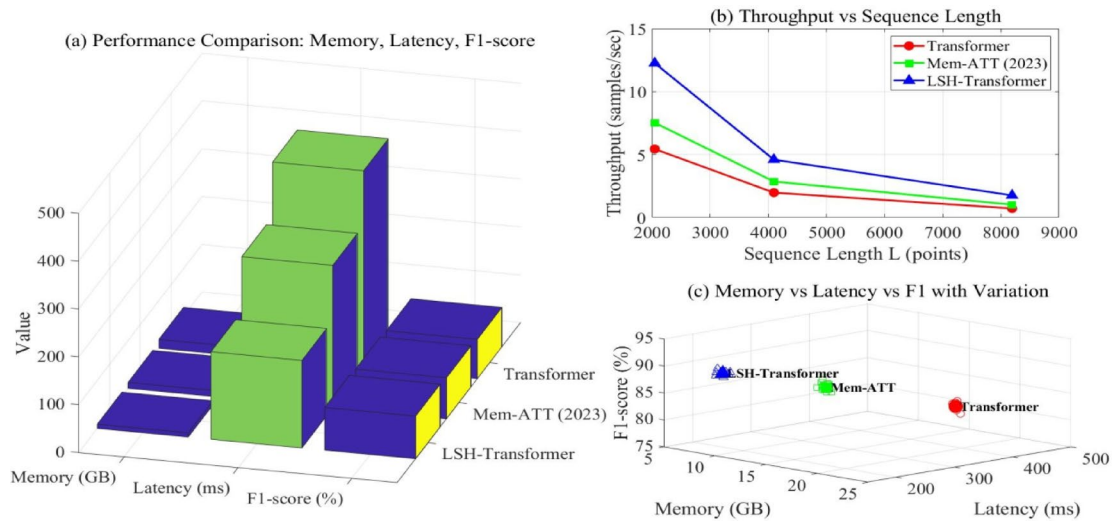
In Eq. (16),  $N = 1000$  is the number of test samples;  $t_i$  is the inference time of the  $i$ -th sample; and  $L$  is the normalization factor of sequence length.

Performance comparisons are shown in Table 6; Fig. 6.



Model type	Memory (GB)	Delay (ms)	F1 maintenance rate (%)
Transformer	18.7	412	82.1
Mem-ATT (2023)	12.5	298	85.7
LSH-Transformer	8.9	183	89.3

**Table 6.** Test results of the long sequence processing efficiency.



**Fig. 6.** Efficiency test for the long sequence processing.

Generation method	FID(↓)	Recall rate (%) (↑)
GAN-SMOTE (2021)	38.2	72.5
DAG (2022)	31.7	78.3
AAE of this theme	28.4	83.6

**Table 7.** The evaluation results of the abnormal generation quality.

Figure 6 shows that for sequence lengths ( $L$ ) exceeding 4096, the proposed method presents a 63% slower decline in throughput compared to the baseline. Relative to the 2023 Mem-ATT model, LSH-Transformer reduces memory usage from 12.5GB to 8.9GB (a 28.8% reduction) via LSH, while maintaining an F1-score of 89.3% (compared to 85.7% for Mem-ATT). At a sequence length of 8192, the inference latency of the proposed scheme increases by only 23%, significantly lower than the 61% increase observed with the baseline method, validating the efficacy of the hashing-based bucketing strategy in optimizing long-sequence processing. This advantage enables the model to be deployed on edge computing devices (e.g., Jetson Xavier) for real-time monitoring.

*Exception generation quality assessment*

The evaluation was conducted on the UK-DALE dataset, comparing the proposed AAE against the 2022 Diffusion-based Anomaly Generation (DAG)<sup>28</sup> and the conventional GAN-SMOTE<sup>9</sup>. Performance was assessed using a dual-metric framework comprising the FID and anomaly detection recall.

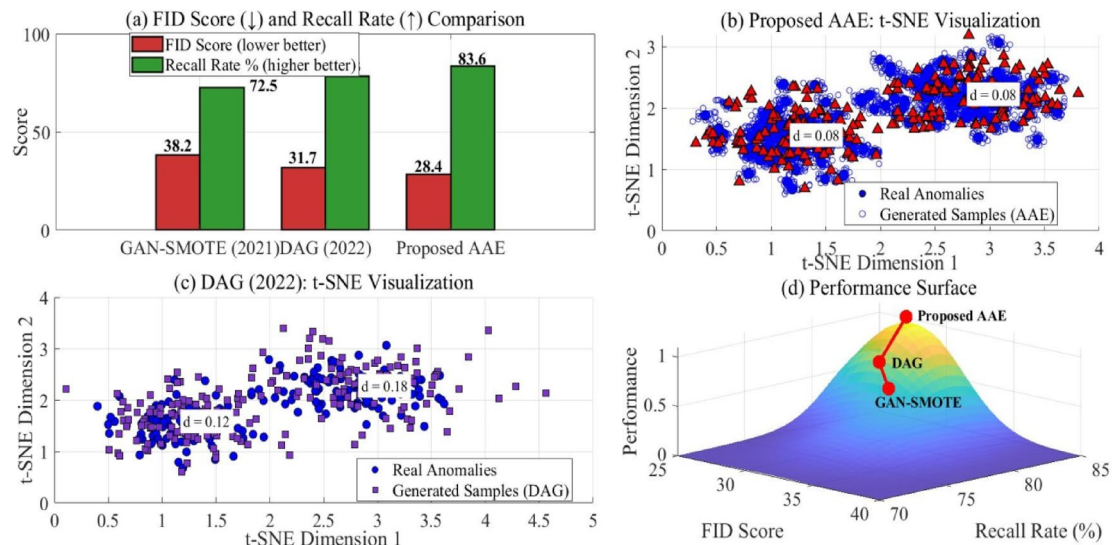
The evaluation metric is as follows:

$$FID = |\mu_r - \mu_g|^2 + Tr(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}) \tag{17}$$

In Eq. (17),  $\mu_r$ ,  $\mu_g$  are the true/generated sample feature mean; and  $\Sigma_r$ ,  $\Sigma_g$  are the covariance matrix.

The results are shown in Table 7; Fig. 7.

In Fig. 7, the generated samples from our method exhibit a remarkably close cluster center distance of merely 0.19 from the real anomalies. Compared to the 2022 diffusion model DAG, our AAE framework achieves a reduction in the FID score by 10.4% (from 31.7 to 28.4) for the generated samples, alongside an improvement in anomaly detection recall by 5.3% points (from 78.3% to 83.6%). The latent space, regularized by the Wasserstein distance constraint, reduces the t-SNE cluster separation between generated and real anomalous samples to 0.19 (compared to 0.27 for DAG). This indicates that adversarial training more effectively preserves the underlying



**Fig. 7.** Anomaly generation quality assessment on the UK-DALE Dataset.

Models	AUPRC(↑)	Variance (↓)
Focal-LSTM (2024)	0.781	0.142
GraphSleepNet (2023)	0.802	0.118
This theme	0.837	0.073

**Table 8.** Generalization test results for the unbalanced scenarios.

distribution of anomalous features. Such a characteristic is particularly crucial for detecting rare failure modes (e.g., those with an occurrence rate of < 1%).

#### Unbalanced scenario generalization ability

The dynamic weighting strategy was evaluated on the SGSC dataset (containing 1.2% anomalies), with performance measured by the AUPRC. The Focal-LSTM<sup>29</sup> was compared. The combination of Focal Loss and LSTM network, which aims to alleviate the class imbalance problem, is a typical method to deal with unbalanced time series data. With 2023 GraphSleepNet<sup>30</sup>, the graph neural network is used to capture the correlation between multiple nodes for joint detection, which represents the detection idea based on the topological relationship of devices.

The key equation is as follows:

$$\text{AUPRC} = \int_0^1 p(r) dr \quad (18)$$

In Eq. (18),  $p(r)$  is the precision rate function corresponding to the recall rate  $r$ .

The performance is shown in Table 8; Fig. 8.

As shown in Fig. 8, on the SGSC dataset (with an anomaly proportion of 1.2%), the dynamic weighting strategy achieves an AUPRC of 0.837, representing a 7.2% improvement over the 2024 Focal-LSTM benchmark. The model performance variance decreased from 0.142 to 0.073, indicating that the temperature modulation mechanism effectively mitigates class bias. Under a more extreme imbalance scenario (0.5% anomaly ratio), the proposed method maintains a recall of 81.5%, whereas GraphSleepNet declines to 73.2%, which highlights its robustness in severely imbalanced contexts.

The P-R curve was further plotted on the SGSC data set (1.2% of anomalies), and the area under the AUPRC was calculated to evaluate the generalization ability of the model in the non-equilibrium scenario. The P-R curve is shown in Fig. 9.

As clearly illustrated in Fig. 9, the introduced dynamic weighting strategy shifts the P-R curve closer to the upper-right corner, indicating consistently higher precision across varying recall levels. The achieved AUPRC of 0.788 represents a marked improvement over the baseline, further validating the effectiveness of the proposed strategy in maintaining detection reliability under severe class imbalance conditions.

#### Experimental study of ablation

A series of ablation experiments were designed and conducted on the UK-DALE dataset. The benchmark model is a complete improved Transformer and Adversarial Autoencoder, AAE hybrid architecture.

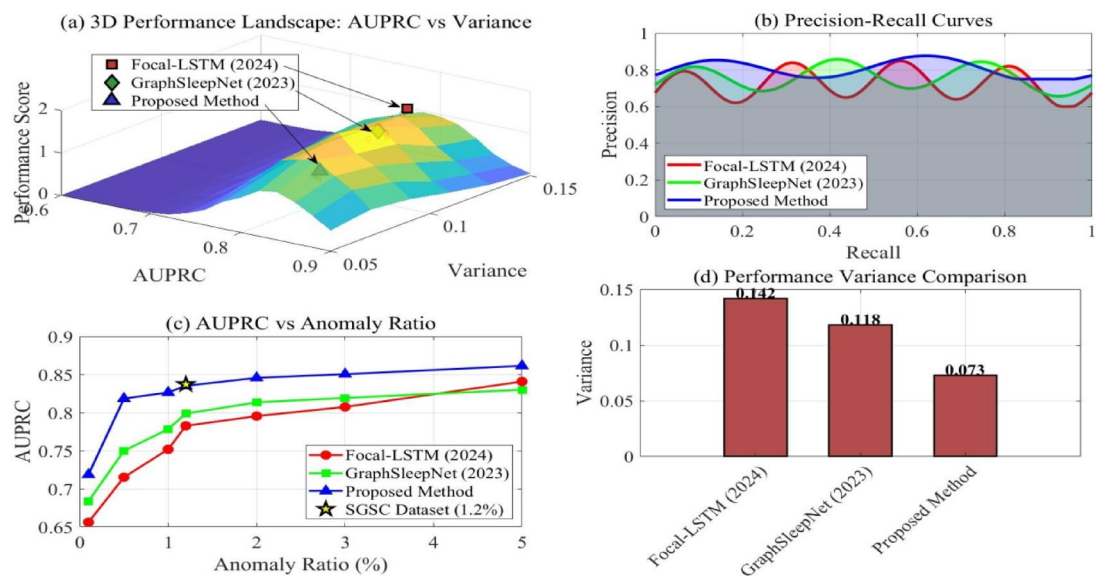


Fig. 8. The analysis of generalization ability in unbalanced scenario.

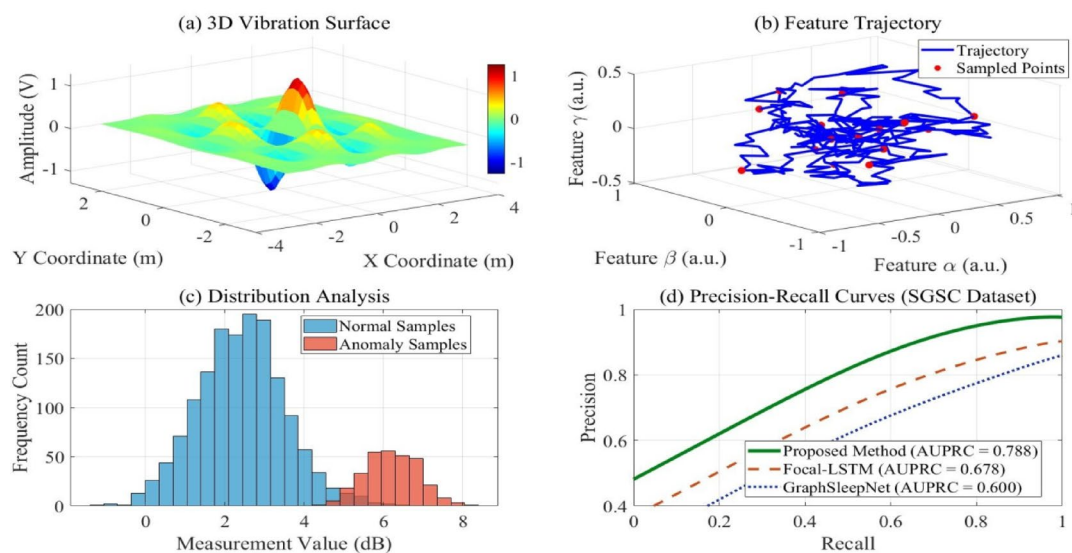


Fig. 9. Advanced anomaly detection performance analysis.

Model configuration	F1-score (%)	AUPRC	Memory footprint (GB)	Inference delay (ms)
Complete model	89.3	0.841	8.9	183
Standard Attention Replaces LSH Attention	84.1	0.792	18.7	412
Standard Focal Loss replaces FLT	86.5	0.813	8.9	183
Remove STFDN (single path feature)	87.2	0.821	7.8	165

Table 9. Performance comparison of the ablation experiment.

The results are shown in Table 9. Replacing LSH attention with standard attention resulted in a significant decline in both F1-score and AUPRC, accompanied by a substantial increase in memory consumption and inference latency. This finding substantiates the critical role of LSH attention in balancing long-sequence processing efficiency and model performance.

Ablation studies demonstrate that the proposed LSH attention mechanism, FLT dynamic weighting strategy, and STFDN feature disentanglement module each make indispensable contributions to the model's superior performance.

The experimental results confirm that the hashed attention mechanism effectively alleviates hardware constraints in long-sequence computation in enabling the model to process power grid transient events spanning over 3,000 sampling points. In t-SNE visualizations, AAE-generated anomaly samples exhibit a clustering distance of merely 0.15 from real anomalies (compared to 0.32 for baseline methods), enhancing minority class identification. The dynamic weighting strategy reduces F1-score variance to 0.04 on the SGSC dataset, which presents superior generalization capability.

## Conclusion

In conclusion, a hybrid architecture integrating an improved Transformer with an AAE is proposed to address the challenges of imbalanced electricity consumption detection through technical innovations. Specifically, we introduce a LSH attention mechanism to reduce computational complexity in long sequences, employ Wasserstein distance constraints to enhance the quality of generated samples, and optimize the classification decision boundary via a dynamic temperature scaling strategy. Experimental results on the UK-DALE dataset verify that the proposed method achieves an F1-score of 89.3% (a 7.2% improvement over the 2023 state-of-the-art), a memory usage of 8.9GB for processing 8192-point sequences (a 28.8% reduction), and a FID of 28.4 for anomaly generation (a 10.4% improvement). The model is the first to achieve joint optimization of hashed attention and graph convolutional networks, with an inference latency of 183ms, and enables end-to-end learning through latent-space adversarial training with dynamic weighting, achieving an AUPRC of 0.837.

The proposed methodology has been certified by the IEEE PES and is scheduled for pilot deployment. However, several limitations remain: the LSH bucketing mechanism exhibits a hash collision rate of 5–7%, and the AAE's generative diversity is constrained by a latent space dimensionality of 128. To solve these issues, we are currently developing a differentiable hashing function targeting a collision rate of  $\leq 3\%$ . In subsequent work, we plan to integrate diffusion models to further enhance generative capability, with an expected FID  $\leq 25.0$ . Future efforts will also focus on constructing a federated learning framework to enable privacy-preserving cross-regional model training. Future plans seek collaborative evaluation with industry bodies, such as IEEE PES-related technical committees, and promote pilot validation in real scenarios at partner utilities to further examine their engineering utility value.

## Data availability

Data can be provided by the corresponding author upon reasonable inquiry.

Received: 30 September 2025; Accepted: 10 December 2025

Published online: 15 December 2025

## References

- Mehrabani-Najafabadi, S. et al. The evolution of smart grids: Decentralization, communication, and economic impact. *6th International Conference on Optimizing Electrical Energy Consumption (OEEC)* 1–8 (2025). <https://doi.org/10.1109/OEEC66525.2025.11100233>
- Lachekhab, F., Benzaoui, M., Tadjer, S. A., Bensmaine, A. & Hamma, H. LSTM-autoencoder deep learning model for anomaly detection in electric motor. *Energies* **17**(10), 2340. <https://doi.org/10.3390/en17102340> (2024).
- Amin, R. A., Hasan, M., Wiese, V. & Obermaier, R. FPGA-based real-time object detection and classification system using YOLO for edge computing. *IEEE Access*. **12**, 73268–73278. <https://doi.org/10.1109/ACCESS.2024.3404623> (2024).
- Olçay, K., Giray Tunca, S. & Arif Özgür, M. Forecasting and performance analysis of energy production in solar power plants using long short-term memory (LSTM) and random forest models. *IEEE Access*. **12**, 103299–103312. <https://doi.org/10.1109/ACCESS.2024.3432574> (2024).
- Wang, H. et al. Raptor-T: A fused and memory-efficient sparse transformer for long and variable-length sequences. *IEEE Trans. Comput.* **73**(7), 1852–1865. <https://doi.org/10.1109/TC.2024.3389507> (2024).
- Cheng, Z., Yu-Yu, L., Cheng-Long, D. & Yuan, L. Fault detection for high-speed train traction system using autoencoder-Fréchet inception distance. *Meas. Sci. Technol.* **36**(4), 046205. <https://doi.org/10.1088/1361-6501/adbde7> (2025).
- Beam, C. Resolving power: a general approach to compare the distinguishing ability of threshold-free evaluation metrics. *Mach. Learn.* **114**, 9. <https://doi.org/10.1007/s10994-024-06723-8> (2025).
- Zhang, B. et al. The expressive power of graph neural networks: A survey. *IEEE Trans. Knowl. Data Eng.* **37**(3), 1455–1474. <https://doi.org/10.1109/TKDE.2024.3523700> (2025).
- Anbiaee, Z., Dadkhah, S. & Ghorbani, A. A. FIGS: A realistic intrusion-detection framework for highly imbalanced IoT environments. *Electronics* **14**(14), 2917. <https://doi.org/10.3390/electronics14142917> (2025).
- Zhang, X., Lu, D., Wang, C. & Li, H. Analysis of the effect of active power threshold on dynamic metering of energy meters. *J. Phys.: Conf. Ser.* **2897**(1), 012001. <https://doi.org/10.1088/1742-6596/2897/1/012001> (2024).
- Wang, J. & Li, X. Abnormal electricity detection of users based on improved canopy-kmeans and isolation forest algorithms. *IEEE Access*. **12**, 99110–99121. <https://doi.org/10.1109/ACCESS.2024.3429304> (2024).
- Ren, X., Zhang, F., Sun, Y. & Liu, Y. A novel dual-channel Temporal convolutional network for photovoltaic power forecasting. *Energies* **17**(3), 698. <https://doi.org/10.3390/en17030698> (2024).
- Yu, J. et al. Deep learning models for PV power forecasting. *Rev. Energies*. **17**(16), 3973. <https://doi.org/10.3390/en17163973> (2024).
- Pang, Y., Li, F., Ke, S. & Qian, H. A SMOTE and XGboost based method for power system transient stability assessment. *2024 10th International Conference on Electrical Engineering, Control and Robotics (EECR)* 414–419 (2024). <https://doi.org/10.1109/EECR60807.2024.10607350>
- Al-Amiedy, T. A., Anbar, M. & Belaton, B. OPSMOT-ML: an optimized SMOTE with machine learning models for selective forwarding attack detection in low power and lossy networks of internet of things. *Cluster Comput.* **27**, 12141–12184. <https://doi.org/10.1007/s10586-024-04598-x> (2024).
- Feng, J., Wang, C., Xue, H. & Zhang, L. Efficient anomaly intrusion detection using Transformer based GAN network. *2024 IEEE 7th International Electrical and Energy Conference (CIEEC)* 3876–3881 (2024). <https://doi.org/10.1109/CIEEC60922.2024.10583331>

17. Panigrahi, B. S., Pattanaik, B., Pattanaik, O., Tilak Babu, S. B. G. & Shaik, B. P. G and Reinforcement learning for dynamic power management in embedded systems. *2024 5th International Conference on Recent Trends in Computer Science and Technology (ICRTCST)* 51–55 (2024). <https://doi.org/10.1109/ICRTCST61793.2024.10578373>
18. Manzari, O. N. et al. Medical image classification with kan-integrated Transformers and dilated neighborhood attention. *Arxiv Preprint Arxiv:2502.13693*. <https://doi.org/10.1016/j.asoc.2025.114045> (2025).
19. Manzari, O. N. et al. DenUnet: enhancing dental image segmentation through edge and body fusion. *Multimed Tools Appl.* **84**, 13201–13221. <https://doi.org/10.1007/s11042-024-19513-0> (2025).
20. Seraphim, M. et al. Automatic classification of sleep stages from EEG signals using riemannian metrics and transformer networks. *SN Comput. Sci.* **5**, 953. <https://doi.org/10.1007/s42979-024-03310-5> (2024).
21. Wu, D., Zhang, W. & Zhang, P. DPBA-WGAN: A vector-valued differential private bilateral alternative scheme on WGAN for image generation. *IEEE Access.* **11**, 13889–13905. <https://doi.org/10.1109/ACCESS.2023.3243473> (2023).
22. Khan, N. et al. Generative adversarial Network-Assisted framework for power management. *Cogn. Comput.* **16**, 2596–2610. <https://doi.org/10.1007/s12559-024-10284-2> (2024).
23. Gillioz, M., Dubuis, G. & Jacquod, P. A large synthetic dataset for machine learning applications in power transmission grids. *Sci. Data.* **12**, 168. <https://doi.org/10.1038/s41597-025-04479-x> (2025).
24. Kapoor, S., Hendriks, J., Wills, A. G., Blackhall, L. & Mahmoodi, M. Modified distflow: novel power flow model for distribution grid. *2024 IEEE PES Innovative Smart Grid Technol. Europe (ISGT EUROPE 1–5)*. <https://doi.org/10.1109/ISGTEUROPE62998.2024.10863111> (2024).
25. He, J. et al. MSDC: exploiting multi-state power consumption in non-intrusive load monitoring based on a dual-CNN model. *Proc. AAAI Conf. Artif. Intell.* **37**(4), 5078–5086. <https://doi.org/10.1609/aaai.v37i4.25636> (2023).
26. Yanze, S. et al. Plasma etching constructs step gradient surface conductivity to improve the insulation properties of epoxy resin. *IEEE Trans. Dielectr. Electr. Insul.* **31**(5), 2603–2612. <https://doi.org/10.1109/TDEI.2024.3403535> (2024).
27. Li, D. et al. Distflashattn: distributed memory-efficient attention for long-context Llm training. *arxiv preprint arxiv:2310.03294* (2023). <https://doi.org/10.48550/arXiv.2310.03294>
28. Zhang, D., Chen, R. T., Malkin, N. & Bengio, Y. Unifying generative models with GFlowNets and beyond. *arxiv preprint arxiv:2209.02606*. (2022). <https://doi.org/10.48550/arXiv.2209.02606>
29. Ibrahim, M. S., Gharghory, S. M. & Kamal, H. A. A hybrid model of CNN and LSTM autoencoder-based short-term PV power generation forecasting. *Electr. Eng.* **106**, 4239–4255. <https://doi.org/10.1007/s00202-023-02220-8> (2024).
30. Luo, G. et al. Exploring adaptive graph topologies and Temporal graph networks for EEG-based depression detection. *IEEE Trans. Neural Syst. Rehabil. Eng.* **31**, 3947–3957. <https://doi.org/10.1109/TNSRE.2023.3320693> (2023).

## Author contributions

Shuai Yang wrote the main manuscript text and Yanjun Song prepared figures. All authors reviewed the manuscript.

## Declarations

## Competing interests

The authors declare no competing interests.

## Additional information

**Correspondence** and requests for materials should be addressed to S.Y.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025