

# Knowledge Aggregation Transformer Network for Multivariate Time Series Classification

Zhiwen Xiao, *Member, IEEE*, Huanlai Xing, *Member, IEEE*, Rong Qu, *Senior Member, IEEE*, Hui Li, Huagang Tong, Shouxi Luo, *Member, IEEE*, Jing Song, Li Feng, and Qian Wan

**Abstract**—Over the years, various sophisticated deep learning algorithms have surfaced for multivariate time series classification (MTSC), notably the dual-network-based model. This model comprises two parallel networks tailored to time series data: one for local feature extraction and the other for global relation extraction. However, effectively integrating these dual networks poses a significant challenge. To address this, we propose a knowledge aggregation transformer network (KATN) for MTSC. KATN, composed of four aggregation transformer blocks, extracts abundant regularizations and connections hidden within the data. Each block incorporates a modified residual network (MResNet) for local feature extraction and a multi-head attention network for global relation extraction. Initially, the block merges MResNet’s output feature with that of the multi-head attention network through an additive operation. Subsequently, it aligns features with a fully connected (i.e., dense) layer and activates neural units using the Gaussian error linear unit function. This strategic feature aggregation allows for capturing long-range dependencies among multiple variables in multivariate time series data. Experimental results demonstrate that KATN significantly outperforms 6 state-of-the-art transformer variants, achieving a ‘win’/‘tie’/‘lose’ record of 9/6/15 and securing the lowest AVG\_rank score. Furthermore, when evaluated against 18 existing MTSC algorithms across 13 UEA datasets, KATN consistently delivers superior performance, attaining the lowest AVG\_rank score among all compared methods.

**Index Terms**—Data Mining, Deep Learning, Feature Aggregation, Multivariate Time Series Classification, Transformer

## I. INTRODUCTION

MULTIVARIATE time series data has been widely applied in a variety of domains, e.g., heart failure risk analysis [1], industrial activity recognition [2], rumor detection [3], and anomaly detection [4], [5]. Different from image,

This work was partially supported by the Natural Science Foundation of Hebei Province (No. F2022105027), the Doctoral Innovation Fund Program of Southwest Jiaotong University (No. CX-2025ZD09), Southwest Jiaotong University, China, and the Fundamental Research Funds for the Central Universities, P. R. China (Corresponding Author: Huanlai Xing).

Z. Xiao, H. Xing, S. Luo, L. Feng, and Q. Wan are with the School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu 610031, China, and also with the Tangshan Institute of Southwest Jiaotong University, Tangshan 063000, China (Emails: xiao1994zw@163.com; hxx@home.swjtu.edu.cn; sxluo@swjtu.edu.cn; fengli@swjtu.edu.cn; qianwan@my.swjtu.edu.cn).

R. Qu is with the School of Computer Science, University of Nottingham, Nottingham NG7 2RD 455356, UK (Email: rong.qu@nottingham.ac.uk).

H. Li is with the School of Mathematics and Statistics, Xi’an Jiaotong University, Xi’an 710049, China (Email: lihui10@mail.xjtu.edu.cn).

H. Tong is with the College of Economic and Management, Nanjing Tech University, Puzhu Road(S), Nanjing 211816, China (Email: huagang-tong@gmail.com).

J. Song is with the Informatization Research Institute, Southwest Jiaotong University, Chengdu 610031, China (Email: jesen811206@126.com).

text, and video data, multivariate time series consist of sequential data points that are organized chronologically and correspond to multiple time-dependent variables, incorporating both local and global patterns. An algorithm designed for multivariate time series classification (MTSC) aims to extract distinct local and global features from each univariate time series (UTS) while simultaneously identifying interconnections among these UTS sequences [6].

In recent times, deep learning (DL) models have garnered significant interest within the MTSC community. These algorithms effectively model the internal data representation hierarchy, capturing the inherent connections among representations [6], [7], [8]. Single-network-based and dual-network-based models are two popular DL-based streams for MTSC. A single-network-based model generally employs a unified (frequently hybridized) network architecture to perform both feature extraction and the identification of relationships within the data. For instance, He *et al.* [9] presented a general neural network model based on convolutional filtering, called Rel-CNN, to extract both local and global features in time series. Ma *et al.* [10] introduced a robust generation time series method, called AJ-RNN, which employs an adversarial joint-learning structure combined with recurrent neural network (RNN). Chen *et al.* [11] devised a paralleling attention structure that gracefully integrates two class of attention variants for MTSC. Chen *et al.* [12] put forward TAR-GAN, a rule mining-based approach that employs the generative adversarial technique to search useful shapelets of time series. On the other hand, a dual-network-based model usually comprises two parallel networks: one dedicated to extracting local features and the other aimed at capturing global relationships. The local feature network typically employs convolutional neural networks (CNNs) to discern local patterns, while the global relation network, commonly based on RNNs or attention-based networks, is responsible for uncovering connections among the extracted representations. Notable dual-network-based models includes the robust temporal feature network (RTFN) containing a temporal feature network and a long short-term memory (LSTM)-based attention network [13], LSTM-fully convolutional network (LSTM-FCN) consisting of LSTM-based network and FCN [14], residual attention net (ResNet-Transformer) composed of a residual network and a transformer-based network [15], robust semisupervised model (SelfMatch) comprising a residual network and an LSTM-based attention network [16], and densely knowledge-aware network (DKN) consisting of a residual multi-head convolutional network and a transformer-based network [17]. Unlike

single-network-based models, dual-network-based algorithms account for the characteristics of multivariate time series data, employing a divide-and-conquer approach to achieve commendable performance across various time series problems. However, dual-network-based algorithms may face the following two challenges:

- *Most dual-network-based algorithms use the concatenation feature aggregation method to integrate the output feature of the local feature network with that of the global relation network before reaching the final prediction layer. Nevertheless, this method may lead to an inadequate alignment between diverse feature types extracted by these two distinct networks. In particular, the lack of interaction between different feature types extracted at lower levels impedes the relationship exploration among various variables within multivariate time series data.*
- *As known, a learning model's performance usually hinges on quality of the semantic information extracted from lower and higher levels within the representation hierarchy [18], [19]. Thus, for an arbitrary dual-network-based model, enhancing the aggregation between features derived from the local feature network and those from the global relation network at each level significantly enriches the extracted semantic information, consequently enhancing the model's performance.*

Recently, a number of researchers have focused their efforts on aggregating features across various dimensions, enhance model's performance in diverse computer vision (CV) tasks [20]. Some studies considered the attention method into CNNs to combine the features on both spatial and channel dimensions, such as dual-attention network [21] and SCA-CNN [22]. Dosovitskiy *et al.* [23] proposed a spatial self-attention module to enhance long-range dependencies between image pixels. Chen *et al.* [24] devised Mixformer based on channel attention to fuse spatial and channel information in vision tasks. Chen *et al.* [25] developed a dual aggregation transformer network based on spatial window and channel-wise self-attention for image super-resolution. However, unlike CV, the feature aggregation methods in time series exhibit notable limitations, listed below.

- *Unlike image data, multivariate time series entails a sequence of chronologically arranged data points linked to several time-dependent variables. Thus, outstanding feature aggregation methods used in CV may not be directly transferable or suitable for time series domain.*
- *Currently, a significant gap exists in tailored methodologies capable of deeply integrating distinct feature types, e.g., local and global patterns, to capture long-range dependencies among multiple variables within multivariate time series data. This limitation is prominent in most dual-network-based algorithms, particularly in efficiently hybridizing the local feature network with the global relation network.*

To overcome the limitations above, we propose a knowledge aggregation transformer network (KATN) for MTSC. Different from most of the dual-network-based algorithms combining disparate feature types at higher levels through concatenation,

KATN deeply aggregates various feature types within each level. This aggregation can enrich the semantic information substantially, thereby effectively extracting long-range dependencies among multiple variables in multivariate time series data.

Our primary contributions are outlined below.

- This paper designs a deep feature aggregation network for MTSC, called KATN. KATN comprises four aggregation transformer blocks adept at amalgamating various feature types at lower and higher levels, exploring abundant regularizations and connections within the data, as depicted in Fig. 1.
- In each aggregation transformer block, a modified residual network (MResNet) and a multi-head attention network are responsible for local feature and global relation extraction, respectively. The block uses an additive operation to connect the output feature of MResNet with that of the multi-head attention network. A fully connected (i.e., dense) layer is employed for feature alignment. Meanwhile, this block activates neural units via the Gaussian error linear unit (GeLU) [26] function.
- The experiments demonstrate that in comparison to 6 state-of-the-art (SOTA) transformer algorithms, KATN achieves 9 wins, 6 ties, and 15 losses, along with the lowest AVG\_rank score of 2.400. Moreover, KATN surpasses 18 existing MTSC algorithms based on both 'win'/'tie'/'lose' metrics and AVG\_rank, as measured by top-1 accuracy. Specifically, among the 30 datasets, KATN secures victory in 13, achieving the lowest AVG\_rank score of 4.167.

The paper's structure for the remaining sections is: in Section II, we review the SOTA methods in MTSC and transformer-based algorithms, highlighting their strengths and limitations. Section III introduces the proposed KATN framework, elaborating on its architecture and core components. Following this, Section IV discusses the experimental setup, results, and comparative analysis. Finally, Section V concludes the paper by summarizing key findings and outlining potential future directions.

## II. RELATED WORK

This section reviews a number of MTSC and transformer algorithms.

### A. MTSC Algorithms

Traditional and DL-based algorithms are two main streams for MTSC.

1) *Traditional Algorithms*: Distance- and feature-based approaches are two typical representatives of traditional algorithms for MTSC [6], [11]. The classical distance-based method employs nearest neighbor (NN) and dynamic time warping (DTW) to measure the similarities between spatial features in the data, such as, dependent DTW, adaptive DTW, and independent DTW [27]. The ensemble learning algorithm with DTW and NN has been developed to address various time series problems, e.g., the elastic ensemble approach [28], collective of transformation-based ensemble (COTE) method [29], hierarchical voting COTE (HIVE-COTE) [30], improved

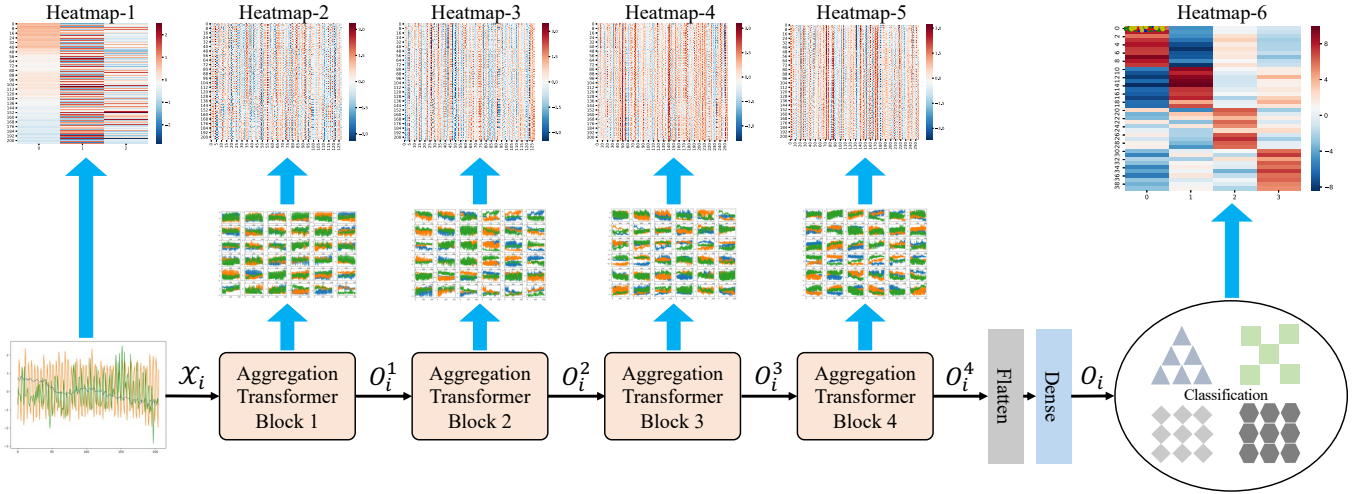


Fig. 1. Overview of KATN. The proposed KATN consists of four aggregation transformer blocks for local feature and global relation extraction in multivariate time series. We visualize the representations learned by each network block in the KATN, where the top picture visualizes the heatmap representation of each network layer. From these heatmaps, a noticeable pattern emerges: as the network’s depth increases, samples sharing similar characteristics tend to cluster together more prominently. Note:  $X_i$ ,  $i = 1, 2, \dots, N$ , is the  $i$ -th input sample, where  $N$  represents the number of input samples.  $O_i^j$ ,  $i = 1, 2, \dots, N$ ,  $j = 1, 2, 3, 4$ , denotes the  $i$ -th output feature vector of  $j$ -th aggregation transformer block, while  $O_i$  stands for the  $i$ -th output feature vector of KATN.

meta HIVE-COTE (HIVE-COTE 2.0) [31], and explainable-by-design ensemble approach [32].

Feature-based algorithms concentrating on mining the representative features from the input. For instance, Li *et al.* [33] introduced an efficient shapelet-based discovery approach called *bspcover* to construct a collection of high-quality shapelets. Baldán and Benítez [34] proposed an interpretable representation approach to handle a variety of MTSC applications. Shifaz *et al.* [35] devised a scalable and accurate forest method, namely TS-Chief, to achieve excellent performance on various time series tasks. The pattern similarity method [36], hidden-unit logistic model [37], time series forest [38], online rule-based classifier learning [39], autoregressive tree-based ensemble [40], bag of symbolic Fourier approximation symbols [41], WEASEL+MUSE [42], and fuzzy cognitive map [43] are all feature-based.

2) *DL-based Algorithms*: DL-based algorithms aim to model the internal data representation hierarchy, mining the inherent connections among representations [6]. Single- and dual-network-based models present research streams [7], [8]. The flexible multi-head linear attention (FMLA) [44], InceptionTime [45], dynamic temporal pooling network [46], fully-convolutional network [47], DA-Net [11], multi-process collaborative architecture [48], TAR-GAN [12], Rel-CNN [9], shapelet-neural network [49], AJ-RNN [10], deep contrastive representation learning with self-distillation [50], reservoir computing [51], ROCKET [52], dynamic graph attention autoencoder [53], dynamic component alignment [54], echo state network [55], and MiniROCKET [56] are well known single-network-based models. On the other hand, typical dual-network-based models include ResNet-Transformer [15], Self-Match [16], DKN [17], LSTM-FCN [14], SelfMatch [16], robust neural temporal search (RNTS) [57], perceptive and lightweight capsule models [58], and attentional prototypical

network (TapNet) [59].

### B. Transformer Algorithms

Since 2017, the Transformer model has gained popularity across various domains, such as time series, CV, natural language processing, and information retrieval, due to its capability to capture relationships among features at diverse positions [60]. The self-attention based transformer [61] serves as a trailblazer in the realm of transformers, with numerous enhanced transformer architectures emerging in this domain. For instance, to reduce the self-attention complexity in both time and space, Wang *et al.* [62] introduced Linformer based on linear attention. Liu *et al.* [63] proposed a hierarchical transformer using shifted windows to model different image scales. Ding *et al.* [64] presented a simple and effective dual attention transformer to capture global context while maintaining computational efficiency. Li *et al.* [65] put forward a dilated convolutional transformer-based generative adversarial network for time series anomaly detection. The transformer with spatial self-attention [23], DKN [17], Mixformer [24], FMLA [44], dual aggregation transformer [25], intention-aware dynamic transformer [66], dynamic graph transformer [67], and temporal graph transformer [68] are typical transformer models.

## III. THE PROPOSED KATN

This section first introduces the structure of KATN and its key components. Then, it describes the loss function.

### A. Overview

KATN comprises four aggregation transformer blocks responsible for amalgamating various feature types at lower and higher levels, mining abundant relationships and rules

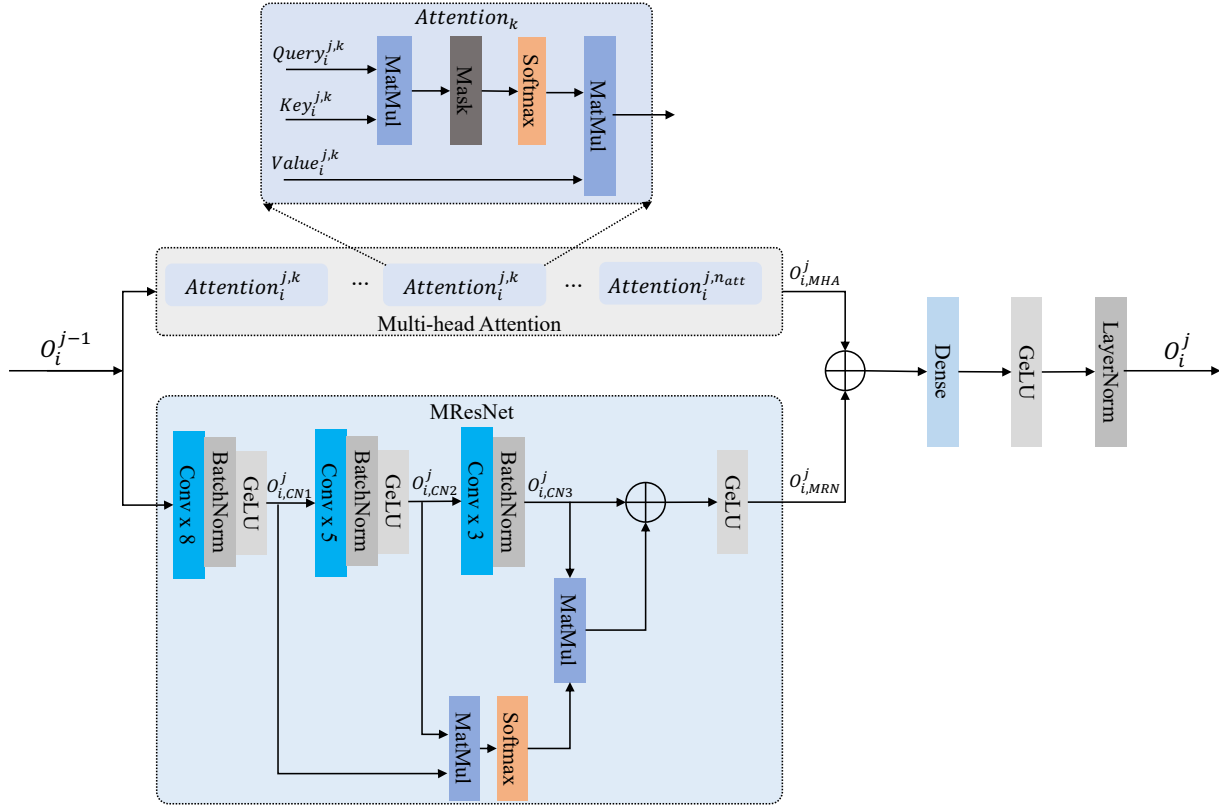


Fig. 2. Architecture of aggregation transformer block. This block integrates a modified residual network (MResNet) with a multi-head attention network in parallel. Similar to the vanilla ResNet [47], MResNet mainly consists of three 1-dimensional convolutional layers. However, MResNet adopts the attention-based method to combine the output features from these convolutional layers. Rather than the typically used rectified linear unit (ReLU) function, MResNet uses the Gaussian error linear unit (GeLU) [26] function for neural unit activation. The multi-head attention network consists of  $n_{att}$  attention modules mining relationships among features across diverse locations. Note: “MatMul” denotes the matrix multiplication operation. “Conv x 8” is a 1-dimensional convolutional layer with a kernel size of 8. “BatchNorm” and “LayerNorm” are the batch normalization and layer normalization operations, respectively.  $O_i^{j-1}$  and  $O_i^j$  are the input and output feature vectors of the  $j$ -th aggregation transformer block,  $j = 1, 2, 3, 4$ , associated with the input sample,  $\mathcal{X}_i$ ,  $i = 1, 2, \dots, N$ , where  $N$  represents the number of input samples. In particular,  $O_i^{j-1} = \mathcal{X}_i$ , when  $j = 1$ .

hidden in multivariate time series. The structure of KATN is shown in Fig. 1. Within each aggregation transformer block, MResNet and the multi-head attention network handle local feature and global relation extraction, respectively. This block connects the output feature of MResNet with that of the multi-head attention network using additive feature aggregation. It employs a fully connected (i.e., dense) layer for feature alignment. Neural units are activated by the GeLU function.

### B. Aggregation Transformer

Each aggregation transformer block adeptly combines local patterns from MResNet and global patterns from the multi-head attention network, capturing extensive dependencies among multiple variables within multivariate time series data. The architecture of an aggregation transformer block is depicted in Fig. 2.

Let  $\mathcal{X}_i$ ,  $i = 1, 2, \dots, N$ , denote the  $i$ -th input sample, where  $N$  represents the number of input samples. Given input,  $\mathcal{X}_i$ ,  $O_i^{j-1}$  and  $O_i^j$  are the input and output feature vectors of the  $j$ -th aggregation transformer block,  $j = 1, 2, 3, 4$ , respectively. Note that  $O_i^{j-1}$  is equal to  $\mathcal{X}_i$ , when  $j = 1$ .

1) *MResNet*: Similar to the vanilla ResNet [47], MResNet primarily comprises three 1-dimensional convolutional layers. But, MResNet employs an attention-based approach to integrate the output features from these convolutional layers. Different from ResNet that typically adopts the rectified linear unit (ReLU) function, MResNet utilizes GeLU activation for neural units.

Let  $O_{i,CN1}^j$ ,  $O_{i,CN2}^j$ , and  $O_{i,CN3}^j$  be the output feature vectors of three 1-dimensional convolutional layers in MResNet, respectively. These vectors are defined as:

$$\begin{aligned} O_{i,CN1}^j &= GeLU(BN(Conv(O_i^{j-1}))) \\ O_{i,CN2}^j &= GeLU(BN(Conv(O_{i,CN1}^j))) \\ O_{i,CN3}^j &= BN(Conv(O_{i,CN2}^j)) \end{aligned} \quad (1)$$

where,  $Conv()$ ,  $BN()$ , and  $GeLU()$  are the 1-dimensional convolution, batch normalization, and GeLU functions, respectively.

MResNet transforms the output feature vectors, namely  $O_{i,CN1}^j$ ,  $O_{i,CN2}^j$ , and  $O_{i,CN3}^j$ , from three 1-dimensional convolutional layers to an output vector,  $O_{i,MRN}^j$ , by an attention-

TABLE I  
OVERVIEW OF 30 UEA MULTIVARIABLE TIME SERIES BENCHMARK DATASETS. NOTES: AUDIO SPECTRA (AS), ELECTROCARDIOGRAM (ECG), ELECTROENCEPHALOGRAPH (EEG), HUMAN ACTIVITY RECOGNITION (HAR), AND MAGNETOENCEPHALOGRAPHY (MEG).

Dataset Index	Dataset Name	Dimensions	SeriesLength	TrainSize	NumClasses	TestSize	Type
AWR	ArticularyWordRecognition	9	144	275	25	300	Motion
AF	AtrialFibrillation	2	640	15	3	15	ECG
BM	BasicMotions	6	100	40	4	40	HAR
CT	CharacterTrajectories	3	182	1422	20	1436	Motion
CK	Cricket	6	1197	108	12	72	HAR
DDG	DuckDuckGeese	1345	270	50	5	50	AS
EW	EigenWorms	6	17984	128	5	131	Motion
EP	Epilepsy	3	206	137	4	138	HAR
EC	EthanolConcentration	3	1751	261	4	263	HAR
ER	ERing	4	65	30	6	270	Other
FD	FaceDetection	144	62	5890	2	3524	EEG/MEG
FM	FingerMovements	28	50	316	2	100	EEG/MEG
HMD	HandMovementDirection	10	400	160	4	74	EEG/MEG
HW	Handwriting	3	152	150	26	850	HAR
HB	Heartbeat	61	405	204	2	205	AS
IW	InsectWingbeat	200	30	30000	10	20000	AS
JV	JapaneseVowels	12	29	270	9	370	AS
LIB	Libras	2	45	180	15	180	HAR
LSST	LSST	6	36	2459	14	2466	Other
MI	MotorImagery	64	3000	278	2	100	EEG/MEG
NATO	NATOPS	24	51	180	6	180	HAR
PD	PenDigits	2	8	7494	10	3498	EEG/MEG
PEMS	PEMS-SF	963	144	267	7	173	EEG/MEG
PS	PhonemeSpectra	11	217	3315	39	3353	AS
RS	RacketSports	6	30	151	4	152	HAR
SRS1	SelfRegulationSCP1	6	896	268	2	293	EEG/MEG
SRS2	SelfRegulationSCP2	7	1152	200	2	180	EEG/MEG
SAD	SpokenArabicDigits	13	93	6599	10	2199	AS
SWJ	StandWalkJump	4	2500	12	3	15	ECG
UW	UWaveGestureLibrary	3	315	120	8	320	HAR

based approach.  $O_{i,MRN}^j$  is defined in Eq. (2).

$$O_{i,MRN}^j = GeLU(O_{i,CAF}^j + O_{i,CN3}^j) \quad (2)$$

where,

$$O_{i,CAF}^j = Softmax\left(\frac{O_{i,CN1}^j \cdot (O_{i,CN2}^j)^T}{\sqrt{d_{i,MRN}^{j,k}}}\right) \cdot O_{i,CN3}^j \quad (3)$$

where,  $(O_{i,CN2}^j)^T$  and  $d_{i,MRN}^{j,k}$  stand for the transpose and dimension of  $O_{i,CN2}^j$ , respectively.  $Softmax()$  outputs the mathematical possibilities of a give vector.

2) *Multi-head attention*: The multi-head attention network is composed of  $n_{att}$  attention modules designed to mine relationships among the features across diverse locations. In the  $j$ -th aggregation transformer block, attention module,  $Attention_{i,k}^{j,k}$ ,  $k = 1, 2, \dots, n_{att}$ , maps a query,  $Query_{i,k}^{j,k}$ , and a connection of key-value pairs,  $Key_{i,k}^{j,k}$ - $Value_{i,k}^{j,k}$ , to an output,  $O_{i,att}^{j,k}$ .  $O_{i,att}^j$  is defined as:

$$O_{i,att}^{j,k} = Softmax\left(\frac{Query_{i,k}^{j,k} \cdot (Key_{i,k}^{j,k})^T}{\sqrt{d_{i,MHA}^{j,k}}}\right) \cdot Value_{i,k}^{j,k} \quad (4)$$

where,  $(Key_{i,k}^{j,k})^T$  and  $d_{i,MHA}^{j,k}$  represent the transpose and dimension of  $Key_{i,k}^{j,k}$ , respectively.

Let  $O_{i,MHA}^j$  denote the output feature vector of the multi-head attention network in the  $j$ -th aggregation transformer block.  $O_{i,MHA}^j$  is calculated in Eq. (5).

$$O_{i,MHA}^j = CONCAT([O_{i,att}^{j,1}, O_{i,att}^{j,2}, \dots, O_{i,att}^{j,n_{att}}]) \quad (5)$$

where,  $CONCAT()$  is the concatenation function.

The final feature aggregation within each aggregation transformer block is realized through a structured fusion of local and global feature representations. Specifically, the local patterns, denoted as  $O_{i,MRN}^j$ , capture short-range dependencies and contextual semantics through a hierarchical convolutional process enhanced by internal attention mechanisms. In parallel, the global relationships captured by the multi-head attention output  $O_{i,MHA}^j$  encode long-term temporal dependencies and inter-variable interactions across the entire time window.

To merge these complementary representations, we adopt an additive feature fusion strategy:

$$O_i^j = LN(GeLU(Dense(O_{i,MHA}^j + O_{i,MRN}^j))) \quad (6)$$

where,  $O_i^j$  represent the output feature vector of the  $j$ -th aggregation transformer block.  $LN()$  and  $Dense()$  present the layer normalization and fully-connected (i.e., dense) functions, respectively.

This formulation ensures that both modalities contribute jointly to the final feature space, while the dense layer performs dimensional alignment and nonlinear transformation. The subsequent GeLU activation introduces smooth nonlinearity, and the layer normalization stabilizes the learning process by preventing internal covariate shift.

As the aggregation transformer blocks are stacked hierarchically, the output of each block  $O_i^j$  becomes the input to the next, progressively enriching the feature hierarchy. Ultimately, the last aggregation block outputs a comprehensive representation that integrates multi-scale dependencies—both local and global—across time and variables. This final representation

---

**Algorithm 1** Training and Inference Procedure of KATN
 

---

**Input:**  $\mathcal{D} = (\mathcal{D}_{train}, \mathcal{D}_{val}, \mathcal{D}_{test})$ ;     $\triangleright$  Training, validation, and test sets.

**Output:**  $\hat{Y}$ ;     $\triangleright$  Predicted labels produced by KATN.

- 1: Randomly initialize model parameters  $\theta_0$ ;
- 2: **for**  $m = 1$  to  $EPS$  **do**     $\triangleright EPS$ : number of training epochs

- 3:    **Begin Forward Pass:**

Feed each training sample  $\mathcal{X}_i \in \mathcal{D}_{train}$  through the KATN architecture (Fig. 1), consisting of four sequential aggregation transformer blocks.

- 4:    **for**  $j = 1$  to 4 **do**  $\triangleright$  Each block performs local-global feature fusion

- 5:       Obtain  $O_{i,CN1}^j$ ,  $O_{i,CN2}^j$ , and  $O_{i,CN3}^j$  via Eq. (1);

- 6:       Compute attention-weighted local feature:  $O_{i,CAF}^j$  using Eq. (3);

- 7:       Fuse local output:  $O_{i,MRN}^j = GeLU(O_{i,CAF}^j + O_{i,CN3}^j)$  via Eq. (2);

- 8:       For each head  $k$ , compute global attention:  $O_{i,att}^{j,k}$  using Eq. (4);

- 9:       Concatenate heads:  $O_{i,MHA}^j$  via Eq. (5);

- 10:       Aggregate local-global features:  $O_i^j = LN(GeLU(Dense(O_{i,MHA}^j + O_{i,MRN}^j)))$  via Eq. (6);

- 11:    **end for**

- 12:    **End Forward Pass**

- 13:    **Begin Loss Computation:**

Using Eq. (7), evaluate the objective:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N y_i \log(O_i) + \epsilon \|\theta\|_2^2$$

- 14:    **End Loss Computation**

- 15:    **Begin Parameter Update:**

Perform gradient descent:

$$\theta_m = \theta_{m-1} - \eta \nabla_{\theta_{m-1}} \mathcal{L}(\theta_{m-1})$$

where  $\eta$  is the learning rate.

- 16:    **End Parameter Update**

- 17:    **if**  $m > 1$  **then**

- 18:       Evaluate model on  $\mathcal{D}_{val}$  to monitor generalization.

- 19:    **end if**

- 20: **end for**

- 21: **Begin Inference:**

Deploy trained model on  $\mathcal{D}_{test}$  to generate final predictions  $\hat{Y}$ .

- 22: **End Inference**

---

is then passed to the classification head for downstream prediction.

### C. Loss Function

The loss function of KATN, denoted as  $\mathcal{L}$ , is formulated using the cross-entropy method. It quantifies the disparity between the ground-truth labels and the predicted vectors, as

expressed in Eq. (7).

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N y_i \log(O_i) + \epsilon \|\theta\|_2^2 \quad (7)$$

where,  $O_i$  stands for the  $i$ -th output feature vector of KATN.  $y_i$  is the  $i$ -th ground truth label.  $\theta$  denotes the KATN's parameters.  $\epsilon$  represents the coefficient of  $\|\theta\|_2^2$  (i.e.,  $L_2$  regularization). Following [13], [16], we set  $\epsilon = 0.0005$  in this paper. The implementation details of the proposed KATN framework are formally outlined in Algorithm 1, which provides a step-by-step description of its computational workflow.

## IV. PERFORMANCE AND EVALUATION

This section begins by detailing the experimental setup, including performance metrics and ablation study analysis. Subsequently, it evaluates the performance and computational efficiency of KATN through comprehensive experiments. The section concludes with an in-depth representation visualization analysis to further illustrate the practical applicability of the proposed method.

### A. Experimental Setup

1) *Dataset Description:* As suggested in [11], [13], [14], [32], [34], [40], we use the University of East Anglia (UEA) multivariate time series archive in 2018 [69] for algorithmic performance evaluation. This archive covers a range of categories, varying from 2 to 39, with time series lengths ranging from 8 to 17,984. It includes data from 7 application scenarios, such as human activity recognition and motion analysis. Table I shows more details about the UEA archive.

TABLE II  
HYPER-PARAMETER SETTINGS OF FOUR AGGREGATION TRANSFORMER BLOCKS.

Aggregation Transformer No.	$n_{att}$	MResNet's Channels	Dense Layer's Units
1	8	128	128
2	8	128	128
3	16	256	256
4	16	256	256

2) *Implementation Details:* The hyper-parameters for the four aggregation transformer blocks are summarized in Table II. For optimization, we employ the Adam Optimizer with the following configurations: a momentum term of 0.95, an initial learning rate of 0.001, and a decay value of 0.9. All experiments are performed on a computational platform running Ubuntu 18.04 OS and Python 3.6, powered by an Nvidia RTX 2080Ti GPU, Tensorflow 1.18, and an AMD R5 1400 CPU with 32GB RAM. To illustrate the training dynamics of KATN, Fig. 3 visualizes the loss curves across multiple datasets, including CharacterTrajectories, Cricket, FaceDetection, Heartbeat, LSST, PenDigits, PhonemeSpectra, SpokenArabicDigits, and StandWalkJump.

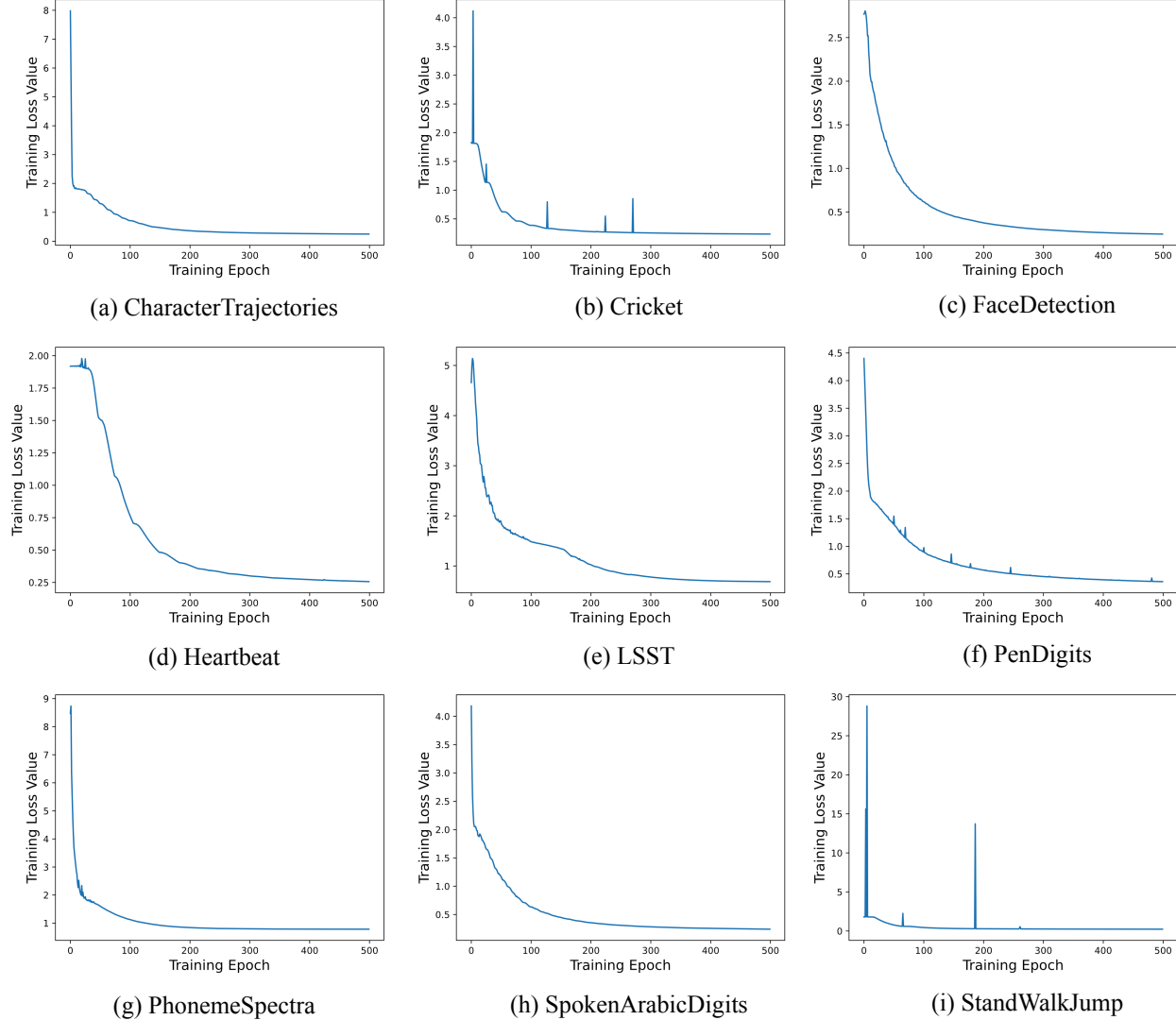


Fig. 3. Training loss values obtained during training on the CharacterTrajectories, Cricket, FaceDetection, Heartbeat, LSST, PenDigits, PhonemeSpectra, SpokenArabicDigits, and StandWalkJump datasets.

### B. Performance Metrics

To validate the effectiveness of the proposed KATN, two commonly used metrics, ‘win’/‘tie’/‘lose’ and AVG\_rank, both based on the top-1 accuracy, are considered. As suggested in [6], [7], [10], [11], [13], [14], [46], for any MTSC algorithm, the ‘win’, ‘tie’, and ‘lose’ scores indicate on how many datasets an algorithm performs better than, equivalent to, or worse than compared algorithms, respectively. The ‘best’ score represents the sum of the ‘win’ and ‘tie’ scores of the algorithm. In addition, similar to prior studies in [6], [7], [13], [14], [15], [16], [46], [57], [59], we use AVG\_rank, which distinguishes between different algorithms based on the Wilcoxon signed-rank test with Holm’s alpha correction at a significance level of 5%.

### C. Ablation Study

To comprehensively validate the design choices of the proposed KATN model, we conduct an integrated ablation

study encompassing both hyper-parameter sensitivity analysis and architectural component dissection. This dual-perspective analysis enables a rigorous examination of KATN’s robustness and design effectiveness across 30 UEA benchmark MTSC datasets. Specifically, we partition the study into two subparts: (1) hyper-parameter sensitivity, focusing on structural configuration choices; and (2) architectural component analysis, isolating key modules to quantify their individual contributions.

1) *Hyper-parameter Sensitivity Analysis*: We begin by evaluating the impact of critical structural hyper-parameters that were systematically optimized during model development. This includes (i) the number of aggregation transformer blocks, and (ii) the activation function employed within each block. These parameters fundamentally affect the model’s representation capacity, learning dynamics, and overall expressivity.

a) *Effect of aggregation transformer depth*: The number of aggregation transformer blocks determines the hierarchical

TABLE III  
TOP-1 ACCURACY RESULTS WITH DIFFERENT AGGREGATION  
TRANSFORMER BLOCKS ACROSS 30 UEA BENCHMARK DATASETS.

Dataset Index	KATN-(1)	KATN-(2)	KATN-(3)	KATN	KATN-(5)
AWR	0.953	0.970	<b>0.993</b>	<b>0.993</b>	<b>0.993</b>
AF	0.333	0.400	0.467	<b>0.533</b>	<b>0.533</b>
BM	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
CT	0.741	0.935	0.951	<b>0.993</b>	0.986
CK	0.917	0.958	0.986	<b>1.000</b>	<b>1.000</b>
DDG	0.600	0.560	0.600	0.660	<b>0.680</b>
EW	0.527	0.727	0.684	<b>0.733</b>	<b>0.733</b>
EP	0.935	0.935	0.986	<b>0.913</b>	<b>0.913</b>
EC	0.321	0.349	0.372	0.399	<b>0.411</b>
ER	0.778	0.822	0.907	0.956	<b>0.968</b>
FD	0.518	0.528	0.614	<b>0.664</b>	<b>0.664</b>
FM	0.500	0.570	0.590	<b>0.620</b>	0.610
HMD	0.270	0.392	0.649	<b>0.688</b>	<b>0.688</b>
HW	0.192	0.294	0.287	<b>0.320</b>	<b>0.320</b>
HB	0.724	0.717	0.761	<b>0.762</b>	<b>0.762</b>
IW	0.228	0.302	0.491	0.515	<b>0.567</b>
JV	0.941	0.949	<b>0.978</b>	0.968	0.968
LIB	0.478	0.744	0.772	<b>0.839</b>	<b>0.839</b>
LSST	0.337	0.456	<b>0.652</b>	0.407	0.407
MI	0.530	0.600	0.580	0.620	<b>0.630</b>
NATO	0.850	0.889	<b>0.900</b>	0.889	0.889
PD	0.970	0.977	0.977	<b>0.982</b>	<b>0.982</b>
PEMS	0.832	0.888	0.939	<b>0.953</b>	<b>0.953</b>
PS	0.288	0.325	<b>0.466</b>	0.366	<b>0.466</b>
RS	0.796	0.829	<b>0.914</b>	0.875	0.875
SRS1	0.805	0.805	0.839	<b>0.908</b>	<b>0.908</b>
SRS2	0.511	0.550	0.561	<b>0.600</b>	<b>0.600</b>
SAD	0.872	0.953	0.963	<b>0.983</b>	<b>0.983</b>
SWJ	0.267	0.400	0.426	<b>0.600</b>	<b>0.600</b>
UW	0.759	0.868	0.868	0.881	<b>0.897</b>
Win	0	0	4	2	<b>7</b>
Tie	1	1	3	<b>17</b>	<b>17</b>
Lose	29	29	23	11	<b>6</b>
Best	1	1	7	19	<b>24</b>
Mean Accuracy	0.626	0.690	0.739	0.754	<b>0.761</b>

TABLE IV  
TOP-1 ACCURACY RESULTS WITH DIFFERENT ACTIVATION FUNCTIONS  
ACROSS 30 UEA BENCHMARK DATASETS.

Dataset Index	ReLU	Leaky ReLU	EReLU	PreLU	GeLU
AWR	0.987	<b>0.993</b>	<b>0.993</b>	0.987	<b>0.993</b>
AF	0.400	0.467	0.467	0.467	<b>0.533</b>
BM	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
CT	0.979	0.983	<b>0.997</b>	0.983	0.993
CK	0.986	<b>1.000</b>	0.986	0.993	<b>1.000</b>
DDG	0.620	0.640	0.640	0.600	<b>0.660</b>
EW	0.712	0.712	<b>0.733</b>	0.727	<b>0.733</b>
EP	0.935	0.935	0.935	<b>0.952</b>	0.913
EC	0.349	0.352	0.352	0.372	<b>0.399</b>
ER	0.933	0.933	0.933	0.919	<b>0.956</b>
FD	0.631	0.614	0.634	<b>0.664</b>	<b>0.664</b>
FM	0.590	0.600	0.600	0.610	<b>0.620</b>
HMD	0.662	0.662	0.650	0.650	<b>0.688</b>
HW	0.308	0.310	<b>0.320</b>	<b>0.320</b>	<b>0.320</b>
HB	0.727	0.771	<b>0.765</b>	0.724	0.762
IW	0.491	<b>0.515</b>	<b>0.515</b>	<b>0.515</b>	<b>0.515</b>
JV	0.941	<b>0.978</b>	0.965	0.941	0.968
LIB	0.833	0.833	0.833	<b>0.850</b>	0.839
LSST	0.389	0.416	0.407	<b>0.421</b>	0.407
MI	0.590	0.610	<b>0.630</b>	0.610	0.620
NATO	<b>0.900</b>	0.872	0.872	0.872	0.889
PD	0.977	0.980	0.980	0.981	<b>0.982</b>
PEMS	0.930	0.930	<b>0.953</b>	<b>0.953</b>	<b>0.953</b>
PS	0.342	0.339	0.359	0.315	<b>0.366</b>
RS	0.856	0.856	0.868	0.803	<b>0.875</b>
SRS1	0.874	0.889	0.841	0.853	<b>0.908</b>
SRS2	0.533	0.533	0.533	0.500	<b>0.600</b>
SAD	0.972	0.963	0.978	0.980	<b>0.983</b>
SWJ	0.427	0.400	0.533	0.400	<b>0.600</b>
UW	0.869	0.859	0.903	<b>0.881</b>	<b>0.881</b>
Win	1	1	3	3	<b>13</b>
Tie	1	4	6	6	<b>9</b>
Lose	28	25	21	21	<b>8</b>
Best	2	5	9	9	<b>22</b>
Mean Accuracy	0.725	0.731	0.739	0.728	<b>0.754</b>

depth at which multi-scale temporal and cross-view dependencies are extracted. To investigate its effect, we construct five model variants with increasing aggregation depth:

- **KATN-(1)**: employs a single aggregation transformer block.
- **KATN-(2)**: integrates two aggregation transformer blocks.
- **KATN-(3)**: integrates three aggregation transformer blocks.
- **KATN**: the default version of our model, comprising four aggregation transformer blocks.
- **KATN-(5)**: extends the model to five aggregation transformer blocks.

The experimental results, summarized in Table III, reveal a consistent trend: increasing the number of transformer blocks yields progressively better classification accuracy, particularly for complex datasets such as InsectWingbeat. For instance, KATN-(5) achieves the highest accuracy across several challenging benchmarks, confirming the benefit of deeper aggregation in modeling intricate temporal patterns.

Nevertheless, the performance gain between KATN-(5) and the default KATN is marginal, with a mean accuracy improvement of only 0.007 across the 30 UEA datasets. Furthermore, the parameter count increases substantially from 3.368M (KATN) to 4.814M (KATN-(5)) on InsectWingbeat. This observation reveals a point of diminishing returns, wherein additional blocks yield negligible performance improvements while incurring higher computational costs. Thus, the default

configuration with four aggregation transformer blocks is adopted as it offers a desirable balance between accuracy and efficiency.

*b) Effect of activation function choice:* Activation functions regulate non-linear transformations within each transformer block, thus playing a pivotal role in determining the expressive capacity of the network. To evaluate their effect, we test KATN under five widely used activation functions: ReLU, Leaky ReLU, Elastic ReLU (EReLU), Parametric ReLU (PreLU), and GeLU. The performance results across the 30 UEA datasets are presented in Table IV.

Among the candidates, GeLU exhibits the most stable and superior performance, with a ‘win’/‘tie’/‘lose’ count of 13/9/8 and a mean accuracy of 0.754. Its stochastic, smooth activation profile enhances gradient flow and enables finer feature selectivity—traits particularly valuable for modeling highly dynamic multivariate sequences. Consequently, GeLU is adopted as the default activation function in the final model architecture.

*c) Hyper-parameter synthesis and interpretation:* Taken together, the above analyses suggest that KATN exhibits high resilience to architectural variations, maintaining competitive performance across a spectrum of design choices. Although incremental increases in depth and activation sophistication do contribute to performance gains, the improvements quickly plateau as architectural complexity grows. This observation underscores KATN’s structural efficiency: its default configuration—comprising four aggregation transformer blocks



TABLE V  
TOP-1 ACCURACY RESULTS OBTAINED BY VARIOUS KATN VARIANTS  
ACROSS 30 UEA BENCHMARK DATASETS.

Dataset Index	KATN-w/o-MResNet	KATN-w/o-MHA	KATN-w/o-FC	KATN-MatMul	KATN-CONCAT	KATN-INTERN	KATN
AWR	0.973	0.970	0.980	0.953	<b>0.993</b>	0.992	<b>0.993</b>
AF	0.467	0.400	0.500	0.200	<b>0.533</b>	0.467	<b>0.533</b>
BM	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
CT	0.986	0.964	0.969	0.931	0.986	0.986	<b>0.993</b>
CK	0.943	0.944	0.986	0.958	0.951	0.986	<b>1.000</b>
DDG	0.520	0.540	0.600	0.600	0.560	0.620	<b>0.660</b>
EW	0.618	0.549	0.727	0.712	0.811	<b>0.847</b>	0.733
EP	0.807	0.666	0.935	<b>1.000</b>	0.979	<b>1.000</b>	0.913
EC	0.332	0.293	0.228	<b>0.487</b>	0.372	0.380	0.399
ER	0.941	0.930	0.930	0.859	0.933	<b>0.981</b>	0.956
FD	0.519	0.519	0.631	0.508	0.631	0.631	0.664
FM	0.560	0.550	0.550	0.570	<b>0.630</b>	0.450	0.620
HMD	0.338	0.278	0.378	0.312	0.662	0.392	<b>0.688</b>
HW	0.451	0.382	0.316	0.308	<b>0.511</b>	0.296	0.320
HB	0.619	0.619	0.727	0.518	0.765	<b>0.771</b>	0.762
IW	0.327	0.128	0.100	0.491	0.362	<b>0.595</b>	0.515
JV	0.918	0.924	0.949	0.941	0.930	<b>0.989</b>	0.968
LIB	0.889	0.833	0.817	0.806	<b>0.900</b>	0.878	0.839
LSST	<b>0.548</b>	0.456	0.337	0.265	0.347	0.456	0.407
MI	0.510	0.510	0.530	0.550	0.570	0.550	<b>0.620</b>
NATO	0.822	0.850	<b>0.900</b>	0.839	0.872	0.850	0.889
PD	0.939	0.973	0.969	0.831	0.948	0.965	<b>0.982</b>
PEMS	0.874	0.705	0.959	<b>0.994</b>	0.930	0.522	0.953
PS	0.215	0.104	0.247	0.269	0.325	0.292	<b>0.366</b>
RS	0.796	0.868	0.809	0.823	<b>0.879</b>	0.868	0.875
SRS1	0.843	0.771	0.874	0.724	<b>0.913</b>	0.874	0.908
SRS2	0.533	0.483	0.500	0.494	0.533	0.522	<b>0.600</b>
SAD	0.963	0.967	0.979	0.959	0.963	0.729	<b>0.983</b>
SWJ	0.426	0.200	0.400	0.267	0.533	0.333	<b>0.600</b>
UW	0.872	0.881	0.833	0.684	0.897	<b>0.916</b>	0.881
Win	1	0	1	2	5	6	<b>10</b>
Tie	1	1	1	2	<b>3</b>	2	<b>3</b>
Lose	28	29	28	26	22	22	<b>17</b>
Best	2	1	2	4	8	8	<b>13</b>
AVG_rank	4.833	5.383	4.283	5.183	2.867	3.333	<b>2.117</b>

and GeLU activation—not only delivers a robust accuracy-complexity trade-off but also demonstrates the model’s scalability and generalization stability under diverse task conditions.

2) *Architectural Component Analysis*: Beyond hyperparameters, we assess the significance of core architectural modules through controlled ablations. This includes the multi-head attention mechanism, MResNet backbone, the fully connected projection layer, and the feature aggregation strategy.

a) *Impact of multi-head attention and MResNet*: We begin by evaluating the respective contributions of the multi-head attention module and the MResNet backbone. To this end, we examine two reduced variants of KATN:

- KATN-w/o-MHA: a variant in which the multi-head attention mechanism is removed.
- KATN-w/o-MResNet: a variant excluding the MResNet component.

The performance comparison, presented in Table V, reveals that the complete KATN architecture surpasses KATN-w/o-MHA on 27 datasets and outperforms KATN-w/o-MResNet on 26. These results highlight the complementary roles of global attention and local convolution in capturing temporal dynamics across varying scales. The attention mechanism enables the model to uncover long-range dependencies that are often elusive to convolutional filters, whereas MResNet preserves local continuity essential for fine-grained pattern recognition.

Despite the gains in accuracy, these modules introduce distinct computational costs. As summarized in Table VI, the

average CPU inference time across the 30 UEA benchmark testing datasets increases from 90.773s (without MHA) and 75.423s (without MResNet) to 106.620s in the full model. This increase, however, is offset by substantial performance benefits, illustrating that the integration of attention and residual learning yields a synergistic effect critical to robust MTSC.

b) *Effectiveness of fully connected layer*: To assess the necessity of the fully connected (dense) layer employed within each aggregation transformer block, we introduce a variant termed KATN-w/o-FC, in which this layer is eliminated. The fully connected layer serves to harmonize the latent spaces of MResNet and the attention mechanism by projecting heterogeneous features into a unified representation space.

As indicated in Table V, the absence of the fully connected layer results in degraded performance on 28 out of 30 datasets, demonstrating that direct feature fusion without alignment significantly impairs the model’s ability to integrate local and global information. In particular, the fully connected layer mitigates semantic incompatibility between disparate feature sources, thereby enabling coherent and discriminative joint representations.

The computational overhead of incorporating the fully connected layer is minimal. Based on evaluations on the 30 UEA benchmark testing datasets in Table VI, the average CPU inference time increases marginally from 101.716s to 106.620s. Given its role in facilitating effective cross-modal integration, this cost is well-justified by the substantial accuracy gains.

c) *Effectiveness of additive feature aggregation*: To investigate the impact of the feature fusion strategy within each aggregation block, we compare KATN with three alternative aggregation designs:

- KATN-MatMul: KATN with multiplicative feature aggregation, as shown in Fig. 4 (a).
- KATN-CONCAT: KATN with concatenation feature aggregation, as depicted in Fig. 4 (b).
- KATN-INTERN: KATN with interactive feature aggregation, as detailed in Fig. 4 (c).

The quantitative results in Table V demonstrate that KATN, utilizing additive aggregation, achieves superior accuracy on 26, 18, and 20 datasets when compared to KATN-MatMul, KATN-CONCAT, and KATN-INTERN, respectively. These outcomes underline the advantage of additive fusion, which offers a lightweight yet semantically balanced integration of local and global descriptors without exacerbating feature dimensionality or causing representational entanglement.

The efficiency comparisons on the 30 UEA benchmark testing datasets further reinforce the merits of this design. According to Table VI, KATN records an average CPU inference time of 106.620s, whereas KATN-MatMul and KATN-CONCAT require 108.602s and 111.526s on CPU, respectively. Notably, KATN-INTERN incurs a prohibitive cost of 663.790s due to its intricate feature interaction operations. These findings suggest that although alternative schemes may enable more complex representation interactions, they do so at the expense of efficiency and, in many cases, accuracy.

Moreover, while concatenation increases feature dimensionality and often causes dominance of one representation stream over another, multiplicative and interaction-based approaches



TABLE VII  
TOP-1 ACCURACY RESULTS OBTAINED BY VARIOUS TRANSFORMER  
VARIANTS ACROSS 30 UEA BENCHMARK DATASETS.

Dataset Index	VanTrans	SwinTrans	Linformer	FMLA	DualAggTrans	DualAttTrans	KATN
AWR	0.973	0.980	0.987	0.992	<b>0.993</b>	<b>0.993</b>	<b>0.993</b>
AF	0.467	0.467	0.467	<b>0.533</b>	0.467	0.467	<b>0.533</b>
BM	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
CT	0.966	0.970	<b>0.997</b>	0.979	0.989	0.986	0.993
CK	0.943	0.958	0.958	0.986	0.986	0.951	<b>1.000</b>
DDG	0.540	0.580	0.600	0.620	<b>0.660</b>	0.640	<b>0.660</b>
EW	0.618	0.847	0.489	<b>0.954</b>	0.733	0.628	0.733
EP	0.807	0.978	0.971	0.979	0.935	<b>1.000</b>	0.913
EC	0.332	0.228	0.323	0.380	0.352	0.372	<b>0.399</b>
ER	0.400	0.930	0.859	<b>0.981</b>	0.933	0.933	0.956
FD	0.519	0.583	0.556	0.631	0.634	0.631	<b>0.664</b>
FM	0.560	0.460	0.530	0.550	0.600	<b>0.630</b>	0.620
HMD	0.338	0.284	0.378	0.392	0.650	0.662	<b>0.688</b>
HW	0.451	0.187	0.357	<b>0.511</b>	0.320	0.294	0.320
HB	0.619	0.727	0.751	<b>0.771</b>	0.765	0.765	0.762
IW	0.327	0.100	0.208	0.362	0.515	<b>0.595</b>	0.515
JV	0.918	0.778	0.965	0.965	<b>0.989</b>	0.930	0.968
LIB	0.889	0.817	0.850	0.850	<b>0.878</b>	0.833	0.839
LSST	0.548	<b>0.652</b>	0.568	0.643	0.407	0.347	0.407
MI	0.510	0.500	0.590	0.550	<b>0.630</b>	0.620	0.620
NATO	0.822	0.800	0.850	0.850	0.850	0.872	<b>0.889</b>
PD	0.939	0.939	0.980	0.965	0.980	0.948	<b>0.982</b>
PEMS	0.874	<b>0.959</b>	0.751	0.522	0.953	0.930	0.953
PS	0.215	0.247	0.175	0.292	0.359	0.325	<b>0.366</b>
RS	0.796	0.809	0.868	0.868	0.823	0.829	<b>0.875</b>
SRS1	0.843	0.771	0.652	0.874	0.841	<b>0.913</b>	0.908
SRS2	0.533	0.450	0.550	0.522	0.533	<b>0.600</b>	<b>0.600</b>
SAD	0.963	0.979	<b>0.983</b>	0.959	0.978	0.963	<b>0.983</b>
SWJ	0.426	0.267	0.400	0.333	0.533	0.533	<b>0.600</b>
UW	0.872	0.728	0.894	<b>0.916</b>	0.903	0.897	0.881
Win	0	2	1	5	3	3	<b>9</b>
Tie	1	1	2	2	3	4	<b>6</b>
Lose	29	27	27	23	24	23	<b>15</b>
Best	1	3	3	7	6	7	<b>15</b>
AVG_rank	5.483	5.500	4.500	3.467	3.167	3.483	<b>2.400</b>

2) *Comparison with existing MTSC algorithms:* To examine the performance of KATN, we compare it with 18 existing MTSC algorithms, listed below.

- MLP: the multi-layer perceptron model applied to MTSC tasks [47].
- FCN: the fully convolutional network model applied to MTSC tasks [47].
- ResNet: the residual neural network model applied to MTSC tasks [47].
- InceptionTime: the neural network model based on Inception blocks applied to MTSC tasks [45].
- Three distance-based baseline algorithms:  $ED_I$ ,  $DTW_I$ , and  $DTW_D$  [6].
- WM: the statistical feature selection approach based on bag-of-pattern (WEASEL+MUSE) [42].
- CBOSS: the contractable approximation symbols method using bag of symbolic Fourier [41].
- MLCN: the multivariate LSTM-FCN [14].
- TSF: the time series forest approach for MTSC [38].
- TapNet: the attentional prototype network model that gracefully embeds traditional classification and DL-based methods [59].
- XEM: the explainable-by-design ensemble method based on boosting-bagging and bias-variance [32].
- CMFM+SVM: the complexity measure-and-feature approach combined with SVM applied to MTSC tasks [34].
- MiniROCKET: the transform approach using fast deterministic technique, an improved version of ROCKET

[56].

- DA-Net: the dual attention-based network model that deeply incorporates squeeze-excitation and sparse attention variants [11].
- Conv-GRU: the convolutional network model combined with gated linear unit structures [70].
- DKN: the densely knowledge-aware network [17].

Table VIII presents the top-1 accuracy results for the various MTSC algorithms evaluated. KATN outperforms all other algorithms, achieving a ‘win’/‘tie’/‘loss’ score of 6/7/17, and securing the lowest AVG\_rank score of 4.167. DKN follows closely behind, ranking second in both the ‘best’ metric and AVG\_rank. MiniROCKET ranks as the second-best algorithm according to the ‘best’ metric, while XEM achieves the third-best position in terms of AVG\_rank. In contrast, MLP delivers the poorest performance among the algorithms tested.

The following provides an explanation for the results presented above. KATN effectively combines the local patterns extracted by MResNet and the global patterns extracted by the multi-head attention network, capturing long-range dependencies across multiple variables. DKN facilitates the transfer of knowledge between lower- and higher-level semantic information, effectively regularizing the model and contributing to its strong performance. MiniROCKET, on the other hand, employs simple linear classifiers alongside randomly initialized convolutional kernels to capture multi-scale subsequence features in time series, which enhances its efficiency in time series classification. XEM employs explicit boosting-bagging techniques and addresses the bias-variance trade-off, effectively uncovering potential correlations across different dimensions. In contrast, MLP, characterized by a shallow neural network structure, fail to mine the inherent relationships within the data, especially the interconnections among variables.

### E. Computational Complexity Assessment

Based on the guidelines outlined in [71], we conduct a comprehensive comparison between the proposed KATN model and 13 representative DL algorithms in terms of parameter count, FLOPs, and inference time, using the 30 UEA benchmark testing datasets as the evaluation foundation. These comparative methods consist of six transformer-based models, namely VanTrans [61], SwinTrans [63], Linformer [41], FMLA [44], DualAggTrans [25], and DualAttTrans [64], as well as seven MTSC baseline models. The baseline group includes four single-network-based methods—InceptionTime [45], ResNet [47], MLP [47], and FCN [47]—and three dual-network-based architectures—TapNet [59], MLCN [14], and DKN [17]. The aggregated results of this comparison are summarized in Table IX.

A close examination reveals that transformer-based models generally impose substantially higher computational complexity. Across the 30 UEA benchmark testing datasets, most transformer variants exhibit larger parameter counts and higher FLOPs than the seven MTSC baseline models. The proposed KATN, while demonstrating robust representational capacity, presents a relatively high computational footprint on the majority of datasets, with only DualAggTrans and DualAttTrans

TABLE VIII  
TOP-1 ACCURACY RESULTS OBTAINED BY VARIOUS MTSC ALGORITHMS ACROSS 30 UEA BENCHMARK DATASETS.

Dataset Index	MLP	FCN	ResNet	Inception Time	$ED_I$	$DTW_D$	$DTW_I$	CBOSS	WM	TSF	Mini ROCKET	XEM	TapNet	MLCN	CMFM + SVM	Conv-GRU	DA-Net	DKN	KATN
AWR	0.043	0.823	0.943	0.897	0.970	0.987	0.980	0.990	<b>0.993</b>	0.953	0.992	<b>0.993</b>	0.987	0.957	0.973	0.973	0.980	<b>0.993</b>	<b>0.993</b>
AF	0.400	0.200	0.200	0.267	0.267	0.220	0.267	0.267	0.267	0.200	0.133	0.467	0.333	0.333	0.267	0.467	0.467	0.467	<b>0.533</b>
BM	0.875	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.676	0.975	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.875	0.975	<b>1.000</b>	0.925	<b>1.000</b>	<b>1.000</b>
CT	0.056	0.741	0.977	0.935	0.964	0.989	0.969	0.986	0.990	0.931	0.065	0.979	0.997	0.917	0.970	0.966	<b>0.998</b>	0.986	0.993
CK	0.111	0.917	0.958	0.958	0.944	<b>1.000</b>	0.986	N/A	0.986	N/A	0.986	0.986	0.958	N/A	0.958	0.943	0.861	0.951	<b>1.000</b>
DDG	0.360	0.600	0.600	0.560	0.275	0.600	0.550	0.480	0.575	0.460	0.650	0.375	0.575	0.380	0.420	0.540	0.520	0.560	<b>0.660</b>
EW	0.233	0.684	0.833	0.727	0.549	0.618	N/A	0.511	0.890	0.712	<b>0.954</b>	0.527	0.489	0.330	0.847	0.811	0.489	0.628	0.733
EP	0.312	0.935	0.964	0.935	0.666	0.964	0.978	0.979	0.993	<b>1.000</b>	<b>1.000</b>	0.986	0.971	0.732	0.978	0.978	0.883	0.979	0.913
EC	0.300	0.349	0.317	0.321	0.293	0.323	0.304	0.304	0.316	<b>0.487</b>	0.380	0.372	0.323	0.373	0.228	0.332	0.338	0.372	0.399
ER	0.159	0.778	0.907	0.822	0.133	0.133	0.133	0.919	0.133	0.859	<b>0.981</b>	0.200	0.133	0.941	0.930	0.400	0.874	0.933	0.956
FD	0.565	0.518	0.534	0.528	0.519	0.529	N/A	0.513	0.545	0.508	0.631	0.614	0.556	0.555	0.583	0.640	0.648	0.631	<b>0.664</b>
FM	0.510	0.510	0.460	0.500	0.550	0.530	0.520	0.519	0.540	0.562	0.450	0.590	0.530	0.580	0.460	0.580	0.510	0.600	<b>0.620</b>
HMD	0.216	0.270	0.216	0.392	0.278	0.231	0.306	0.292	0.378	0.312	0.392	0.649	0.378	0.544	0.284	0.338	0.365	0.662	<b>0.688</b>
HW	0.038	0.192	0.382	0.294	0.200	0.286	0.316	0.504	<b>0.531</b>	0.191	0.511	0.287	0.357	0.305	0.187	0.451	0.159	0.231	0.320
HB	0.665	0.724	0.716	0.717	0.619	0.717	0.658	0.564	0.727	0.518	<b>0.771</b>	0.761	0.751	0.458	0.727	0.746	0.624	0.765	0.762
IW	0.104	0.491	0.231	0.302	0.128	N/A	N/A	N/A	N/A	N/A	<b>0.595</b>	0.228	0.208	N/A	0.100	0.208	0.567	0.362	0.515
JV	0.114	0.941	0.924	0.949	0.924	0.949	0.959	N/A	0.978	N/A	0.989	0.978	0.965	N/A	0.778	<b>0.991</b>	0.938	0.930	0.968
LIB	0.078	0.478	0.844	0.744	0.833	0.870	0.894	0.894	0.894	0.806	0.878	0.772	0.850	0.850	0.817	0.889	0.800	<b>0.900</b>	0.839
LSST	0.326	0.337	0.232	0.290	0.456	0.551	0.575	0.458	0.628	0.265	0.643	<b>0.652</b>	0.568	0.390	<b>0.652</b>	0.548	0.560	0.347	0.407
MI	0.530	0.590	0.560	0.530	0.510	N/A	N/A	0.390	0.500	0.550	0.550	0.600	0.590	0.510	0.500	0.512	0.500	<b>0.620</b>	<b>0.620</b>
NATO	0.167	0.900	0.878	0.889	0.850	0.883	0.850	0.850	0.883	0.839	0.928	0.916	<b>0.939</b>	0.900	0.800	0.916	0.878	0.872	0.889
PD	0.211	0.970	0.973	0.977	0.973	0.977	0.939	0.939	0.969	0.831	0.965	0.977	0.980	0.979	0.665	0.939	0.980	0.948	<b>0.982</b>
PEMS	0.340	0.832	0.828	0.888	0.705	0.711	0.734	0.730	N/A	<b>0.994</b>	0.522	0.942	0.751	0.745	0.959	0.874	0.867	0.930	0.953
PS	0.414	0.466	0.466	0.466	0.104	0.151	0.151	0.151	0.190	0.269	0.292	0.288	0.175	0.151	0.247	0.215	0.093	0.525	0.366
RS	0.276	0.796	0.836	0.829	0.868	0.803	0.842	0.854	0.914	0.823	0.868	<b>0.941</b>	0.868	0.856	0.809	0.888	0.803	0.879	0.875
SRS1	0.686	0.805	0.761	0.805	0.771	0.775	0.765	0.765	0.744	0.724	0.874	0.839	0.652	0.908	0.771	0.843	<b>0.924</b>	0.913	0.908
SRS2	0.456	0.511	0.511	0.561	0.483	0.539	0.533	0.533	0.522	0.494	0.522	0.550	0.550	0.506	0.450	0.566	0.561	<b>0.600</b>	<b>0.600</b>
SAD	0.108	0.729	0.932	0.872	0.967	0.963	0.959	N/A	0.982	N/A	0.100	0.973	<b>0.983</b>	N/A	0.979	0.963	0.980	0.963	<b>0.983</b>
SWJ	0.200	0.267	0.133	0.133	0.200	0.200	0.333	0.333	0.333	0.267	0.333	0.400	0.400	0.400	0.267	0.426	0.400	0.533	<b>0.600</b>
UW	0.131	0.497	0.759	0.544	0.881	0.903	0.868	0.869	0.903	0.684	0.916	0.897	0.894	0.859	0.728	<b>0.919</b>	0.833	0.897	0.881
Win	0	0	0	0	0	1	0	0	1	2	4	1	1	0	0	2	2	2	6
Tie	0	1	1	1	0	1	1	1	2	2	2	3	2	0	1	1	0	4	7
Lose	30	29	29	29	30	28	29	29	27	26	24	26	27	30	29	27	28	24	17
Best	0	1	1	1	0	2	1	1	3	4	6	4	3	0	1	3	2	6	13
AVG_rank	15.817	11.717	11.183	10.667	13.467	10.900	11.700	11.983	8.200	12.833	6.717	6.333	7.917	11.500	11.767	7.250	9.933	5.900	4.167

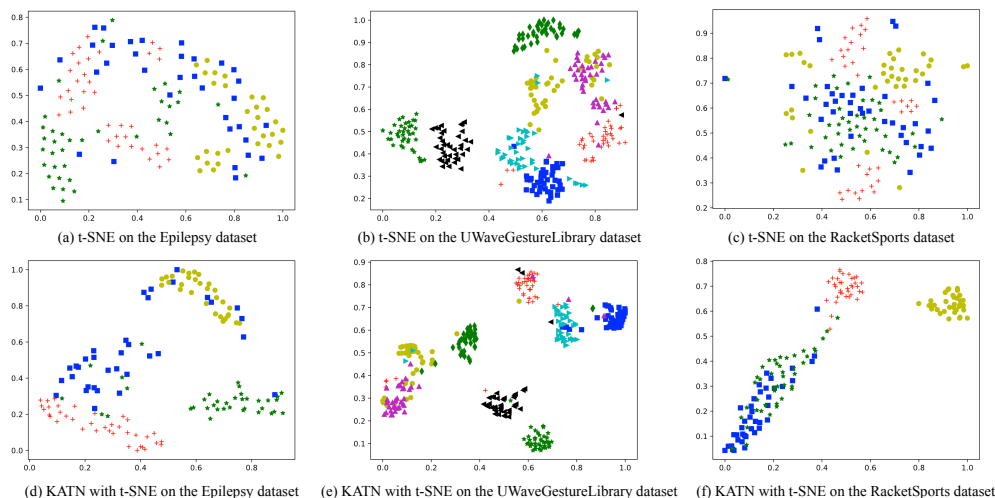


Fig. 5. Visualization of representations learned by t-SNE and KATN with t-SNE on the Epilepsy, UWaveGestureLibrary, and RacketSports datasets.

incurring greater complexity. This outcome reflects a deliberate architectural trade-off, where KATN prioritizes expressive temporal-spatial modeling while accepting a moderate increase in computational cost to enhance its modeling effectiveness.

The computational burden becomes more evident when assessing CPU-based inference efficiency. Across the 30 UEA benchmark testing datasets, all transformer variants, including KATN, exhibit substantially longer CPU inference times compared to the seven MTSC baseline models. For instance, the slowest MTSC baseline method, DKN, achieves an average inference time of 15.004s, whereas the fastest transformer

model, Linformer, still requires 83.859s. This pronounced disparity highlights the considerably higher deployment cost associated with transformer architectures, particularly in CPU-limited environments.

Focusing specifically on GPU inference time across the same testing datasets, the differences among the seven transformer variants become marginal. For instance, Linformer achieves the shortest GPU average inference time (3.766s), while DualAttnTrans incurs the highest (4.714s). This narrow variance suggests that GPU parallelism effectively mitigates the architectural overhead of transformers to a large extent.



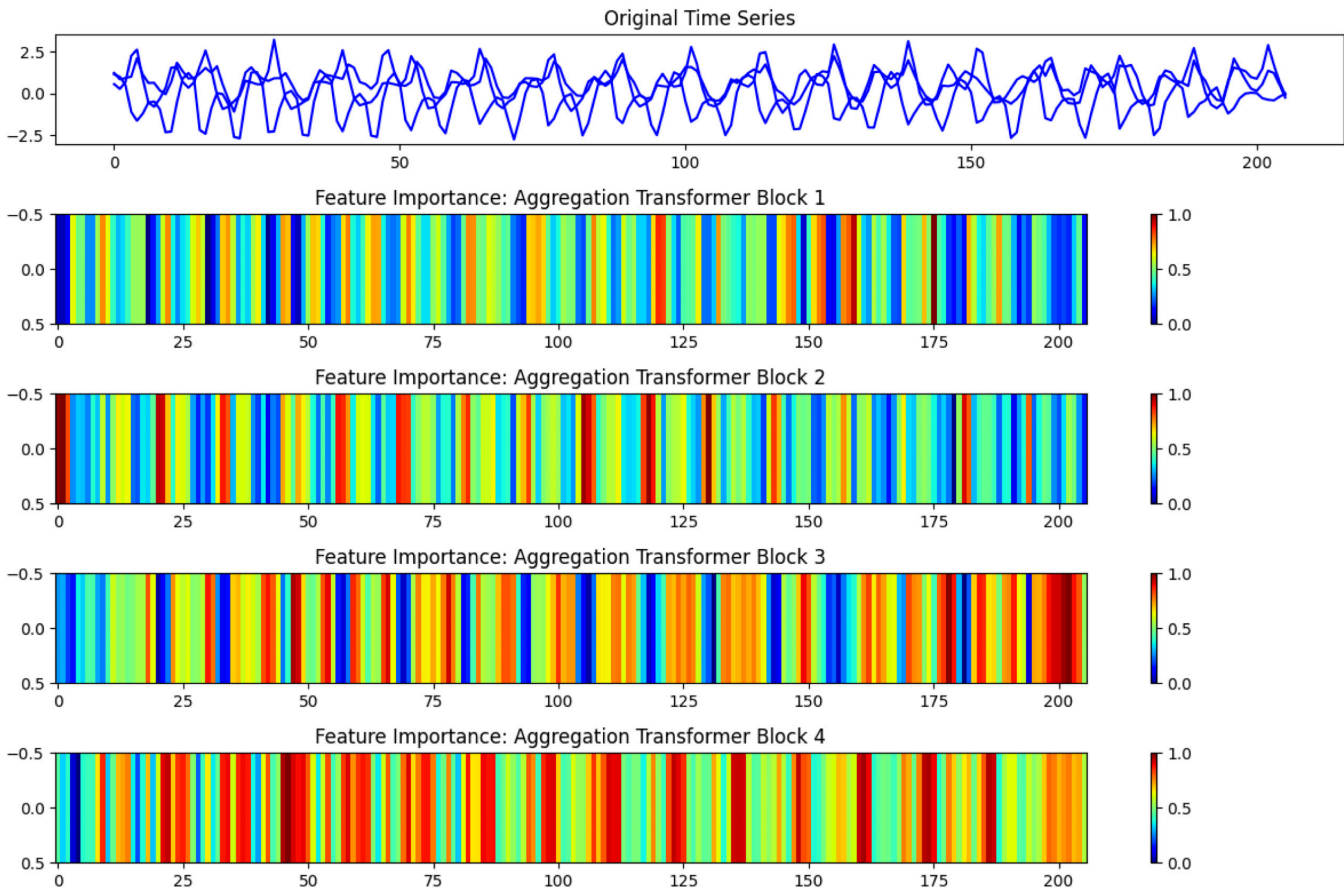


Fig. 7. Visualization of temporal receptive field expansion across aggregation transformer blocks in KATN on the Epilepsy dataset. This figure illustrates the hierarchical evolution of temporal activation patterns as input signals propagate through KATN’s aggregation transformer blocks. The color intensity represents the normalized variance of feature responses at each time step, where warmer hues indicate time regions of greater importance for seizure characterization. The effective activation widths for blocks 1 to 4 are measured as 81, 95, 137, and 143 time steps, respectively, under a threshold of 0.5. This gradual increase reflects the model’s ability to integrate information over progressively extended temporal ranges. By capturing both fine-grained local variations and long-range dependencies, KATN enables more comprehensive temporal reasoning, which is essential for accurately identifying complex seizure dynamics in multivariate time series.

perform a two-dimensional projection of its high-dimensional latent features using t-distributed stochastic neighbor embedding (t-SNE) [72], a nonlinear manifold learning method well-suited for visualizing complex feature spaces. As shown in Fig. 5, this projection enables a comparative analysis of embeddings across the Epilepsy, UWaveGestureLibrary, and RacketSports datasets. Notably, the embeddings produced by KATN exhibit more distinct and tightly grouped clusters than those derived from t-SNE applied directly to raw features. This improvement is particularly salient in the UWaveGestureLibrary dataset (Figs. 5(b) and 5(e)), where intra-class cohesion and inter-class boundaries are clearly delineated—reflecting the model’s ability to extract semantically aligned, task-relevant representations.

To complement this analysis, we present the input and output heatmaps in Fig. 6, which illustrate how KATN hierarchically refines temporal features across successive layers. These visualizations reveal a selective enhancement of salient dynamics while diminishing the influence of spurious or redundant patterns, thereby improving representational com-

compactness and class separability.

Beyond these spatial representations, we further investigate KATN’s capacity for modeling long-range temporal dependencies—an essential trait for MTSC tasks involving delayed or distributed temporal cues, such as seizure recognition. Following the methodology in [73], we visualize the progressive expansion of temporal receptive fields across the model’s four aggregation transformer blocks on the Epilepsy dataset (Fig. 7). The heatmaps, constructed via normalized feature variance at each time step, highlight regions with heightened predictive relevance, with warmer tones indicating greater salience.

A clear hierarchical pattern emerges: as the depth of the network increases, so too does the temporal span of effective activation. Specifically, the receptive fields expand from 81 time steps in the first transformer block to 95, 137, and 143 in subsequent layers, using a fixed salience threshold of 0.5. This progressive broadening indicates a gradual shift from localized temporal processing toward integration of extended temporal contexts. The integration of multi-head self-attention with

MResNet facilitates adaptive focusing on temporally distant yet semantically correlated events. Meanwhile, the additive feature aggregation design enables the seamless fusion of temporal information across multiple scales.

Through this multi-stage refinement, KATN constructs a temporally expressive representation space that unifies short-term cues and long-term dependencies. Such hierarchical modeling is essential for uncovering complex, non-local patterns in real-world time series. The observed temporal activation dynamics thus provide compelling evidence of KATN's ability to internalize and operationalize long-range temporal structures, enhancing its utility for challenging tasks such as epileptic event detection.

## V. CONCLUSION

The proposed KATN incorporates four aggregation transformer blocks, each merging the local patterns extracted by MResNet with the global patterns derived from the multi-head attention network. These blocks are able to capture intricate long-range dependencies prevalent among multiple variables. Through extensive experiments, KATN demonstrates its superiority over 6 SOTA transformer algorithms, i.e., VanTrans, SwinTrans, Linformer, FMLA, DualAggTrans, and DualAttTrans, in terms of 'win'/'tie'/'lose' and AVG\_rank. Furthermore, when compared to 18 MTSC algorithms on 30 UEA datasets, KATN is winner of 13 datasets, attaining an AVG\_rank score of 4.167. These results signify KATN's promising potential in addressing various MTSC applications in real-world applications.

While KATN offers significant benefits, it also presents challenges, particularly due to its resource-intensive nature, which may hinder practical deployment. To mitigate this, we plan to incorporate network compression techniques, reducing computational complexity by eliminating redundant nodes. Additionally, we aim to improve the model's robustness through automatic hyper-parameter optimization, exploring methods such as grid search, random search, and Bayesian optimization to identify optimal configurations for diverse datasets and applications. Furthermore, we intend to leverage neural architecture search to optimize both hyper-parameters and the network architecture in an integrated manner, ensuring enhanced performance and adaptability.

## REFERENCES

- [1] C. Chen, L. Li, M. Beetz, A. Banerjee, R. Gupta, and V. Grau, "Large language model-informed ecg dual attention network for heart failure risk prediction," *IEEE Trans. Big Data*, pp. 1–13, 2025.
- [2] A. Hussain, T. Hussain, W. Ullah, S. U. Khan, M. J. Kim, K. Muhammad, J. D. Ser, and S. W. Baik, "Big data analysis for industrial activity recognition using attention-inspired sequential temporal convolution network," *IEEE Trans. Big Data*, pp. 1–12, 2024.
- [3] Y. Shi, X. Zhang, Y. Shang, and N. Yu, "Don't be misled by emotion! disentangle emotions and semantics for cross-language and cross-domain rumor detection," *IEEE Trans. Big Data*, vol. 10, no. 3, pp. 249–259, 2024.
- [4] F. Ding, B. Li, X. Ben, J. Zhao, and H. Zhou, "Alad: A new unsupervised time series anomaly detection paradigm based on activation learning," *IEEE Trans. Big Data*, pp. 1–13, 2024.
- [5] A. Chen, X. Zhou, Y. Fan, and H. Chen, "Anomaly detection in multi-level model space," *IEEE Trans. Big Data*, pp. 1–12, 2025.
- [6] H. I. Fawaz, G. Forestier, J. Weber, L. Idoumghar, and P.-A. Muller, "Deep learning for time series classification: a review," *Data Min. Knowl. Disc.*, vol. 33, pp. 917–963, 2019.
- [7] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, pp. 436–444, 2015.
- [8] H. Xing, Z. Xiao, R. Qu, Z. Zhu, and B. Zhao, "An efficient federated distillation learning system for multitask time series classification," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–12, 2022.
- [9] F. He, T.-Y. Fu, and W.-C. Lee, "Rel-cnn: Learning relationship features in time series for classification," *IEEE Trans. Knowl. Data En.*, vol. 35, no. 7, pp. 7412–7426, 2023.
- [10] Q. Ma, S. Li, and G. W. Cottrell, "Adversarial joint-learning recurrent neural network for incomplete time series classification," *IEEE Trans. Pattern Anal.*, vol. 44, no. 4, pp. 1765–1776, 2022.
- [11] R. Chen, X. Yan, S. Wang, and G. Xiao, "Da-net: Dual-attention network for multivariate time series classification," *Inf. Sci.*, vol. 610, pp. 472–487, 2022.
- [12] G. He, L. Dai, Z. Yu, and C. L. P. Chen, "Gan-based temporal association rule mining on multivariate time series data," *IEEE Trans. Knowl. Data En.*, pp. 1–13, 2023.
- [13] Z. Xiao, X. Xu, H. Xing, S. Luo, P. Dai, and D. Zhan, "Rtfn: A robust temporal feature network for time series classification," *Inf. Sci.*, vol. 571, pp. 65–86, 2021.
- [14] F. Karim, S. Majumdar, H. Darabi, and S. Harford, "Multivariate lstm-fcns for time series classification," *Neural Networks*, vol. 116, pp. 237–245, 2019.
- [15] S. Huang, L. Xu, and C. Jiang, "Artificial intelligence and advanced time series classification: Residual attention net for cross-domain modeling," *Fintech with Artificial Intelligence, Big Data, and Blockchain, Blockchain Technologies*, 2021.
- [16] H. Xing, Z. Xiao, D. Zhan, S. Luo, P. Dai, and K. Li, "Selfmatch: Robust semisupervised time-series classification with self-distillation," *Int. J. Intell. Syst.*, pp. 1–28, 2022.
- [17] Z. Xiao, H. Xing, R. Qu, L. Feng, S. Luo, P. Dai, B. Zhao, and Y. Dai, "Densely knowledge-aware network for multivariate time series classification," *IEEE Trans. Syst. Man Cy-S.*, pp. 1–13, 2024.
- [18] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Trans. Pattern Anal.*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [19] P. Munro, "Backpropagation," *Sammur, C., Webb, G. I. (eds) Encyclopedia of Machine Learning*, 2011.
- [20] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit *et al.*, "Mlp-mixer: An all-mlp architecture for vision," in *Proc. NeurIPS*, pp. 24 261–24 272, 2021.
- [21] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proc. CVPR*, 2019, pp. 3146–3154.
- [22] L. Chen, H. Zhang, J. Xiao, L. Nie, J. Shao, W. Liu, and T.-S. Chua, "Sea-cnn: Spatial and channel-wise attention in convolutional networks for image captioning," in *Proc. CVPR*, 2017, pp. 5659–5667.
- [23] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [24] Q. Chen, Q. Wu, J. Wang, Q. Hu, T. Hu, E. Ding, J. Cheng, and J. Wang, "Mixformer: Mixing features across windows and dimensions," in *Proc. CVPR*, 2022, pp. 5249–5259.
- [25] Z. Chen, Y. Zhang, J. Gu, L. Kong, X. Yang, and F. Yu, "Dual aggregation transformer for image super-resolution," in *Proc. CVPR*, 2023, pp. 12 312–12 321.
- [26] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelus)," *arXiv preprint arXiv:1606.08415*, 2016.
- [27] A. Bagnall, J. Lines, A. Bostrom, J. Large, and E. Keogh, "The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances," *Data Min. Knowl. Disc.*, vol. 31, pp. 1–55, 2017.
- [28] J. Lines and A. Bagnall, "Time series classification with ensembles of elastic distance measures," *Data Min. Knowl. Disc.*, vol. 29, pp. 565–592, 2015.
- [29] A. Bagnall, J. Lines, J. Hills, and A. Bostrom, "Time series classification with cote: the collective of transformation-based ensembles," in *Proc. ICDE 2016*, pp. 1548–1549, 2016.
- [30] J. Lines, S. Taylor, and A. Bagnall, "Time series classification with hive-cote: the hierarchical of transformation-based ensembles," *ACM Trans. Knowl. Discov. D.*, vol. 21, no. 52, pp. 1–35, 2018.

- [31] M. Middlehurst, J. Large, M. Flynn, J. Lines, A. Bostrom, and A. Bagnall, "Hive-cote 2.0: a new meta ensemble for time series classification," *Machine Learning*, vol. 110, pp. 3211–3243, 2021.
- [32] K. Fauvel, É. Fromont, V. Masson, P. Faverdin, and A. Termier, "Xem: An explainable-by-design ensemble method for multivariate time series classification," *Data Min. Knowl. Disc.*, vol. 36, pp. 917–957, 2022.
- [33] G. Li, B. Choi, J. Xu, S. S. Bhowmick, K.-P. Chun, and G. L.-H. Wong, "Efficient shapelet discovery for time series classification," *IEEE Trans. Knowl. Data En.*, vol. 34, no. 3, pp. 1149–1163, 2022.
- [34] F. J. Baldán and J. M. Benítez, "Multivariate times series classification through an interpretable representation," *Inf. Sci.*, vol. 569, pp. 596–614, 2021.
- [35] A. Shifaz, C. Pelletier, F. Petitjean, and G. I. Webb, "Ts-chief: A scalable and accurate forest algorithm for time series classification," *Data Min. Knowl. Disc.*, vol. 34, pp. 742–775, 2020.
- [36] M. G. Baydogan and G. Runger, "Time series representation and similarity based on local auto patterns," *Data Min. Knowl. Disc.*, vol. 30, pp. 476–509, 2016.
- [37] W. Pei, H. Dibeklioglu, D. M. J. Tax, and L. van der Maaten, "Multivariate time-series classification using the hidden-unit logistic model," *IEEE Trans. Neur. Net. Lear.*, vol. 29, no. 4, pp. 920–931, 2018.
- [38] H. Deng, G. Runger, E. Tuv, and M. Vladimir, "A time series forest for classification and feature extraction," *Inf. Sci.*, vol. 239, pp. 142–153, 2013.
- [39] G. He, X. Xin, R. Peng, M. Han, J. Wang, and X. Wu, "Online rule-based classifier learning on dynamic unlabeled multivariate time series data," *IEEE Trans. Syst. Man Cy-S.*, vol. 52, no. 2, pp. 1121–1134, 2022.
- [40] K. S. Tuncel and M. G. Baydogan, "Autoregressive forests for multivariate time series modeling," *Pattern Recogn.*, vol. 73, pp. 202–215, 2018.
- [41] J. Large, A. Bagnall, S. Malinowski, and R. Tavenard, "From bop to boss and beyond: time series classification with dictionary based classifier," *arXiv preprint arXiv:1809.06751*, 2018.
- [42] P. Schäfer and U. Leser, "Multivariate time series classification with weasel+muse," *arXiv preprint arXiv:1711.11343*, 2017.
- [43] K. Wu, K. Yuan, Y. Teng, J. Liu, and L. Jiao, "Broad fuzzy cognitive map systems for time series classification," *App. Soft Comput.*, vol. 128, pp. 1–13, 2022.
- [44] B. Zhao, H. Xing, X. Wang, F. Song, and Z. Xiao, "Rethinking attention mechanism in time series classification," *Inf. Sci.*, vol. 627, pp. 97–114, 2023.
- [45] H. I. Fawaz, B. Lucas, G. Forestier, C. Pelletier, D. Schmidt, J. Weber, G. I. Webb, L. Idoumghar, P.-A. Muller, and F. Petitjean, "Inceptiontime: finding alexnet for time series classification," *Data Min. Knowl. Disc.*, vol. 34, pp. 1936–1962, 2020.
- [46] D. Lee, S. Lee, and H. Yu, "Learnable dynamic temporal pooling for time series classification," *In Proc. AAAI*, vol. 35, no. 9, pp. 8288–8296, 2021.
- [47] H. I. Fawaz, G. Forestier *et al.*, "Time series classification from scratch with deep neural networks: A strong baseline," *In Proc. IEEE IJCNN 2017*, pp. 1578–1585, 2017.
- [48] Z. Xiao, X. Xu, H. Zhang, and E. Szczerbicki, "A new multi-process collaborative architecture for time series classification," *Knowl.-Based Syst.*, vol. 220, pp. 1–11, 2021.
- [49] G. Li *et al.*, "Shapenet: A shapelet-neural network approach for multivariate time series classification," *In Proc. AAAI*, vol. 35, pp. 8375–8383, 2021.
- [50] Z. Xiao, H. Xing, B. Zhao, R. Qu, S. Luo, P. Dai, K. Li, and Z. Zhu, "Deep contrastive representation learning with self-distillation," *IEEE Trans. Emerg. Top. Comput. Intell.*, pp. 1–13, 2023.
- [51] F. M. Bianchi, S. Scardapane, S. Løkse, and R. Jenssen, "Reservoir computing approaches for representation and classification of multivariate time series," *IEEE Trans. Neur. Net. Learn.*, vol. 32, no. 5, pp. 2169–2179, 2021.
- [52] A. Dempster, F. Petitjean, and G. I. Webb, "Rocket: exceptionally fast and accurate time series classification using random convolutional kernels," *Data Min. Knowl. Disc.*, vol. 34, pp. 1454–1495, 2020.
- [53] L. Sun, C. Li, Y. Ren, and Y. Zhang, "A multitask dynamic graph attention autoencoder for imbalanced multilabel time series classification," *IEEE Trans. Neur. Net. Lear.*, pp. 1–14, 2024.
- [54] Y. Cheng, D. Lee, H. Oberhauser, and H. Li, "Generalized time series classification via component decomposition and alignment," *IEEE Trans. Big Data*, pp. 1–16, 2025.
- [55] Z. Huang, C. Yang, X. Chen, X. Zhou, G. Chen, T. Huang, and W. Gui, "Functional deep echo state network improved by a bi-level optimization approach for multivariate time series classification," *App. Soft Comput.*, vol. 106, pp. 1–12, 2021.
- [56] A. Dempster, D. F. Schmidt, and G. I. Webb, "Minirocket: A very fast (almost) deterministic transform for time series classification," *In Proc. KDD'21*, pp. 248–257, 2021.
- [57] Z. Xiao, X. Xu, H. Xing, R. Qu, F. Song, and B. Zhao, "Rnts: Robust neural temporal search for time series classification," *In Proc. IJCNN 2021*, pp. 1–8, 2021.
- [58] Z. Xiao, X. Xu, H. Xing, B. Zhao, X. Wang, F. Song, R. Qu, and L. Feng, "Dtcn: Deep transformer capsule mutual distillation for multivariate time series classification," *IEEE Trans. Cogn. Dev. Syst.*, pp. 1–17, 2024.
- [59] X. Zhang, Y. Gao, J. Lin, and C.-T. Lu, "Tapnet: Multivariate time series classification with attentional prototypical network," *In Proc. AAAI 2020*, pp. 6845–6852, 2020.
- [60] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler, "Efficient transformers: A survey," *ACM Comput. Surv.*, vol. 55, no. 6, 2022.
- [61] A. Vaswani, N. Shazeer, J. U. N. Parmar, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *In Proc. NeurIPS 2017*, pp. 5998–6008, 2017.
- [62] S. Wang, B. Z. Li, M. Khabza, H. Fang, and H. Ma, "Linformer: Self-attention with linear complexity," *arXiv preprint arXiv:2006.04768*, 2020.
- [63] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *In Proc. ICCV*, 2021, pp. 9992–10002.
- [64] M. Ding, B. Xiao, N. Codella, P. Luo, J. Wang, and L. Yuan, "Davvit: Dual attention vision transformers," *In Proc. ECCV*. Springer, 2022, pp. 74–92.
- [65] Y. Li, X. Peng, J. Zhang, Z. Li, and M. Wen, "Dct-gan: Dilated convolutional transformer-based gan for time series anomaly detection," *IEEE Trans. Knowl. Data En.*, vol. 35, no. 4, pp. 3632–3644, 2023.
- [66] Z. Lin, S. Zang, R. Wang, Z. Sun, J. Senthilnath, C. Xu, and C. K. Kwok, "Attention over self-attention: Intention-aware re-ranking with dynamic transformer encoders for recommendation," *IEEE Trans. Knowl. Data En.*, vol. 35, no. 8, pp. 7782–7795, 2023.
- [67] Y. Liu, S. Pan, Y. G. Wang, F. Xiong, L. Wang, Q. Chen, and V. C. Lee, "Anomaly detection in dynamic graphs via transformer," *IEEE Trans. Knowl. Data En.*, vol. 35, no. 12, pp. 12081–12094, 2023.
- [68] L. Xia, C. Huang, Y. Xu, and J. Pei, "Multi-behavior sequential recommendation with temporal graph transformer," *IEEE Trans. Knowl. Data En.*, vol. 35, no. 6, pp. 6099–6112, 2023.
- [69] A. Bagnall, H. Dau, J. Lines, M. Flynn, J. Large, A. Bostrom, P. Southam, and E. Keogh, "The uea multivariate time series classification archive, 2018," *arXiv preprint arXiv:1811.00075*, 2018.
- [70] C. Liu, J. Zhen, and W. Shan, "Time series classification based on convolutional network with a gated linear units kernel," *Eng. Appl. Artif. Intel.*, vol. 123, pp. 1–11, 2023.
- [71] X. Hu, L. Chu, J. Pei, W. Liu, and J. Bian, "Model complexity of deep learning: A survey," *Knowl. Inf. Syst.*, vol. 63, pp. 2585–2619, 2021.
- [72] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne," *J. Mach. Learn. Res.*, vol. 9, no. 11, 2008.
- [73] A. A. Ismail, M. Gunady, H. C. Bravo, and S. Feizi, "Benchmarking deep learning interpretability in time series predictions," in *Proc. the 34th International Conference on Neural Information Processing Systems*, ser. NIPS '20, 2020.



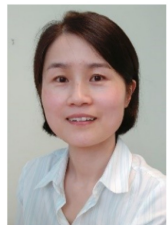
**Zhiwen Xiao** (Member, IEEE) received the B.Eng. degree in network engineering from the Chengdu University of Information Technology, Chengdu, China, in 2019, and the M.Eng. degree in computer science from the Northwest A&F University, Yangling, China, in 2023. He is currently pursuing the Ph.D. degree in computer science with Southwest Jiaotong University, Chengdu. His research interests include data mining, time series analysis, federated learning, representation learning, semantic communication, and computer vision.





**Huanlai Xing** (Member, IEEE) received Ph.D. degree in computer science from University of Nottingham, Nottingham, U.K., in 2013. He was a Visiting Scholar in Computer Science, The University of Rhode Island, USA. Supervisor: Dr. Haibo He (Robert Haas Endowed Chair Professor, IEEE Fellow, <https://www.ele.uri.edu/faculty/he/>) in 2020-2021. Huanlai Xing is an Associate Professor and PhD Supervisor with the School of Computing and Artificial Intelligence, Southwest Jiaotong University. He is on Editorial Board (Young Scientists

Committee) of SCIENCE CHINA INFORMATION SCIENCES. He was a member of several international conference program and senior program committees, such as IJCAI, ECML-PKDD, MobiMedia, ISCIT, ICCV, IJCNN, TrustCom, and ICSINC. His research interests include semantic communication, representation learning, data mining, reinforcement learning, machine learning, network function virtualization, and software defined networking.



**Rong Qu** (Senior Member, IEEE) is a full Professor at the School of Computer Science, University of Nottingham. She received her B.Sc. in Computer Science and Its Applications from Xidian University, China in 1996 and Ph.D. in Computer Science from The University of Nottingham, U.K. in 2003. Her research interests include the modeling and optimization for logistics transport scheduling, personnel scheduling, network routing, portfolio optimization and timetabling problems by using evolutionary algorithms, mathematical programming, constraint

programming in operational research and artificial intelligence. These computational techniques are integrated with knowledge discovery, machine learning and data mining to provide intelligent decision support on logistic fleet operations at SMEs, workforce scheduling at hospitals, policy making in education, and cyber security for connected and autonomous vehicles.

Dr. Qu is an associated editor at Engineering Applications of Artificial Intelligence, IEEE Computational Intelligence Magazine, IEEE Transactions on Evolutionary Computation, Journal of Operational Research Society, and PeerJ Computer Science. She is a Senior IEEE Member since 2012 and the Vice-Chair of Evolutionary Computation Task Committee since 2019 and Technical Committee on Intelligent Systems Applications (2015-2018) at IEEE Computational Intelligence Society. She has guest edited special issues on the automated design of search algorithms and machine learning at the IEEE Transactions on Pattern Analysis and Machine Intelligence and IEEE Computational Intelligence Magazine.



**Hui Li** received the B.Sc. and M.Sc. degrees in applied mathematics from the School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an, China, in 1999 and 2002, respectively, and the Ph.D. degree in computer science from the University of Essex, Colchester, U.K., in 2008. He is currently a Professor with the School of Mathematics and Statistics, Xi'an Jiaotong University. His current research interests include evolutionary computation, multiobjective optimization, and machine learning.

Dr. Li was a recipient of the 2010 IEEE TRANSACTIONS ON EVOLUTIONARY COMPUTATION Outstanding Paper Award as one of the inventors for MOEA/D.



**Huagang Tong**, received his Ph.D. degree in management sciences from Nanjing University of Aeronautics and Astronautics in 2022. He is an Associate Professor with College of Economic and Management, Nanjing Tech University, Puzhu Road(S), Nanjing 211816, China. His research interests are knowledge discovery, operation research, and data mining.



**Shouxi Luo** (Member, IEEE) received the bachelor's degree in communication engineering and the Ph.D. degree in communication and information systems from the University of Electronic Science and Technology of China, Chengdu, China, in 2011 and 2016, respectively. He is an Associate Professor with the School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu. His research interests include data center networks, software-defined networking, and networked systems.



**Song Jing** is a Senior Engineer with the Informatization Research Institute, Southwest Jiaotong University, and Deputy Director of the Intelligent Network Laboratory, National Engineering Research Center for High-Speed Railway and Urban Rail Transit. He was selected into Tianfu Ten-thousand Talents Program. He have involved in various projects and engineering tasks such as the National Railway Group's Railway 5G Network Security and Passenger Dedicated Line Passenger Service System Security Assurance Platform. He focuses on cyber-

security in cyberspace and information communication networks.



**Li Feng** received his PhD degree from Xi'an Jiaotong University, in 2005, under the supervision of Prof. Xiaohong Guan (Academian of CAS, IEEE Fellow). He is a Research Professor and PhD supervisor with the School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu. His research interests include artificial intelligence, cyber security and its applications.



**Qian Wan** received the B.Eng. degree in computer science from the Wuhan University, Wuhan, China, and the M.Des. degree in interaction design from the China University of Geosciences, Wuhan, China. He is pursuing the Ph.D. degree in computer science at Southwest Jiaotong University, Chengdu, China. His research interests include data model and mining, virtual reality, and augmented reality.