
Transparent Networks for Multivariate Time Series

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transparent models, which are machine learning models that produce inherently
2 interpretable predictions, are receiving significant attention in high-stakes domains.
3 However, despite much real-world data being collected as time series, there is a lack
4 of studies on transparent time series models. To address this gap, we propose a novel
5 transparent neural network model for time series called Generalized Additive Time
6 Series Model (GATSM). GATSM consists of two parts: 1) independent feature
7 networks to learn feature representations, and 2) a transparent temporal module to
8 learn temporal patterns across different time steps using the feature representations.
9 This structure allows GATSM to effectively capture temporal patterns and handle
10 dynamic-length time series while preserving transparency. Empirical experiments
11 show that GATSM significantly outperforms existing generalized additive models
12 and achieves comparable performance to black-box time series models, such as
13 recurrent neural networks and Transformer. In addition, we demonstrate that
14 GATSM finds interesting patterns in time series. The source code is available at
15 <https://anonymous.4open.science/r/GATSM-78F4/>.

16 1 Introduction

17 Artificial neural networks excel at learning complex representations and demonstrate remarkable
18 predictive performance across various fields. However, their complexity makes interpreting the
19 decision-making processes of neural network models challenging. Consequently, post-hoc explainable
20 artificial intelligence (XAI) methods, which explain the predictions of trained black-box models,
21 have been widely studied in recent years [1, 2, 3, 4]. XAI methods are generally effective at
22 providing humans with understandable explanations of model predictions. However, they may
23 produce incorrect and unfaithful explanations of the underlying black-box model and cannot provide
24 actual contributions of input features to model predictions [5, 6]. Therefore, their applicability to
25 high-stakes domains—such as healthcare and fraud detection, where faithfulness to the underlying
26 model and actual contributions of features are important—is limited.

27 Due to these limitations, transparent (i.e., inherently interpretable) models are attracting attention as
28 alternatives to XAI in high-stakes domains [7, 8, 9]. Modern transparent models typically adhere to
29 the *generalized additive model* (GAM) framework [10]. A GAM consists of independent functions,
30 each corresponding to an input feature, and makes predictions as a linear combination of these
31 functions (e.g., the sum of all functions). Therefore, each function reflects the contribution of its
32 respective feature. For this reason, interpreting GAMs is straightforward, making them widely used in
33 various fields, such as healthcare [11, 12], survival analysis [13], and model bias discovery [7, 14, 15].
34 However, despite much real-world data being collected as time series, research on GAMs for time
35 series remains scarce. Consequently, the applicability of GAMs in real-world scenarios is still limited.

36 To overcome this limitation, we propose a novel transparent model for multivariate time series
37 called Generalized Additive Time Series Model (GATSM). GATSM consists of independent feature
38 networks to learn feature representations and a transparent temporal module to learn temporal patterns.

39 Since employing distinct networks across different time steps requires a massive amount of learnable
40 parameters, the feature networks in GATSM share the weights across all time steps, while the
41 temporal module independently learns temporal patterns. GATSM then generates final predictions by
42 integrating the feature representations with the temporal information from the temporal module. This
43 strategy allows GATSM to effectively capture temporal patterns and handle dynamic-length time
44 series while preserving transparency. Additionally, this approach facilitates the separate extraction of
45 time-independent feature contributions, the importance of individual time steps, and time-dependent
46 feature contributions through the feature functions, temporal module, and final prediction. To
47 demonstrate the effectiveness of GATSM, we conducted empirical experiments on various time series
48 datasets. The experimental results show that GATSM significantly outperforms existing GAMs
49 and achieves comparable performances to black-box time series models, such as recurrent neural
50 networks and Transformer [16]. In addition, we provide visualizations of GATSM’s predictions to
51 demonstrate that GATSM finds interesting patterns in time series.

52 2 Related Works

53 Various XAI studies have been conducted over the past decade [7, 8, 9, 17, 18]; however, they are
54 less relevant to the transparent model that is the subject of this study. Therefore, we refer readers to
55 [19, 20] for more detailed information on recent XAI research. In this section, we review existing
56 transparent models closely related to our GATSM and discuss their limitations.

Table 1: Advantages of GATSM.

	Time series input	Temporal pattern	Dynamic time series
existing GAMs			
NATM	✓		
GATSM (our)	✓	✓	✓

57 The simple linear model is designed to fit the conditional expectation $g(\mathbb{E}(y | \mathbf{x})) = \sum_{i=1}^M x_i w_i$,
58 where $g(\cdot)$ is a link function, M indicates the number of input features, y is the target value for the
59 given input features $\mathbf{x} \in \mathbb{R}^M$, and $w_i \in \mathbb{R}$ is the learnable weight for x_i . This model captures only
60 linear relationships between the target y and the inputs \mathbf{x} . To address this limitation, GAM [10]
61 extends the simple linear model to the generalized form as follows:

$$g(\mathbb{E}(y | \mathbf{x})) = \sum_{i=1}^M f_i(x_i), \quad (1)$$

62 where each $f_i(\cdot)$ is a function that models the effect of a single feature, referred as a feature function.
63 Typically, $f_i(\cdot)$ becomes a non-linear function such as a decision tree or neural network to capture
64 non-linear relationships.

65 Originally, GAMs were fitted via the backfitting algorithm using smooth splines [10, 21]. Later, Yin
66 Lou et al. [22] and Harsha Nori et al. [23] have proposed boosted decision tree-based GAMs, which
67 use boosted decision trees as feature functions. Spline- and tree-based GAMs have less flexibility
68 and scalability. Thus, extending them to transfer or multi-task learning is challenging. To overcome
69 this problem, various neural network-based GAMs have been proposed in recent years. Potts [24]
70 introduced generalized additive neural network, which employs 2-layer neural networks as feature
71 functions. Similarly, Rishabh Agarwal et al. [7] proposed neural additive model (NAM) that employs
72 multi-layer neural networks. To improve the scalability of NAM, Chun-Hao Chang et al. [8] and
73 Filip Radenovic et al. [9] proposed the neural oblivious tree-based GAM and the basis network-based
74 GAM, respectively. Xu et al. [25] introduced a sparse version of NAM using the group LASSO. One
75 disadvantage of GAMs is their limited predictive power, which stems from the fact that they only
76 learn first-order feature interactions-i.e., relationships between the target value and individual features.
77 To address this, various studies have been conducted to enhance the predictive powers of GAMs by
78 incorporating higher-order feature interactions, while still maintaining transparency. GA²M [26]
79 simply takes pairwise features as inputs to learn pairwise interactions. GAMI-Net [27], a neural
80 network-based GAM, consists of networks for main effects (i.e., first-order interactions) and pairwise
81 interactions. To enhance the interpretability of GAMI-Net, the sparsity and heredity constraints are
82 added, and trivial features are pruned in the training process. Sparse interaction additive network [28]

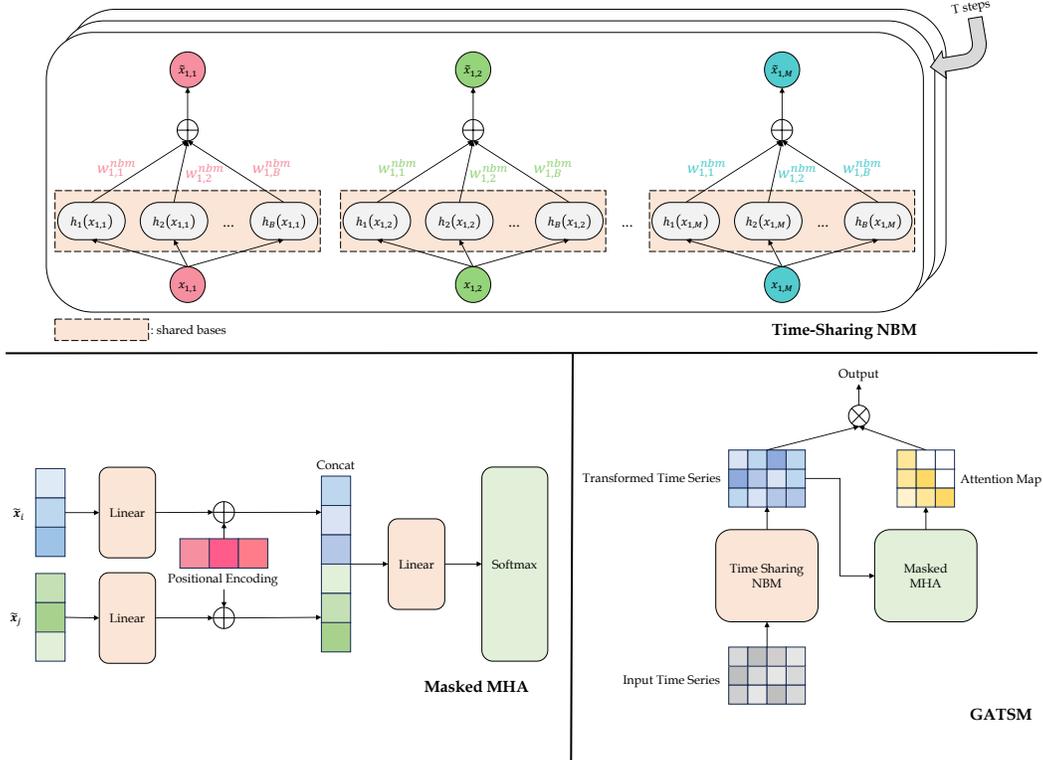


Figure 1: Architecture of GATSM.

83 is a 3-phase method for exploiting higher-order interactions. Initially, a black-box neural network is
 84 trained; subsequently, the top- k important features are identified using explainable feature attribution
 85 methods like LIME [1] and SHAP [2], and finally, NAM is trained with these extracted features.
 86 Dubey et al. [29] introduced scalable polynomial additive model, an end-to-end model that learns
 87 higher-order interactions via polynomials. Similarly, Kim et al. [15] proposed higher-order NAM that
 88 utilizes the feature crossing technique to capture higher-order interactions. Despite their capabilities,
 89 the aforementioned GAMs cannot process time series data, which limits their applicability in real-
 90 world scenarios. Recently, neural additive time series Model (NATM) [30], a time-series adaptation
 91 of NAM, has been proposed. However, NATM handles each time step independently with separate
 92 feature networks. This approach cannot capture effective temporal patterns and only takes fixed-length
 93 time series as input. Our GATSM not only captures temporal patterns but also handles dynamic-length
 94 time series. Table 1 shows the advantages of our GATSM compared to existing GAMs.

95 3 Problem Statement

96 We tackle the problem of the existing GAMs on time series. Equation (1) outlines the GAM framework
 97 for tabular data, which fails to capture the interactions between current and previous observations in
 98 time series. A straightforward method to extend GAM to time series, adopted in NATM, is applying
 99 distinct feature functions to each time step and summing them to produce predictions:

$$g(\mathbb{E}(y_t | \mathbf{X}_{:t})) = \sum_{i=1}^t \sum_{j=1}^M f_{i,j}(x_{i,j}), \quad (2)$$

100 where $\mathbf{X} \in \mathbb{R}^{T \times M}$ is a time series with T time steps and M features, and t is the current time step.
 101 This method can handle time series data as input but fails to capture effective temporal patterns
 102 because the function $f_{i,j}(\cdot)$ still does not interact with previous time steps. To overcome this problem,

103 we suggest a new form of GAM for time series defined as follows:

$$g(\mathbb{E}(y_t | \mathbf{X}_t)) = \sum_{i=1}^t \sum_{j=1}^M f_{i,j}(x_{i,j}, \mathbf{X}_t). \quad (3)$$

104 **Definition 3.1** *GAMs for time series, which capture temporal patterns hold the form of Equation 3.*

105 In Equation (3), the function $f(\cdot, \cdot)$ can capture interactions between current and previous time steps.
 106 Therefore, GAMs adhering to Definition 3.1 are capable of capturing temporal patterns. However,
 107 implementing such a model while maintaining transparency poses challenges. In the following
 108 section, we will describe our approach to implementing a GAM that holds Definition 3.1. To the best
 109 of our knowledge, no existing literature addresses Definition 3.1.

110 4 Our Method: Generalized Additive Time Series Model

111 4.1 Architecture

112 Figure 1 shows the overall architecture of GATSM. Our model has two modules: 1) feature networks,
 113 called time-sharing neural basis model, for learning feature representations, and 2) masked multi-head
 114 attention for learning temporal patterns.

115 **Time-Sharing NBM:** Assume a time series with T time steps and M features. Applying GAMs
 116 to this time series necessitates $T \times M$ feature functions, which becomes problematic when dealing
 117 with large T or M due to increased model size. This limits the applicability of GAMs to real-world
 118 datasets. To overcome this problem, we extend neural basis model (NBM) [9] to time series as:

$$\tilde{x}_{i,j} = f_j(x_{i,j}) = \sum_{k=1}^B h_k(x_{i,j}) w_{j,k}^{nbm}. \quad (4)$$

119 We refer to this extended version of NBM as time-sharing NBM. Time-sharing NBM has B basis
 120 functions, with each basis $h_k(\cdot)$ taking a feature $x_{i,j}$ as input. The feature-specific weight $w_{j,k}^{nbm}$
 121 then projects the basis to the transformed feature $\tilde{x}_{i,j}$. As depicted in Equation 4, the basis functions
 122 are shared across all features and time steps, drastically reducing the number of required feature
 123 functions $T \times M$ to B . We use $B = 100$ and implement $h_k(\cdot)$ using multi-layer perceptron (MLP).

124 **Masked MHA:** GATSM employs multi-head attention (MHA) to learn temporal patterns. Although
 125 the dot product attention [16] is popular, simple dot operation has low expressive power [31].
 126 Therefore, we adopt the 2-layer attention mechanism proposed by [31] to GATSM. We first transform
 127 $\tilde{\mathbf{x}}_i = [\tilde{x}_{i,1}, \tilde{x}_{i,2}, \dots, \tilde{x}_{i,M}] \in \mathbb{R}^M$ produced by Equation 4 as follows:

$$\mathbf{v}_i = \tilde{\mathbf{x}}_i^\top \mathbf{Z} + \mathbf{pe}_i, \quad (5)$$

128 where $\mathbf{Z} \in \mathbb{R}^{M \times D}$ is a learnable weight, $\mathbf{pe}_i = [pe_{i,1}, pe_{i,2}, \dots, pe_{i,D}] \in \mathbb{R}^D$ is the positional
 129 encoding for i -th step, and D indicates the hidden size. The positional encoding is defined as follows:

$$pe_{i,j} = \begin{cases} \sin\left(\frac{i}{10000^{2j/D}}\right) & \text{if } j \bmod 2 = 1, \\ \cos\left(\frac{i}{10000^{2j/D}}\right) & \text{otherwise.} \end{cases} \quad (6)$$

130 The positional encoding helps GATSM effectively capture temporal patterns. While learnable position
 131 embedding also works in GATSM, we recommend positional encoding because position embedding
 132 requires knowledge of the maximum number of time steps, which is often unknown in real-world
 133 settings. After computing \mathbf{v}_i , we calculate the attention scores as follows:

$$e_{k,i,j} = \sigma([\mathbf{v}_i | \mathbf{v}_j]^\top \mathbf{w}_k^{attn}) m_{i,j}, \quad (7)$$

$$a_{k,i,j} = \frac{\exp(e_{k,i,j})}{\sum_{t=1}^T \exp(e_{k,i,t})}, \quad (8)$$

134 where k is attention head index, $\sigma(\cdot)$ is an activation function, $\mathbf{w}_k^{attn} \in \mathbb{R}^{2D}$, and $m_{i,j} \in \mathbb{R}$ is the
 135 mask value used to block future information. The time mask is defined as follows:

$$m_{i,j} = \begin{cases} 1 & \text{if } i \leq j, \\ -\infty & \text{otherwise.} \end{cases} \quad (9)$$

136 **Inference:** The prediction of GATSM is produced by combining the transformed features from
 137 time-sharing NBM with the attention scores from masked MHA.

$$\hat{y}_t = \sum_{k=1}^K \mathbf{a}_{k,t}^\top \tilde{\mathbf{X}} \mathbf{w}_k^{out}, \quad (10)$$

138 where K is the number of attention heads, $\mathbf{a}_{k,t} = [a_{k,i,1}, a_{k,i,2}, \dots, a_{k,i,T}] \in \mathbb{R}^T$ is the attention
 139 map in Equation 8, $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_T] \in \mathbb{R}^{T \times M}$ is the transformed features in Equation 4, and
 140 $\mathbf{w}_k^{out} \in \mathbb{R}^M$ is the learnable output weight.

141 **Interpretability:** We can rewrite Equation 10 as the following scalar form:

$$\begin{aligned} \sum_{k=1}^K \mathbf{a}_{k,t}^\top \tilde{\mathbf{X}} \mathbf{w}_k^{out} &= \sum_{u=1}^t \sum_{m=1}^M \sum_{k=1}^K \sum_{b=1}^B a_{k,t,u} h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out} \\ &= \sum_{u=1}^t \sum_{m=1}^M f_{u,m}(x_{u,m}, \mathbf{X}; t) \end{aligned} \quad (11)$$

142 Equation 11 shows that GATSM satisfying Definition 3.1. We can derive three types of interpretations
 143 from GATSM: 1) $a_{k,t,u}$ indicates the importance of time step u at time step t , 2) $h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out}$
 144 represents the time-independent contribution of feature m , and 3) $a_{k,t,u} h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out}$ repre-
 145 sents the time-dependent contribution of feature m at time step t .

146 5 Experiments

147 5.1 Experimental Setup

148 **Datasets:** We conducted our experiments using eight publicly available real-world time series
 149 datasets. From the Monash repository [32], we sourced three datasets: Energy, Rainfall, and
 150 AirQuality. Another three datasets, Heartbeat, LSST, and NATOPS, were downloaded from the
 151 UCR repository [33]. The remaining two datasets, Mortality and Sepsis, were downloaded from
 152 the PhysioNet [34]. We perform ordinal encoding for categorical features and standardize features
 153 to have zero-mean and unit-variance. For forecasting tasks, target value y is also standardized to
 154 zero-mean and unit-variance. If the dataset contains missing values, we impute categorical features
 155 with their modes and numerical features with their means. The dataset is split into a 60%/20%/20%
 156 ratio for training, validation, and testing, respectively. Table 2 shows the statistics of the experimental
 157 datasets. Further details of the experimental datasets can be found in Appendix B.

Table 2: Dataset statistics.

Dataset	Task	Variable length	# of time series	Avg. length	# of features	# of classes
Energy	1-step FCST	No	137	24	24	-
Rainfall	1-step FCST	No	160,267	24	3	-
AirQuality	1-step FCST	No	16,966	24	9	-
Heartbeat	Binary	No	409	405	61	2
Mortality	Binary	Yes	12,000	49.861	41	2
Sepsis	Binary	Yes	40,336	38.482	40	2
LSST	Multi-class	No	4,925	36	6	14
NATOPS	Multi-class	No	360	51	24	6

FCST: forecasting

158 **Baselines:** We compare our GATSM with 12 baselines, which can be categorized into four groups: 1)
 159 Black-box tabular models include extreme gradient boosting (XGBoost) [35] and MLP. 2) Black-box
 160 time series models include simple recurrent neural network (RNN), gated recurrent unit (GRU), long
 161 short-term memory (LSTM), and Transformer [16]. 3) Transparent tabular models are simple linear
 162 model (Linear), explainable boosting machine (EBM) [23], NAM [7], NodeGAM [8], and NBM [9].
 163 4) NATM [30] is a transparent time series model.

164 **Implementation:** We implement XGBoost and EBM models using the `xgboost` and `interpretml`
 165 libraries, respectively. For NodeGAM, we employ the official implementation provided by its authors
 166 [8]. The remaining models are developed using PyTorch [36]. All models undergo hyperparameter

Table 3: Predictive performance comparison of various models.

Model Type	Model	Energy	Rainfall	AirQuality	Heartbeat	Mortality	Sepsis	LSST	NATOPS	Avg. Rank
Black-box Tabular Model	XGBoost	0.094 (±0.137)	0.002 (±0.002)	0.532 (±0.019)	0.679 (±0.094)	0.707 (±0.015)	0.816 (±0.007)	0.424 (±0.012)	0.200 (±0.049)	8.500 (±4.000)
	MLP	0.459 (±0.101)	0.011 (±0.004)	0.423 (±0.031)	0.654 (±0.082)	0.842 (±0.014)	0.786 (±0.007)	0.417 (±0.008)	0.211 (±0.065)	7.375 (±2.134)
Black-box Time Series Model	RNN	0.320 (±0.122)	0.068 (±0.020)	0.644 (±0.032)	0.661 (±0.078)	0.581 (±0.040)	0.782 (±0.009)	0.422 (±0.029)	0.592 (±0.110)	7.750 (±2.712)
	GRU	0.435 (±0.107)	0.089 (±0.034)	<u>0.701</u> (±0.018)	0.694 (±0.052)	0.818 (±0.014)	0.785 (±0.010)	<u>0.629</u> (±0.013)	0.931 (±0.045)	4.375 (±2.669)
	LSTM	0.359 (±0.112)	<u>0.090</u> (±0.031)	0.683 (±0.026)	0.648 (±0.042)	0.790 (±0.020)	0.779 (±0.008)	0.491 (±0.082)	0.908 (±0.035)	6.375 (±3.623)
	Transformer	0.263 (±0.263)	0.098 (±0.035)	0.711 (±0.027)	0.690 (±0.040)	0.844 (±0.019)	0.789 (±0.010)	0.679 (±0.019)	0.967 (±0.029)	<u>4.000</u> (±3.703)
Transparent Tabular Model	Linear	<u>0.482</u> (±0.112)	0.004 (±0.001)	0.241 (±0.019)	0.637 (±0.070)	0.838 (±0.017)	0.723 (±0.011)	0.311 (±0.010)	0.206 (±0.045)	10.125 (±3.871)
	EBM	-0.200 (±0.409)	0.004 (±0.001)	0.324 (±0.014)	0.666 (±0.056)	0.729 (±0.017)	0.802 (±0.011)	0.408 (±0.016)	0.164 (±0.053)	9.750 (±3.284)
	NAM	0.363 (±0.218)	0.006 (±0.002)	0.300 (±0.013)	0.645 (±0.026)	<u>0.853</u> (±0.014)	0.800 (±0.006)	0.400 (±0.011)	0.242 (±0.040)	7.875 (±3.643)
	NodeGAM	0.398 (±0.195)	0.006 (±0.002)	0.380 (±0.032)	0.681 (±0.046)	0.854 (±0.013)	<u>0.802</u> (±0.007)	0.400 (±0.028)	0.247 (±0.012)	6.375 (±3.623)
	NBM	0.330 (±0.251)	0.007 (±0.003)	0.301 (±0.012)	0.716 (±0.039)	0.852 (±0.014)	0.799 (±0.006)	0.388 (±0.014)	0.189 (±0.029)	7.875 (±3.603)
Transparent Time Series Model	NATM	0.304 (±0.122)	0.038 (±0.011)	0.548 (±0.028)	<u>0.724</u> (±0.043)	N/A	N/A	0.452 (±0.010)	0.878 (±0.058)	5.667 (±2.582)
	GATSM (ours)	0.493 (±0.173)	0.073 (±0.027)	0.583 (±0.026)	0.843 (±0.025)	<u>0.853</u> (±0.015)	0.797 (±0.007)	0.570 (±0.024)	<u>0.956</u> (±0.027)	3.125 (±1.808)

167 tuning via Optuna [37]. The pytorch-based models are optimized with the Adam with decoupled
 168 weight decay (AdamW) [38] optimizer on an NVIDIA A100 GPU. Model training is halted if the
 169 validation loss does not decrease over 20 epochs. We use mean squared error for the forecasting tasks,
 170 and for classification tasks, we use cross-entropy loss. Further details of the model implementations
 171 and hyper-parameters are provided in Appendix C.

172 **5.2 Comparison with baselines**

173 Table 3 shows the predictive performances of the experimental models. We report mean scores
 174 and standard deviations over five different random seeds. For the forecasting datasets, we evaluate
 175 R^2 scores. For the binary classification datasets, we assess the area under the receiver operating
 176 characteristic curve (AUROC). For the multi-class classification datasets, we measure accuracy. We
 177 highlight the best-performing model in **bold** and underline the second-best model. Since the tabular
 178 models cannot handle time series, they only take \mathbf{x}_t to produce y_t .

179 On the Energy and Heartbeat datasets, which are small in size, our GATSM demonstrates the best
 180 performance, indicating strong generalization ability. EBM, XGBoost, and Transformer struggle
 181 with overfitting on the Energy dataset. For the Mortality and Sepsis datasets, there is no significant
 182 performance difference between tabular and time series models, nor between black-box and trans-
 183 parent models. This suggests that these two healthcare datasets lack significant temporal patterns
 184 and feature interactions. It is likely that seasonal patterns are hard to detect in medical data, and
 185 the patient’s current condition already encapsulates previous conditions, making historical data less
 186 crucial. Since these datasets contain variable-length time series, the performance of NATM, which
 187 can only handle fixed-length time series, is not available. On the Rainfall, AirQuality, LSST, and
 188 NATOPS datasets, the time series models significantly outperform the tabular models, indicating
 189 that these datasets contain important temporal patterns that tabular models cannot capture. Addition-
 190 ally, the black-box models outperform the transparent models, suggesting that these datasets have
 191 higher-order feature interactions that transparent models cannot capture. Nevertheless, GATSM is the
 192 best model within the transparent model group and performs comparably to Transformer. Overall,
 193 GATSM achieved the best average rank in the experiments, followed by the Transformer, indicating
 194 GATSM’s superiority. Additional experiments on model throughput and an ablation study on the
 195 basis functions are presented in Appendix D.

Table 4: Ablation study on different feature functions.

Feature Function	Energy	Rainfall	AirQuality	Heartbeat	Mortality	Sepsis	LSST	NATOPS
Linear	0.283(± 0.277)	0.071(± 0.024)	0.563(± 0.019)	0.766(± 0.024)	0.832(± 0.015)	0.735(± 0.012)	0.398(± 0.030)	0.972 (± 0.020)
NAM	0.304(± 0.229)	0.068(± 0.021)	0.564(± 0.019)	0.838(± 0.032)	0.851(± 0.013)	0.801 (± 0.005)	0.553(± 0.023)	0.933(± 0.039)
NBM	0.493 (± 0.173)	0.073 (± 0.027)	0.583 (± 0.026)	0.843 (± 0.025)	0.853 (± 0.015)	0.797(± 0.007)	0.570 (± 0.024)	0.956(± 0.027)

Table 5: Ablation study on the temporal module.

Temporal Module	Energy	Rainfall	AirQuality	Heartbeat	Mortality	Sepsis	LSST	NATOPS
Base	0.452(± 0.087)	0.007(± 0.002)	0.299(± 0.012)	0.661(± 0.043)	0.854 (± 0.013)	0.798(± 0.008)	0.392(± 0.006)	0.192(± 0.027)
Base + PE	0.397(± 0.054)	0.007(± 0.003)	0.299(± 0.012)	0.681(± 0.068)	0.852(± 0.013)	0.799 (± 0.007)	0.385(± 0.027)	0.228(± 0.029)
Base + MHA	0.368(± 0.230)	0.048(± 0.017)	0.555(± 0.020)	0.821(± 0.044)	0.847(± 0.020)	0.779(± 0.033)	0.595 (± 0.013)	0.856(± 0.059)
Base + PE + MHA	0.493 (± 0.173)	0.073 (± 0.027)	0.583 (± 0.026)	0.843 (± 0.025)	0.853(± 0.015)	0.797(± 0.007)	0.570(± 0.024)	0.956 (± 0.027)

196 **5.3 Ablation study**

197 **Choice of feature function:** We evaluate the performance of GATSM by changing the feature
 198 functions using three models: Linear, NAM, and NBM. Table 4 presents the results of this experiment.
 199 The simple linear function performs poorly because it lacks the capability to capture non-linear
 200 relationships. In contrast, NAM, which can capture non-linearity, shows improved performance over
 201 the linear function. However, NBM stands out by achieving the best performance in six out of eight
 202 datasets. This indicates that the basis strategy of NBM is highly effective for time series data.

203 **Design of temporal module:** We evaluate the performance of GATSM by modifying the design of
 204 the temporal module. The results are presented in Table 5. GATSM without the temporal module
 205 (Base) fails to learn temporal patterns and shows poor performance in the experiment. GATSM with
 206 only positional encoding (Base + PE) also shows similar performance to the Base, indicating that
 207 positional encoding alone is insufficient for capturing effective temporal patterns. GATSM with only
 208 multi-head attention (Base + MHA) outperforms the previous two methods, demonstrating that the
 209 MHA mechanism is beneficial for capturing temporal patterns. Finally, our full GATSM (Base + PE +
 210 MHA) significantly outperforms the other methods, suggesting that the combination of PE and MHA
 211 creates a synergistic effect. Consistent with our previous findings in section 5.2, all four methods
 212 show similar performances on the Mortality and Sepsis datasets, which lack significant temporal
 213 patterns.

214 **5.4 Interpretation**

215 In this section, we visualize four interpretations of GATSM’s predictions on the AirQuality dataset.
 216 In addition, interpretations for the Rainfall dataset can be found in Appendix E.

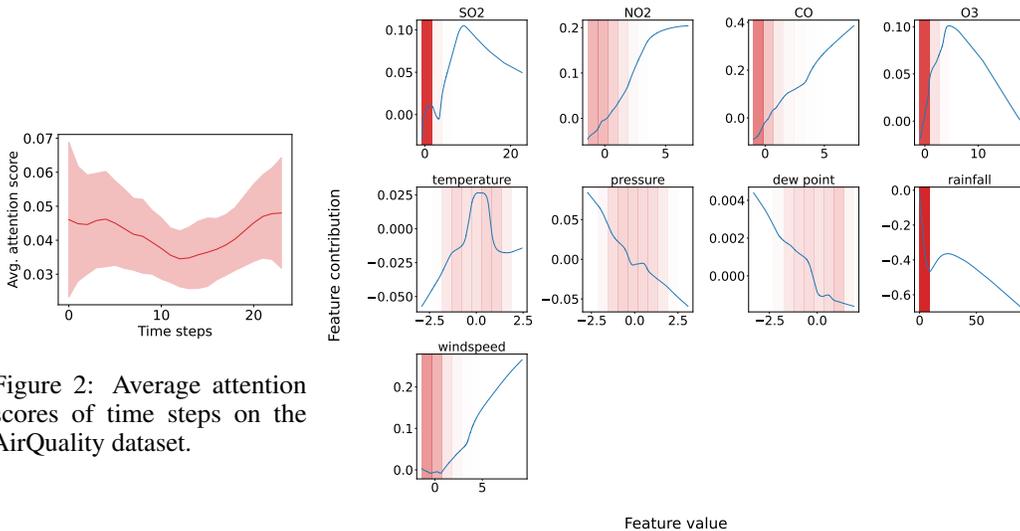


Figure 2: Average attention scores of time steps on the AirQuality dataset.

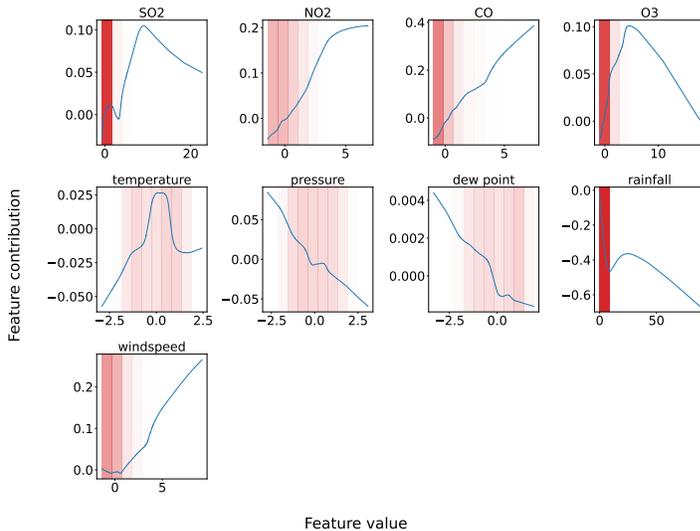


Figure 3: Global interpretations of features in the Air Quality dataset.

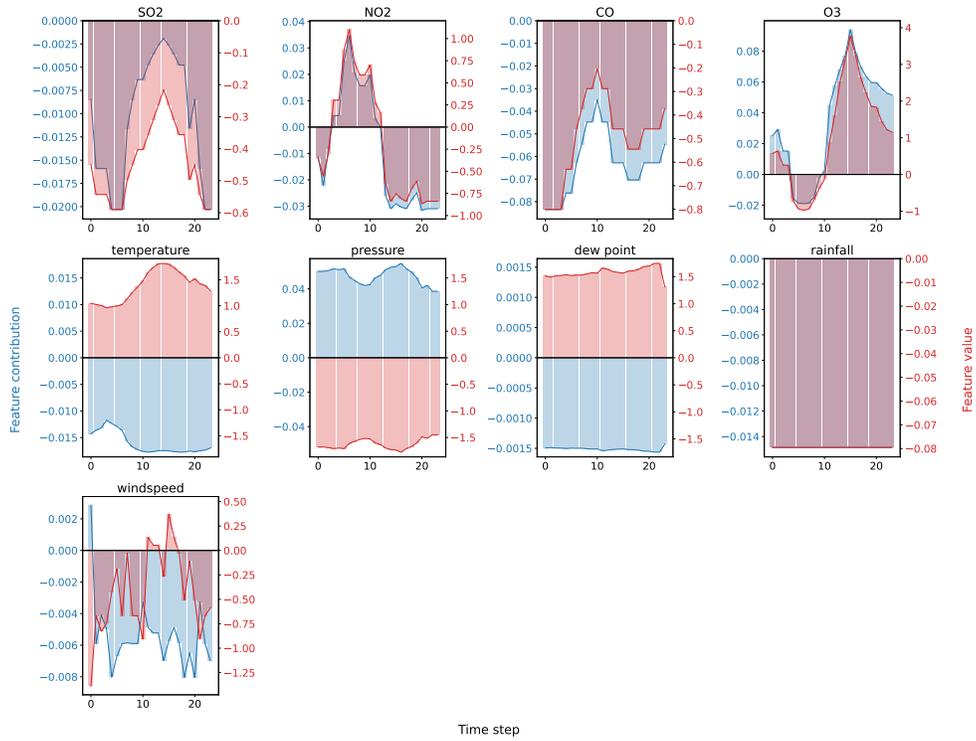


Figure 4: Local time-independent feature contributions.

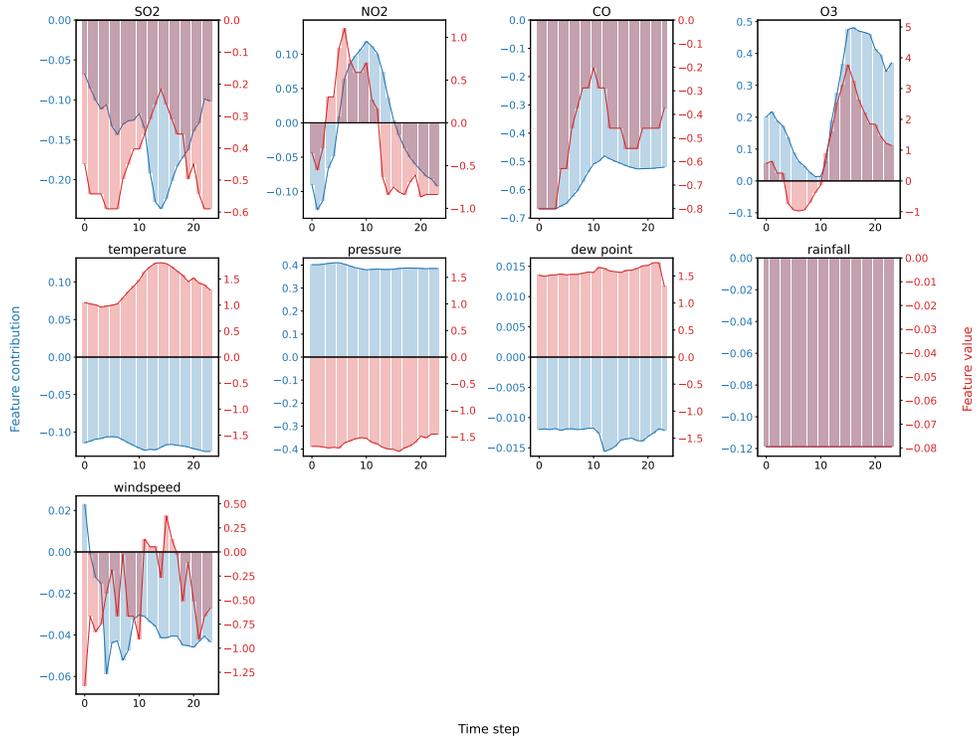


Figure 5: Local time-dependent feature contributions.

217 **Time-step importance:** We plot the average attention scores at the last time step T in Figure 2.
 218 The process for extracting the average attention score of time step u at time step t is formalized as
 219 $\sum_{k=1}^K a_{k,t,u}$. This process is repeated over all data samples, and the results are averaged. Based
 220 on Figure 2, it seems that GATSM pays more attention to the initial and last states than to the
 221 intermediate states. This indicates that the current concentration of particulate matter depends on the
 222 initial state.

223 **Global feature contribution:** Figure 3 illustrates the global behavior of features in the
 224 AirQuality dataset, with red bars indicating the density of training samples. We extract
 225 $\sum_{k=1}^K h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out}$ from GATSM and repeat this process over the range of minimum to
 226 maximum feature values to plot the line. We found that the behavior of SO_2 , O_3 , and $wind\ speed$ is
 227 inconsistent with prior human knowledge. Typically, high levels of SO_2 and O_3 are associated with
 228 poor air quality. However, GATSM learned that particulate matter concentration starts to decrease
 229 when SO_2 exceeds 10 and O_3 exceeds 5. This discrepancy may be due to sparse training samples in
 230 these regions, leading to insufficient training, or there may be interactions with other features. Another
 231 known fact is that high $wind\ speed$ decreases particulate matter concentration. This is consistent when
 232 $wind\ speed$ is below 0.7 in our observation. However, particulate matter concentration drastically
 233 increases when $wind\ speed$ exceeds 0.7, likely due to the wind causing yellow dust.

234 **Local time-independent feature contribution:** To interpret the prediction of a data sample, we
 235 plot the local time-independent feature contributions, $\sum_{k=1}^K h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out}$, in Figure 4. The
 236 main x-axis (blue) represents feature contribution, the sub x-axis (red) represents feature value, and
 237 the y-axis represents time steps. We found that SO_2 , NO_2 , CO , and O_3 have positive correlations.
 238 In contrast, $temperature$, $pressure$, $dew\ point$, and $wind\ speed$ have negative correlations. These are
 239 consistent with the global interpretations shown in Figure 3. Rainfall has the same values across all
 240 time steps.

241 **Local time-dependent feature contribution:** We also visualize the local time-dependent feature con-
 242 tributions, $\sum_{k=1}^K a_{k,t,u} h_b(x_{t,m}) w_{m,b}^{nbm} w_{k,m}^{out}$. Figure 5 illustrates the interpretation of the same data
 243 sample as in Figure 4. The time-dependent interpretation differs slightly from the time-independent
 244 interpretation. We found that there are time lags in SO_2 , NO_2 , CO , and O_3 , meaning previous feature
 245 values affect current feature contributions. For example, in the case of SO_2 , low feature values around
 246 time step 5 lead to low feature contributions around time step 13.

247 6 Future Works & Conclusion

248 Although GATSM achieved state-of-the-art performance within the transparent model category,
 249 it has several limitations. This section discusses these limitations and suggests future work to
 250 address them. GAMs have relatively slower computational times and larger model sizes compared to
 251 black-box models because they require the same number of feature functions as input features. To
 252 address this problem, methods such as the basis strategy can be proposed to reduce the number of
 253 feature functions, or entirely new methods for transparent models can be developed. The attention
 254 mechanism in GATSM may be a bottleneck. Fast attention mechanisms proposed in the literature
 255 [39, 40, 41, 42, 43], or the recently proposed Mamba [44], can help overcome this limitation. Existing
 256 time series models, including GATSM, only handle discrete time series and have limited length
 257 generalization ability, resulting in significantly reduced performance when very long sequences,
 258 unseen during training, are input. Extending GATSM to continuous models using NeuralODE [45]
 259 or HiPPO [46] could address this issue. GATSM still cannot learn higher-order feature interactions
 260 internally and shows low performance on complex datasets. Feature interaction methods proposed
 261 for transparent models may help address this problem [29, 15].

262 In this paper, we proposed a novel transparent model for time series named GATSM. GATSM
 263 consists of time-sharing NBM and the temporal module to effectively learn feature representations
 264 and temporal patterns while maintaining transparency. The experimental results demonstrated that
 265 GATSM has superior generalization ability and is the only transparent model with performance
 266 comparable to Transformer. We provided various visual interpretations of GATSM, demonstrated that
 267 GATSM capture interesting patterns in time series data. We anticipate that GATSM will be widely
 268 adopted in various fields and demonstrate strong performance. The broader impacts of GATSM
 269 across various fields can be found in Appendix A.

270 References

- 271 [1] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining
272 the Predictions of Any Classifier. In *ACM SIGKDD International Conference on Knowledge
273 Discovery and Data Mining*, 2016.
- 274 [2] Scott M. Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. In
275 *Advances in Neural Information Processing Systems*, 2017.
- 276 [3] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
277 Parikh, and Dhruv Batra. Grad-CAM: Visual Explanations From Deep Networks via Gradient-
278 Based Localization. 2017.
- 279 [4] Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. Explaining Machine Learning
280 Classifiers through Diverse Counterfactual Explanations. In *Proceedings of the 2020 Conference
281 on Fairness, Accountability, and Transparency*, 2020.
- 282 [5] Cynthia Rudin. Please Stop Explaining Black Box Models for High Stakes Decisions. In
283 *Advances in Neural Information Processing Systems, Workshop on Critiquing and Correcting
284 Trends in Machine Learning*, 2018.
- 285 [6] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions
286 and use interpretable models instead. *Nature Machine Intelligence*, 1:206–215, May 2019.
- 287 [7] Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich
288 Caruana, and Geoffrey E. Hinton. Neural Additive Models: Interpretable Machine Learning
289 with Neural Nets. In *Advances in Neural Information Processing Systems*, 2021.
- 290 [8] Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. NODE-GAM: Neural Generalized
291 Additive Model for Interpretable Deep Learning. In *International Conference on Learning
292 Representations*, 2022.
- 293 [9] Filip Radenovic, Abhimanyu Dubey, and Dhruv Mahajan. Neural Basis Models for Inter-
294 pretability. In *Advances in Neural Information Processing Systems*, 2022.
- 295 [10] Trevor Hastie and Robert Tibshirani. Generalized Additive Models. *Statistical Science*, 1(3):
296 297–318, August 1986.
- 297 [11] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noémie Elhadad. Intel-
298 ligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission.
299 In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.
- 300 [12] Chun-Hao Chang, Sarah Tan, Ben Lengerich, Anna Goldenberg, and Rich Caruana. How
301 Interpretable and Trustworthy are GAMs? In *ACM SIGKDD International Conference on
302 Knowledge Discovery and Data Mining*, 2021.
- 303 [13] Lev V. Utkin, Egor D. Satyukov, and Andrei V. Konstantinov. SurvNAM: The machine learning
304 survival model explanation. *Neural Networks*, 147:81–102, March 2022.
- 305 [14] Sarah Tan, Rich Caruana, Giles Hooker, and Yin Lou. Distill-and-Compare: Auditing Black-
306 Box Models Using Transparent Model Distillation. In *Proceedings of the 2018 AAAI/ACM
307 Conference on AI, Ethics, and Society*, 2018.
- 308 [15] Minkyu Kim, Hyun-Soo Choi, and Jinho Kim. Higher-order Neural Additive Models: An Inter-
309 pretable Machine Learning Model with Feature Interactions. *arXiv preprint arXiv:2209.15409*,
310 2022.
- 311 [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
312 Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Advances in Neural
313 Information Processing Systems*, 2017.
- 314 [17] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert
315 Müller, and Wojciech Samek. On Pixel-Wise Explanations for Non-Linear Classifier Decisions
316 by Layer-Wise Relevance Propagation. *PLoS ONE*, 10(7), July 2015.
- 317 [18] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning Important Features
318 Through Propagating Activation Differences. In *International Conference on Machine Learning*,
319 2017.
- 320 [19] Sajid Ali, Tamer Abuhmed, Shaker El-Sappagh, Khan Muhammad, Jose M. Alonso-Moral,
321 Roberto Confalonieri, Riccardo Guidotti, Javier Del Ser, Natalia Díaz-Rodríguez, and Francisco

- 322 Herrera. Explainable Artificial Intelligence (XAI): What we know and what is left to attain
323 Trustworthy Artificial Intelligence. *Information Fusion*, 99:101805, November 2023.
- 324 [20] Vikas Hassija, Vinay Chamola, Atmesh Mahapatra, Abhinandan Singal, Divyansh Goel, Kaizhu
325 Huang, Simone Scardapane, Indro Spinelli, Mufti Mahmud, and Amir Hussain. Interpreting
326 Black-Box Models: A Review on Explainable Artificial Intelligence. *Cognitive Computation*,
327 16(1):45–74, January 2024.
- 328 [21] Grace Wahba. *Spline Models for Observational Data*. SIAM, September 1990.
- 329 [22] Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligible Models for Classification and
330 Regression. In *ACM SIGKDD International Conference on Knowledge Discovery and Data
331 Mining*, 2012.
- 332 [23] Harsha Nori, Samuel Jenkins, Paul Koch, and Rich Caruana. InterpretML: A Unified Framework
333 for Machine Learning Interpretability. *arXiv preprint arXiv:1909.09223*, 2019.
- 334 [24] William J. E. Potts. Generalized Additive Neural Networks. In *ACM SIGKDD International
335 Conference on Knowledge Discovery and Data Mining*, 1999.
- 336 [25] Shiyun Xu, Zhiqi Bu, Pratik Chaudhari, and Ian J. Barnett. Sparse Neural Additive Model:
337 Interpretable Deep Learning with Feature Selection via Group Sparsity. In *Joint European
338 Conference on Machine Learning and Knowledge Discovery in Databases*, 2023.
- 339 [26] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. Accurate intelligible models with
340 pairwise interactions. In *ACM SIGKDD International Conference on Knowledge Discovery and
341 Data Mining*, 2013.
- 342 [27] Zebin Yang, Aijun Zhang, and Agus Sudjianto. GAMI-Net: An Explainable Neural Network
343 based on Generalized Additive Models with Structured Interactions. *Pattern Recognition*, 120:
344 108192, December 2021.
- 345 [28] James Enouen and Yan Liu. Sparse Interaction Additive Networks via Feature Interaction
346 Detection and Sparse Selection. *Advances in Neural Information Processing Systems*, 35, 2022.
- 347 [29] Abhimanyu Dubey, Filip Radenovic, and Dhruv Mahajan. Scalable Interpretability via Polyno-
348 mials. *Advances in Neural Information Processing Systems*, 2022.
- 349 [30] Wonkeun Jo and Dongil Kim. Neural additive time-series models: Explainable deep learning
350 for multivariate time-series prediction. *Expert Systems with Applications*, 228:120307, October
351 2023.
- 352 [31] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
353 Bengio. Graph Attention Networks. In *International Conference on Learning Representations*,
354 2018.
- 355 [32] Chang Wei Tan, Christoph Bergmeir, Francois Petitjean, and Geoffrey I. Webb. Monash Univer-
356 sity, UEA, UCR Time Series Extrinsic Regression Archive. *arXiv preprint arXiv:2006.10996*,
357 2020.
- 358 [33] Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom,
359 Paul Southam, and Eamonn Keogh. The UEA multivariate time series classification archive,
360 2018. *arXiv preprint arXiv:1811.00075*, 2018.
- 361 [34] Ary L. Goldberger, Luis A. N. Amaral, Leon Glass, Jeffrey M. Hausdorff, Plamen Ch. Ivanov,
362 Roger G. Mark, Joseph E. Mietus, George B. Moody, Chung-Kang Peng, and H. Eugene Stanley.
363 PhysioBank, PhysioToolkit, and PhysioNet. *Circulation*, 101(23):e215–e220, June 2000.
- 364 [35] Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proceedings
365 of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
366 2016.
- 367 [36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
368 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
369 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
370 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style,
371 High-Performance Deep Learning Library. In *Advances in Neural Information Processing
372 Systems*, 2019.
- 373 [37] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:
374 A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM
375 SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019.

- 376 [38] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International*
377 *Conference on Learning Representations*, 2019.
- 378 [39] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are
379 RNNs: Fast Autoregressive Transformers with Linear Attention. In *International Conference*
380 *on Machine Learning*, 2020.
- 381 [40] Lovish Madaan, Srinadh Bhojanapalli, Himanshu Jain, and Prateek Jain. Treeformer: Dense
382 Gradient Trees for Efficient Attention Computation. In *International Conference on Learning*
383 *Representations*, 2023.
- 384 [41] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-Attention
385 with Linear Complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 386 [42] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,
387 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger,
388 Lucy Colwell, and Adrian Weller. Rethinking Attention with Performers. In *International*
389 *Conference on Learning Representations*, 2021.
- 390 [43] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The Efficient Transformer. In
391 *International Conference on Learning Representations*, 2020.
- 392 [44] Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces,
393 2023.
- 394 [45] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural Ordinary
395 Differential Equations. In *Advances in Neural Information Processing Systems*, 2018.
- 396 [46] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. HiPPO: Recurrent Memory
397 with Optimal Polynomial Projections. In *Advances in Neural Information Processing Systems*,
398 2020.
- 399 [47] Eric J. Pedersen, David L. Miller, Gavin L. Simpson, and Noam Ross. Hierarchical generalized
400 additive models in ecology: an introduction with mgcv. *PeerJ*, 7:e6876, May 2019.
- 401 [48] Trevor Hastie and Robert Tibshirani. Generalized additive models for medical research. *Statisti-*
402 *cal Methods in Medical Research*, 4(3):187–196, September 1995.
- 403 [49] Appliances Energy Dataset, 2020. URL <https://doi.org/10.5281/zenodo.3902637>.
- 404 [50] Australia Rainfall Dataset, 2020. URL <https://doi.org/10.5281/zenodo.3902654>.
- 405 [51] Beijing PM10 Dataset, 2020. URL <https://doi.org/10.5281/zenodo.3902667>.
- 406 [52] Classification of Heart Sound Recordings: The PhysioNet/Computing in Cardiology Challenge
407 2016, 2016. URL <https://physionet.org/content/challenge-2016/1.0.0/>.
- 408 [53] Predicting Mortality of ICU Patients: The PhysioNet/Computing in Cardiology Challenge 2012,
409 2012. URL <https://physionet.org/content/challenge-2012/1.0.0/>.
- 410 [54] Early Prediction of Sepsis from Clinical Data: The PhysioNet/Computing in Cardiology Chal-
411 lenge 2019, 2019. URL <https://physionet.org/content/challenge-2019/1.0.0/>.
- 412 [55] PLAsTiCC Astronomical Classification, 2018. URL [https://www.kaggle.com/c/](https://www.kaggle.com/c/PLAsTiCC-2018)
413 [PLAsTiCC-2018](https://www.kaggle.com/c/PLAsTiCC-2018).
- 414 [56] AALTD’16 Time Series Classification Contest, 2016. URL [https://aaltd16.irisa.fr/](https://aaltd16.irisa.fr/challenge/)
415 [challenge/](https://aaltd16.irisa.fr/challenge/).
- 416 [57] Ignacio Oguiza. tsai - a state-of-the-art deep learning library for time series and sequential data.
417 Github, 2023. URL <https://github.com/timeseriesAI/tsai>.
- 418 [58] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep
419 Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, 2012.
- 420 [59] Saining Xie, Ross Girshick, Piotr Dollar, Zhuowen Tu, and Kaiming He. Aggregated Residual
421 Transformations for Deep Neural Networks. In *Proceedings of the IEEE Conference on*
422 *Computer Vision and Pattern Recognition*, 2017.
- 423 [60] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for Hyper-
424 Parameter Optimization. In *Advances in Neural Information Processing Systems*, 2011.
- 425 [61] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast
426 and Memory-Efficient Exact Attention with IO-Awareness. *arXiv preprint arXiv:2205.14135*,
427 2022.

428 **NeurIPS Paper Checklist**

429 **1. Claims**

430 Question: Do the main claims made in the abstract and introduction accurately reflect the
431 paper's contributions and scope?

432 Answer: [Yes]

433 Justification: The main claims made in the abstract and introduction accurately reflect the
434 paper's contributions and scope.

435 Guidelines:

- 436 • The answer NA means that the abstract and introduction do not include the claims
437 made in the paper.
- 438 • The abstract and/or introduction should clearly state the claims made, including the
439 contributions made in the paper and important assumptions and limitations. A No or
440 NA answer to this question will not be perceived well by the reviewers.
- 441 • The claims made should match theoretical and experimental results, and reflect how
442 much the results can be expected to generalize to other settings.
- 443 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
444 are not attained by the paper.

445 **2. Limitations**

446 Question: Does the paper discuss the limitations of the work performed by the authors?

447 Answer: [Yes]

448 Justification: The limitations of our work are described in section 6.

449 Guidelines:

- 450 • The answer NA means that the paper has no limitation while the answer No means that
451 the paper has limitations, but those are not discussed in the paper.
- 452 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 453 • The paper should point out any strong assumptions and how robust the results are to
454 violations of these assumptions (e.g., independence assumptions, noiseless settings,
455 model well-specification, asymptotic approximations only holding locally). The authors
456 should reflect on how these assumptions might be violated in practice and what the
457 implications would be.
- 458 • The authors should reflect on the scope of the claims made, e.g., if the approach was
459 only tested on a few datasets or with a few runs. In general, empirical results often
460 depend on implicit assumptions, which should be articulated.
- 461 • The authors should reflect on the factors that influence the performance of the approach.
462 For example, a facial recognition algorithm may perform poorly when image resolution
463 is low or images are taken in low lighting. Or a speech-to-text system might not be
464 used reliably to provide closed captions for online lectures because it fails to handle
465 technical jargon.
- 466 • The authors should discuss the computational efficiency of the proposed algorithms
467 and how they scale with dataset size.
- 468 • If applicable, the authors should discuss possible limitations of their approach to
469 address problems of privacy and fairness.
- 470 • While the authors might fear that complete honesty about limitations might be used by
471 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
472 limitations that aren't acknowledged in the paper. The authors should use their best
473 judgment and recognize that individual actions in favor of transparency play an impor-
474 tant role in developing norms that preserve the integrity of the community. Reviewers
475 will be specifically instructed to not penalize honesty concerning limitations.

476 **3. Theory Assumptions and Proofs**

477 Question: For each theoretical result, does the paper provide the full set of assumptions and
478 a complete (and correct) proof?

479 Answer: [NA]

480 Justification: Our work does not include theoretical results.

481 Guidelines:

- 482 • The answer NA means that the paper does not include theoretical results.
- 483 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 484 referenced.
- 485 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 486 • The proofs can either appear in the main paper or the supplemental material, but if
- 487 they appear in the supplemental material, the authors are encouraged to provide a short
- 488 proof sketch to provide intuition.
- 489 • Inversely, any informal proof provided in the core of the paper should be complemented
- 490 by formal proofs provided in appendix or supplemental material.
- 491 • Theorems and Lemmas that the proof relies upon should be properly referenced.

492 4. Experimental Result Reproducibility

493 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

494 perimental results of the paper to the extent that it affects the main claims and/or conclusions

495 of the paper (regardless of whether the code and data are provided or not)?

496 Answer: [Yes]

497 Justification: We provided experimental setup and implementation details in section 5.1 and

498 Appendix C.

499 Guidelines:

- 500 • The answer NA means that the paper does not include experiments.
- 501 • If the paper includes experiments, a No answer to this question will not be perceived
- 502 well by the reviewers: Making the paper reproducible is important, regardless of
- 503 whether the code and data are provided or not.
- 504 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 505 to make their results reproducible or verifiable.
- 506 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 507 For example, if the contribution is a novel architecture, describing the architecture fully
- 508 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 509 be necessary to either make it possible for others to replicate the model with the same
- 510 dataset, or provide access to the model. In general, releasing code and data is often
- 511 one good way to accomplish this, but reproducibility can also be provided via detailed
- 512 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 513 of a large language model), releasing of a model checkpoint, or other means that are
- 514 appropriate to the research performed.
- 515 • While NeurIPS does not require releasing code, the conference does require all submis-
- 516 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 517 nature of the contribution. For example
- 518 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 519 to reproduce that algorithm.
- 520 (b) If the contribution is primarily a new model architecture, the paper should describe
- 521 the architecture clearly and fully.
- 522 (c) If the contribution is a new model (e.g., a large language model), then there should
- 523 either be a way to access this model for reproducing the results or a way to reproduce
- 524 the model (e.g., with an open-source dataset or instructions for how to construct
- 525 the dataset).
- 526 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 527 authors are welcome to describe the particular way they provide for reproducibility.
- 528 In the case of closed-source models, it may be that access to the model is limited in
- 529 some way (e.g., to registered users), but it should be possible for other researchers
- 530 to have some path to reproducing or verifying the results.

531 5. Open access to data and code

532 Question: Does the paper provide open access to the data and code, with sufficient instruc-

533 tions to faithfully reproduce the main experimental results, as described in supplemental

534 material?

535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586

Answer: [Yes]

Justification: We used public datasets and opened our code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We described the experimental setting in section 5.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provided standard deviations with experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 587
- 588
- 589
- 590
- 591
- 592
- 593
- 594
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

595 8. Experiments Compute Resources

596 Question: For each experiment, does the paper provide sufficient information on the com-
597 puter resources (type of compute workers, memory, time of execution) needed to reproduce
598 the experiments?

599 Answer: [Yes]

600 Justification: We provided information on the computational resource used in the experi-
601 ments.

602 Guidelines:

- 603
- 604
- 605
- 606
- 607
- 608
- 609
- 610
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

611 9. Code Of Ethics

612 Question: Does the research conducted in the paper conform, in every respect, with the
613 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

614 Answer: [Yes]

615 Justification: Our work conform with the NeurIPS Code of Ethics.

616 Guidelines:

- 617
- 618
- 619
- 620
- 621
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

622 10. Broader Impacts

623 Question: Does the paper discuss both potential positive societal impacts and negative
624 societal impacts of the work performed?

625 Answer: [Yes]

626 Justification: We discussed the potential impacts of GATSM in Appendix A.

627 Guidelines:

- 628
- 629
- 630
- 631
- 632
- 633
- 634
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647
- 648
- 649
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

650 11. Safeguards

651 Question: Does the paper describe safeguards that have been put in place for responsible
652 release of data or models that have a high risk for misuse (e.g., pretrained language models,
653 image generators, or scraped datasets)?

654 Answer: [NA]

655 Justification: Our work poses no such risks.

656 Guidelines:

- 657
- 658
- 659
- 660
- 661
- 662
- 663
- 664
- 665
- 666
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

667 12. Licenses for existing assets

668 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
669 the paper, properly credited and are the license and terms of use explicitly mentioned and
670 properly respected?

671 Answer: [Yes]

672 Justification: We properly cited the used codes and data.

673 Guidelines:

- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

687 • If this information is not available online, the authors are encouraged to reach out to
688 the asset’s creators.

689 13. New Assets

690 Question: Are new assets introduced in the paper well documented and is the documentation
691 provided alongside the assets?

692 Answer: [Yes]

693 Justification: We opened the source code of GATSM, and the document to run the code is
694 provided along with the code.

695 Guidelines:

- 696 • The answer NA means that the paper does not release new assets.
- 697 • Researchers should communicate the details of the dataset/code/model as part of their
698 submissions via structured templates. This includes details about training, license,
699 limitations, etc.
- 700 • The paper should discuss whether and how consent was obtained from people whose
701 asset is used.
- 702 • At submission time, remember to anonymize your assets (if applicable). You can either
703 create an anonymized URL or include an anonymized zip file.

704 14. Crowdsourcing and Research with Human Subjects

705 Question: For crowdsourcing experiments and research with human subjects, does the paper
706 include the full text of instructions given to participants and screenshots, if applicable, as
707 well as details about compensation (if any)?

708 Answer: [NA]

709 Justification: Our work does not involve crowdsourcing nor research with human subjects.

710 Guidelines:

- 711 • The answer NA means that the paper does not involve crowdsourcing nor research with
712 human subjects.
- 713 • Including this information in the supplemental material is fine, but if the main contribu-
714 tion of the paper involves human subjects, then as much detail as possible should be
715 included in the main paper.
- 716 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
717 or other labor should be paid at least the minimum wage in the country of the data
718 collector.

719 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 720 Subjects

721 Question: Does the paper describe potential risks incurred by study participants, whether
722 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
723 approvals (or an equivalent approval/review based on the requirements of your country or
724 institution) were obtained?

725 Answer: [NA]

726 Justification: Our work does not involve crowdsourcing nor research with human subjects.

727 Guidelines:

- 728 • The answer NA means that the paper does not involve crowdsourcing nor research with
729 human subjects.
- 730 • Depending on the country in which research is conducted, IRB approval (or equivalent)
731 may be required for any human subjects research. If you obtained IRB approval, you
732 should clearly state this in the paper.
- 733 • We recognize that the procedures for this may vary significantly between institutions
734 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
735 guidelines for their institution.
- 736 • For initial submissions, do not include any information that would break anonymity (if
737 applicable), such as the institution conducting the review.

738 **A Broader impact**

739 We discuss the expected impacts of GATSM across various fields.

- 740 • **Time series adaptation:** GATSM extends existing GAMs to time series, enabling tasks that
741 traditional GAMs could not perform in this context - e.g., better performance on time series and
742 finding temporal patterns.
- 743 • **Improved decision-making system:** GATSM can show users their exact decision-making process,
744 providing trust and confidence in its predictions to users. This enables decision-makers to make
745 more informed choices, crucial in high-stakes domains such as healthcare.
- 746 • **Ethical AI:** GATSM can examine that their outcomes are biased or discriminatory by displaying
747 the shape of feature functions. This is particularly important in ethically sensitive domains, such as
748 recidivism prediction.
- 749 • **Scientific discovery:** Transparent models have already been used in various research fields for
750 scientific discovery [47, 48]. GATSM also can be applied to these domains to obtain novel scientific
751 insights.

752 Despite these advantages, it is important to remember that the interpretations of transparent models
753 do not necessarily reflect exact causal relationships. While transparent models provide clear and
754 faithful interpretations, they are still not capable of identifying causal relationships. Causal discovery
755 is a complex task that requires further research.

756 **B Dataset details**

757 We use eight publicly available datasets for our experiments. Three datasets - Energy, Rainfall, and
758 AirQuality - can be downloaded from the Monash repository [32]. Another three datasets - Heartbeat,
759 LSST, and NATOPS - are available from the UCR repository [33]. The remaining two datasets can
760 be downloaded from the PhysioNet [34]. Details of the datasets are provided below:

- 761 • **Energy** [49]: This dataset consists of 24 features related to temperature and humidity from sensors
762 and weather conditions. These features are measured every 10 minutes. The goal of this dataset is
763 to predict total energy usage.
- 764 • **Rainfall** [50]: This dataset consists of temperatures measured hourly. The goal of this dataset is to
765 predict total daily rainfall in Australia.
- 766 • **AirQuality** [51]: This dataset consists of features related to air pollutants and meteorological data.
767 The goal of this dataset is to predict the PM10 level in Beijing.
- 768 • **Heartbeat** [52]: This dataset consists of heart sounds collected from various locations on the body.
769 Each sound was truncated to five seconds, and a spectrogram of each instance was created with a
770 window size of 0.061 seconds with a 70% overlap. The goal of this dataset is to classify the sounds
771 as either normal or abnormal.
- 772 • **Mortality** [53] This dataset consists of records of adult patients admitted to the ICU. The input
773 features include the patient demographics, vital signs, and lab results. The goal of this dataset is to
774 predict the in-hospital death of patients.
- 775 • **Sepsis** [54]: This dataset consists of records of ICU patients. The input features include patient
776 demographics, vital signs, and lab results. The goal of this dataset is to predict sepsis six hours in
777 advance at every time step.
- 778 • **LSST** [55]: This challenge dataset aims to classify astronomical time series. These time series
779 consist of six different light curves, simulated based on the data expected from the Large Synoptic
780 Survey Telescope (LSST).
- 781 • **NATOPS** [56]: This dataset aims to classify the Naval Air Training and Operating Procedures
782 Standardization (NATOPS) motions used to control aircraft movements. It consists of 24 features
783 representing the x, y, and z coordinates for each of the eight sensor locations attached to the body.

784 We used `get_UCR_data()` and `get_Monash_regression_data()` functions in the `tsai` library
785 [57] to load the UCR and Monash datasets.

Table 6: Optimal hyper-parameters for GATSM.

GATSM: [256, 256, 128] hidden dims, 100 basis functions								
Dataset	Batch Size	NBM Batch Norm.	NBM Dropout	Attn. Embedding Size	Attn. Heads	Attn. Dropout	Learning Rate	Weight Decay
Energy	32	False	2.315e-1	110	8	6.924e-2	4.950e-3	1.679e-3
Rainfall	32,768	False	5.936e-3	44	7	1.215e-3	9.225e-3	2.204e-6
AirQuality	4,096	False	2.340e-2	81	8	1.169e-1	6.076e-3	5.047e-6
Heartbeat	64	True	1.749e-1	92	2	1.653e-1	8.061e-3	4.787e-6
Mortality	512	False	7.151e-2	125	8	7.324e-1	7.304e-3	2.181e-4
Sepsis	512	True	6.523e-2	90	6	8.992e-1	4.509e-3	2.259e-2
LSST	1,024	False	2.500e-2	59	7	2.063e-1	5.561e-2	5.957e-3
NATOPS	64	True	4.827e-3	49	8	7.920e-1	8.156e-3	2.748e-2

786 C Implementation details

787 We use 13 models, including GATSM, for our experiments. We implement XGBoost and EBM
788 using the `xgboost` [35] and `interpretml` [23] libraries, respectively. For NodeGAM, we employ
789 the official implementation provided by its authors [8]. The remaining models are developed using
790 PyTorch [36]. In addition, we implement the feature functions in NAM and NBM using grouped
791 convolutions [58, 59] to enhance their efficiency. XGBoost and EBM are trained on two AMD EPYC
792 7513 CPUs, while the other models are trained on an NVIDIA A100 GPU with 80GB VRAM. All
793 models undergo hyperparameter tuning via Optuna [37] with the Tree-structured Parzen Estimator
794 (TPE) algorithm [60] in 100 trials. The hyperparameter search space and the optimal hyperparameters
795 for the models are provided below:

796 • **XGBoost:** We tune the `n_estimators` in the integer interval [1, 1000], `max_depth` in the integer
797 interval [0, 2000], `learning_rate` in the continuous interval [1e-6, 1], `subsample` in the continuous
798 interval [0, 1], and `colsample_bytree` in the continuous interval [0, 1].

799 • **MLP, NAM, NBM and NATM:** We tune the `batchnorm` in the discrete set {False, True}, `dropout`
800 in the continuous interval [0, 0.9], `learning_rate` in the continuous interval [1e-3, 1e-2], and
801 `weight_decay` in the continuous interval [1e-6, 1e-1] on a log scale.

802 • **RNN, GRU and LSTM:** We tune the `hidden_size` in the integer interval [8, 128], `dropout`
803 in the continuous interval [0, 0.9], `learning_rate` in the continuous interval [1e-3, 1e-2], and
804 `weight_decay` in the continuous interval [1e-6, 1e-1] on a log scale.

805 • **Transformer:** We tune the `n_layers` in the integer interval [1, 4], `emb_size` in the integer
806 interval [8, 32], `hidden_size` in the integer interval [8, 128], `n_heads` in the integer interval [1,
807 8], `dropout` in the continuous interval [0, 0.9], `learning_rate` in the continuous interval [1e-3,
808 1e-2], and `weight_decay` in the continuous interval [1e-6, 1e-1] on a log scale.

809 • **Linear:** We tune the `learning_rate` in the continuous interval [1e-3, 1e-2], and `weight_decay`
810 in the continuous interval [1e-6, 1e-1] on a log scale.

811 • **EBM:** We tune `max_bins` in the integer interval [8, 512], `min_samples_leaf` and `max_leaves`
812 in the integer interval [1, 50], `inner_bags` and `outer_bags` in the integer interval [1, 128],
813 `learning_rate` in the continuous interval [1e-6, 100] on a log scale, and `max_rounds` in the
814 integer interval [1000, 10000].

815 • **NodeGAM:** We tune `n_trees` in the integer interval [1, 256], `n_layers` and `depth` in the integer
816 intervals [1, 4], `dropout` in the continuous interval [0, 0.9], `learning_rate` in the continuous
817 interval [1e-3, 1e-2], and `weight_decay` in the continuous interval [1e-6, 1e-1] on a log scale.

818 • **GATSM:** We tune `nbm_batchnorm` in the discrete set {False, True}, `nbm_dropout` in the con-
819 tinuous interval [0, 0.9], `attn_emb_size` in the integer interval [8, 128], `attn_n_heads` in the
820 integer interval [1, 8], `attn_dropout` in the continuous interval [0, 0.9], `learning_rate` in the
821 continuous interval [1e-3, 1e-2], and `weight_decay` in the continuous interval [1e-6, 1e-1] on a
822 log scale. The optimal hyper-parameters for GATSM across all experimental datasets are provided
823 in Table 6.

824 D Additional experiments

825 D.1 Inference speed

826 The inference speed of machine learning models is a crucial metric for real-world systems. We
 827 evaluate the throughput of various models. The results are presented in Table 7. Since the datasets
 828 have fewer features than the number of basis functions in NBM, NAM achieves higher throughput
 829 than NBM. Transparent tabular models typically exhibit fast speeds. However, their throughput
 830 significantly decreases in datasets with many features, such as Heartbeat, Mortality, and Sepsis,
 831 because they require the same number of feature functions as the number of input features. Trans-
 832 former shows higher throughput than the transparent time series models because it does not require
 833 feature functions, which are the main bottleneck of transparent models. Additionally, the PyTorch
 834 implementation of Transformer uses the flash attention mechanism [61] to enhance its efficiency.
 835 NATM has slightly higher throughput than GATSM, as it does not require the attention mechanism
 836 and has fewer feature functions compared to the number of basis functions in GATSM.

Table 7: Inference throughput of different models.

	Energy	Rainfall	AirQuality	Heartbeat	Mortality	Sepsis	LSST	NATOPS
NAM	65.3K	1.8M	5.1M	139.1K	772.2K	23.9K	2.3M	147.9K
NBM	45.5K	1.1M	1.0M	55.9K	375.8K	6.5K	1.6M	85.6K
Transformer	30.9K	240.5K	174.2K	15.7K	161.9K	134.6K	214.4K	68.3K
NATM	5.3K	699.3K	241.3K	1.3K	N/A	N/A	28.6K	19.2K
GATSM	6.1K	350.6K	192.8K	1.2K	4.9K	3.8K	126.5K	12.5K

837 D.2 Number of basis functions

838 We evaluate GATSM by varying the number of basis functions in the time-sharing NBM. The results
 839 for forecasting, binary classification, and multi-class classification datasets are presented in Figure 6.
 840 For the Sepsis dataset, using 200 and 300 basis functions causes the out-of-memory error. For the
 841 Energy and Heartbeat datasets, performance improves up to 100 basis functions but shows no further
 842 benefit when the number of bases exceeds 100. In other datasets, performance changes are not
 843 significant with different numbers of basis functions. In addition, there is a trade-off between the
 844 number of basis functions and computational speed. Therefore, we recommend generally setting the
 845 number of basis functions to 100. Note that the performance of GATSM with this hyper-parameter
 846 depends on the dataset size and complexity. Hence, a larger number of basis functions may benefit
 847 more complex datasets.

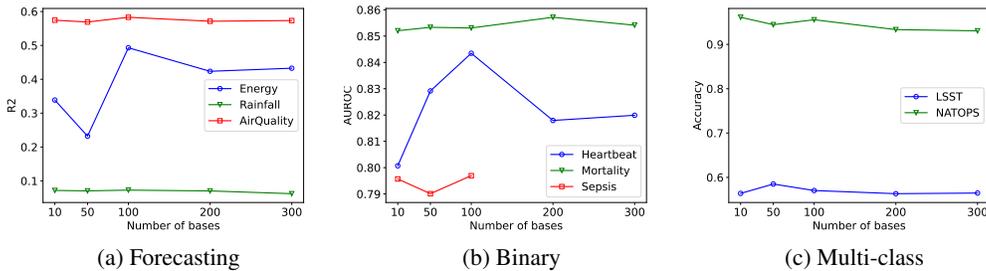


Figure 6: Performances of GATSM on the different number of basis functions.

848 E Additional visualizations

849 In addition to the interpretations on the AirQuality dataset in section 5.4, we present another interesting
 850 interpretations of GATSM on the Rainfall dataset.

851 **Time-step importance:** Figure 7 illustrates the average importance of all time steps at the final time
 852 step. The importance exhibit a cyclical pattern of rising and falling at regular intervals, indicating
 853 that GATSM effectively captures seasonal patterns in the Rainfall dataset.

854 **Global feature contribution:** Figure 8 illustrates the global behavior of features in the Rainfall
855 dataset, with red bars indicating the density of training samples. Our findings indicate that low *Max*
856 *Temperature* and high *Min Temperature* contribute to an increase in rainfall.

857 **Local time-independent feature contribution:** Figure 9 shows the local time-independent feature
858 contributions. Consistent with the global interpretation, *Avg. Temperature* and *Min Temperature* have
859 positive correlations with rainfall, while *Max Temperature* has a negative correlation with rainfall.

860 **Local time-dependent feature contribution:** Figure 10 shows the local time-dependent feature
861 contributions. All features exhibit patterns similar to the local time-independent contributions.
862 However, we found that *Avg. Temperature* and *Min Temperature* have time lags between feature
863 values and contributions.

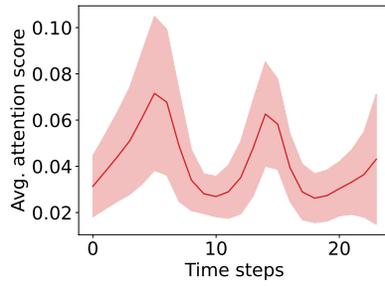


Figure 7: Average attention scores of time steps on the Rainfall dataset.

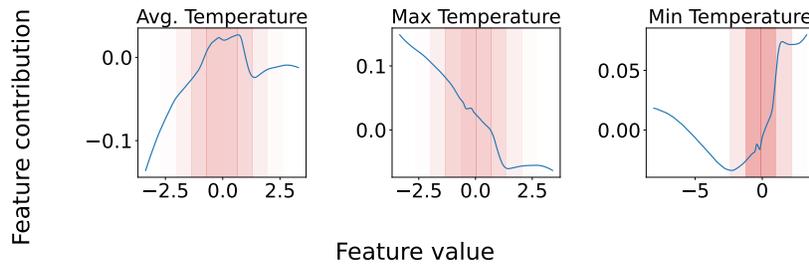


Figure 8: Global interpretations of features in the Rainfall dataset.

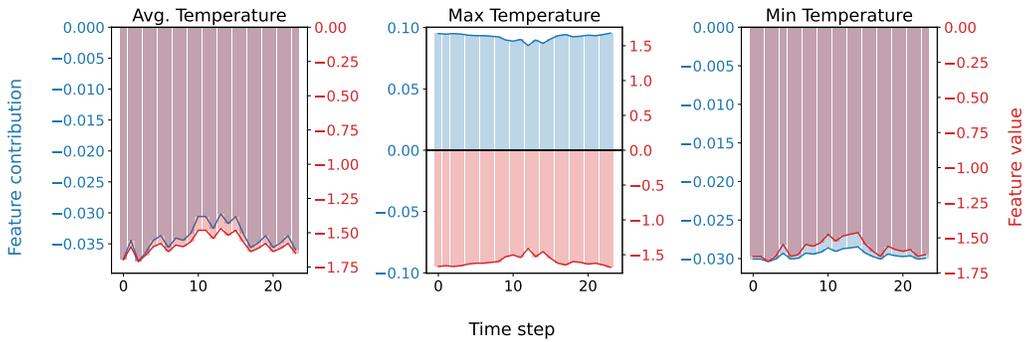


Figure 9: Local time-independent contributions of features in the Rainfall dataset.

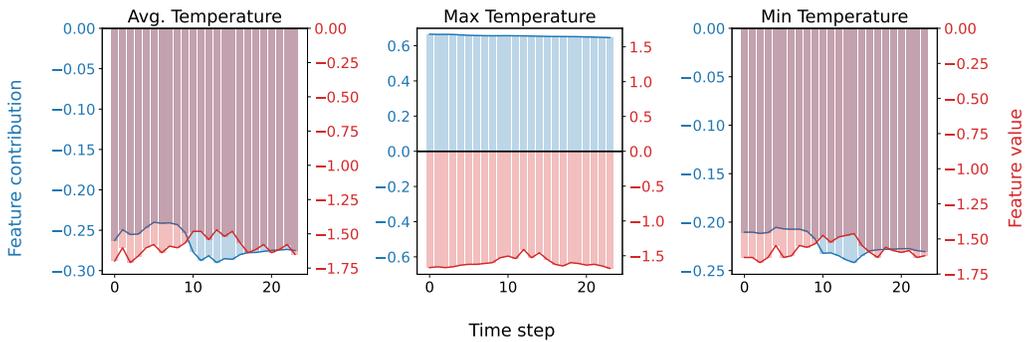


Figure 10: Local time-dependent contributions of features in the Rainfall dataset.