



# MltAuxTSPP: a unified benchmark for deep learning-based traffic state prediction with multi-source auxiliary data<sup>\*#</sup>

Yusong ZHOU<sup>†1,2</sup>, Xiaoyu JIANG<sup>†1,2</sup>, Shu SUN<sup>†1,2</sup>,  
 Xinmin ZHANG<sup>†‡1,2</sup>, Yuanqiu MO<sup>3</sup>, Zhihuan SONG<sup>1,2</sup>

<sup>1</sup>College of Control Science and Engineering, Zhejiang University, Hangzhou 310027, China

<sup>2</sup>State Key Laboratory of Industrial Control Technology, Zhejiang University, Hangzhou 310027, China

<sup>3</sup>School of Mathematics, Southeast University, Nanjing 211189, China

<sup>†</sup>E-mail: zhoushichan@zju.edu.cn; jiangxiaoyu@zju.edu.cn; shusun@zju.edu.cn; xinminzhang@zju.edu.cn

Received Mar. 17, 2025; Revision accepted Aug. 17, 2025; Crosschecked Sept. 18, 2025

**Abstract:** Deep learning has empowered traffic prediction models to integrate diverse auxiliary data sources, such as weather and temporal features, for enhanced forecasting accuracy. However, existing approaches often suffer from limited generality and scalability, and the field lacks a unified benchmark for fair model comparison. This absence hinders consistent performance evaluation, slows the development of robust and adaptable models, and makes it challenging to quantify the incremental benefits of different auxiliary data sources. To address these issues, we present MltAuxTSPP, a unified benchmark framework for deep learning-based traffic state prediction with multi-source auxiliary data. The framework features a standardized data container and a fusion embedding module, enabling consistent utilization of heterogeneous data and improving scalability. It produces unified hidden representations that can be seamlessly adopted by various downstream models, ensuring fair and reproducible comparisons under identical conditions. Extensive experiments on real-world datasets demonstrate that MltAuxTSPP effectively leverages weather and temporal features to improve long-term forecast performance and offers a practical and reproducible foundation for advancing research in traffic state prediction.

**Key words:** Traffic prediction; Benchmark platform; Deep learning; Multi-source auxiliary data

<https://doi.org/10.1631/FITEE.2500169>

**CLC number:** TP18

## 1 Introduction

Accurate traffic state prediction is a cornerstone of modern intelligent transportation systems (ITSs) and is essential for applications ranging from traffic management to autonomous driving. Although deep learning models have excelled at capturing complex spatiotemporal dependencies from traffic data (Do

et al., 2019; Zhu et al., 2019; Lee et al., 2021), their performance is often constrained by relying solely on traffic-related data, such as the fixed position sensor data, trajectory data, network infrastructure data, and trip records data (Nagy and Simon, 2018; Shaygan et al., 2022).

Real-world traffic is significantly influenced by external factors, such as predictable temporal patterns (e.g., weekday rush hours vs. weekends) and adverse weather conditions (Essien et al., 2018; Taghipour et al., 2020). Consequently, recent studies have begun to integrate such auxiliary data to enhance prediction accuracy (Taghipour et al., 2020; Essien et al., 2021). However, the integration of multi-source data remains fraught with challenges.

<sup>‡</sup> Corresponding author

\* Project supported by the National Science and Technology Major Project of China (No. 2022ZD0120003)

# Electronic supplementary materials: The online version of this article (<https://doi.org/10.1631/FITEE.2500169>) contains supplementary materials, which are available to authorized users

ORCID: Yusong ZHOU, <https://orcid.org/0009-0003-4475-623X>; Xinmin ZHANG, <https://orcid.org/0000-0002-4761-3969>

© Zhejiang University Press 2025

Current approaches often handle a limited set of auxiliary variables, which typically lack spatial dimensions. Furthermore, these methods generally lack the generality and scalability required to accommodate diverse and growing data sources. Crucially, the field is hampered by the absence of a unified benchmark platform for systematically evaluating and comparing models under multi-source data conditions.

These limitations collectively pose significant obstacles to advancing the field. The lack of generality hinders the application of existing methods to new datasets or information streams. Poor scalability makes it challenging to handle the increasing volume and variety of data. Crucially, the absence of a unified benchmark leads to fragmented research, impedes objective comparison of models, and ultimately slows ITS innovation.

To overcome these obstacles, we introduce MltAuxTSPP, a novel unified benchmark framework for deep learning-based traffic state prediction using multi-source auxiliary data. At its core, MltAuxTSPP features a unified data container (UDC) and a versatile fusion embedding module (Fig. 1). This design addresses critical challenges in data handling and scalability, enabling seamless integration with various existing spatiotemporal models. This allows researchers to fairly evaluate and compare different model architectures on a standardized basis. The main contributions of this work are as follows:

1. We design and propose MltAuxTSPP, a general and scalable framework for enhancing traffic prediction with multi-source data. Its fusion module standardizes the integration of diverse data types for various downstream models.

2. We construct and release two new multi-

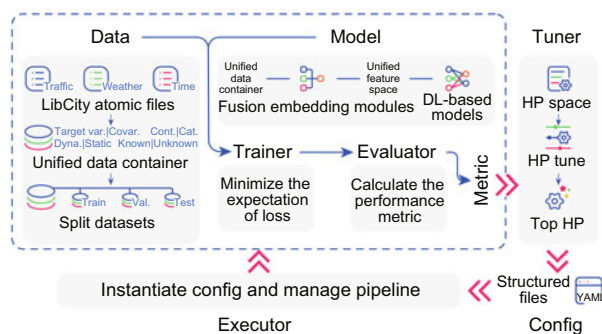
source traffic forecasting datasets by augmenting the widely used METR-LA and PEMS-BAY benchmarks with extensive weather and temporal features, facilitating standardized evaluation. Note that the METR-LA data are listed in the supplementary materials.

3. We conduct extensive experiments to systematically evaluate the impact of auxiliary data on state-of-the-art models and provide novel insights into the significant benefits of multi-source data for long-term forecasting.

## 2 Related works

Traffic forecasting models can be broadly divided into temporal modeling and spatiotemporal modeling. A wide array of temporal models has been applied to traffic prediction, progressing from classical methods, e.g., historical average (HA) and autoregressive integrated moving average (ARIMA) (Liu J and Guan, 2004; Yang et al., 2020), to deep learning architectures like long short-term memory (LSTM) and Transformers. These deep learning architectures have proven highly effective at capturing complex traffic dynamics (Hochreiter and Schmidhuber, 1997; Ma et al., 2015; Vaswani et al., 2017; Zhou et al., 2021). Spatiotemporal models incorporate spatial dependencies primarily through convolutional neural networks (CNNs) or graph convolutional networks (GCNs), which leads to foundational architectures like spatiotemporal GCNs (Yu et al., 2018; Song et al., 2020). Current research, however, focuses on more advanced techniques for simultaneous spatiotemporal modeling. Key innovations include learning dynamic graph structures with adaptive matrices, e.g., Graph WaveNet (Wu ZH et al., 2019), modeling continuous dynamics with ordinary differential equations, e.g., spatial-temporal graph ordinary differential equation neural network (STGODE) (Fang et al., 2021), and developing unified attention mechanisms that concurrently capture spatiotemporal correlations, e.g., T-Graphormer (Bai and Liu, 2025).

Although spatiotemporal models excel at learning from traffic data that are commonly sourced from sensor networks like PeMS (<https://pems.dot.ca.gov>), such as in Yu et al. (2018), Guo et al. (2019), and Song et al. (2020), their performance is inherently limited



**Fig. 1** The benchmarking system in MltAuxTSPP. **DL:** deep learning; **HP:** hyperparameter; **var.:** variable; **covar.:** covariance; **cont.:** continuous; **cat.:** categorical variable; **dyna:** dynamic; **val.:** validate

by overlooking external factors. Consequently, although many studies now incorporate auxiliary data such as weather and time features (see more in the supplementary materials), their fusion methodologies exhibit critical bottlenecks. Specifically, current methods often rely on simplistic techniques like concatenation, failing to capture complex interdependencies. They also tend to treat auxiliary data as global variables, thereby ignoring crucial spatial heterogeneity (e.g., localized weather effects). Furthermore, these approaches typically lack the scalability and generalizability required to handle diverse or novel auxiliary data streams. These challenges underscore the fact that effective multi-source data utilization requires more than just data availability; advanced, scalable, and spatially aware fusion strategies are paramount. As shown in Table 1, although several valuable benchmark platforms exist, they are not tailored for multi-source traffic prediction. Platforms like time-series prediction platform (TSPP), time-series library (TSLib), and basic time series (BasicTS<sup>+</sup>) focus on general time-series forecasting (Wu HX et al., 2023; Bączek et al., 2024; Shao et al., 2025). Conversely, even comprehensive traffic-specific libraries like LibCity (Wang et al., 2021) primarily handle core traffic data and timestamps, offering limited support for integrating diverse auxiliary data sources.

Consequently, a critical gap remains; there is no unified platform designed to systematically evaluate how different models leverage rich auxiliary data (e.g., weather and events) in traffic forecasting. This absence hinders the objective comparison of data fusion strategies, and slows progress in developing truly robust models. To address this, we develop

MltAuxTSPP as a unified and extensible environment specifically for this purpose.

### 3 Multi-source auxiliary data fusion framework

#### 3.1 Task formulation

In this work, the task of traffic state prediction with multi-source auxiliary data uses historical traffic states to forecast future traffic states, incorporating auxiliary data from both historical and future perspectives. Traffic state refers to the conditions of traffic flow on a road network, typically characterized by traffic volume, speed, density, and congestion levels. When no auxiliary data are used, the task reduces to single-source spatiotemporal forecasting. If spatial dependencies are ignored or flattened into features, it further simplifies to time-series forecasting.

Given the historical target variable  $\mathbf{X} \in \mathbb{R}^{T \times n \times f}$ , auxiliary variables  $\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \dots, \mathbf{C}^{(m)}$ , and the road network  $G = (N, E)$ , the task is to predict a future sequence  $\mathbf{X}^* \in \mathbb{R}^{h \times n \times f}$ , where  $T$  is the historical input length,  $h$  is the length of the future horizon,  $n$  is the road network node count,  $f$  is the feature count, and  $m$  is the number of auxiliary variables. Each auxiliary variable  $\mathbf{C}^{(j_c)} \in \mathbb{R}^{T^{(j_c)} \times n^{(j_c)} \times f^{(j_c)}}$  is characterized by its own temporal span  $T^{(j_c)}$ , node count  $n^{(j_c)}$ , and feature count  $f^{(j_c)}$  ( $j_c \in 1, 2, \dots, m$ ). The graph  $G$  with node set  $N$  and edge set  $E$  is encoded as an adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Entries  $A_{a,b}$  indicate either binary connectivity or weighted relationships (e.g., distance and travel time) between nodes  $N_a$  and  $N_b$  ( $a$  and  $b \in [1, n]$ ). Notably,  $\mathbf{A}$  is treated as one of the auxiliary variables.

**Table 1 Comparison of benchmark platforms for time series and traffic prediction**

Framework	Primary focus	Key feature/scope	Auxiliary data support
TSPP	General time series	Including nine models (e.g., temporal fusion Transformers (Lim et al., 2021)) and some traffic datasets	Limited
TSLib	General time series	Providing advanced models like Informer and TimesNet	Limited
BasicTS <sup>+</sup>	General time series	Focusing on multi-variable, spatiotemporal, and long-term forecasting	Limited
LibCity	Traffic prediction	Comprehensive library with 60 models and 35 datasets for various traffic tasks	Limited
Ours	Multi-source traffic prediction	Unified environment with data containers and fusion modules for diverse auxiliary inputs	Yes

The task uses the model  $M : \{\mathbf{X}, \mathbf{C}^{(i)}\} \rightarrow \mathbb{R}^{t \times n \times f}$  ( $i = 1, 2, \dots, m$ ) to predict the future sequence  $\mathbf{X}^*$  and minimize the expected error  $\mathbb{E}[\mathcal{L}(\mathbf{X}^*, \widehat{\mathbf{X}}^*)]$ , where  $\mathcal{L}$  measures the discrepancy between prediction  $\mathbf{X}^*$  and the ground truth  $\widehat{\mathbf{X}}^*$ .

### 3.2 Data collection and preprocessing

Our study leverages two standard traffic benchmark datasets, METR-LA and PEMS-BAY (Li et al., 2018), augmented with external data to investigate their impact. The auxiliary data include three global weather datasets, ERA5-Land (<https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land?tab=overview>), ERA5-HEAT (<https://cds.climate.copernicus.eu/datasets/derived-utci-historical?tab=overview>), and the global hourly-integrated surface database (GH-ISD) (<https://www.ncei.noaa.gov/products/land-based-station/integrated-surface-database#tab=336>), and three discrete temporal features, minute of the hour (MoH), hour of the day (HoD), and day of the week (DoW). To unify the spatiotemporal heterogeneity of these sources, we develop a preprocessing pipeline. First, all weather data are temporally interpolated to a 5-min frequency to match the traffic data. Second, all datasets are spatially converted to point-based representations, with weather data aligned to the nearest traffic sensor. Finally, we develop two enriched datasets in Table 2, and the detailed variable descriptions are provided in the supplementary materials.

### 3.3 Framework of the benchmarking system

We design the benchmarking system in Fig. 1 with the primary goal of integrating multi-source data and fairly comparing the performance of different models. The system adopts a top-down design approach that divides the entire architecture into several modules, each responsible for a specific function to achieve modularity. This design approach not only enhances the maintainability and scalability of the system, but also improves its flexibility and comprehensibility.

The system modules are categorized into two main types: core modules and auxiliary modules. The core modules include the data module, model module, trainer module, evaluator module, and executor module, whereas the auxiliary modules con-

Table 2 Summary of datasets

Source	Interval (min)	METR-LA		
		$N_{\text{step}}$	$N_{\text{node}}$	$N_{\text{var}}$
Traffic	5	34 272	207	2
ERA5-Land			10	10
ERA5-HEAT			4	2
GH-ISD			8	1
Source	Interval (min)	PEMS-BAY		
		$N_{\text{step}}$	$N_{\text{node}}$	$N_{\text{var}}$
Traffic	5	52 116	325	2
ERA5-Land			9	10
ERA5-HEAT			2	2
GH-ISD			5	1

$N_{\text{step}}$ : number of steps;  $N_{\text{node}}$ : number of nodes;  $N_{\text{var}}$ : number of variables

sist of the config and tuner modules. The core modules are implemented in Python, and the interfaces and behaviors of the modules are defined through abstract base classes. Through derived classes, these modules can adapt to a variety of scenario requirements. Specifically, the functions of each module are shown in Table 3.

The core and auxiliary modules are interconnected through the executor and Hydra (Meta Research, 2019). The executor instantiates the modules based on the configuration files, invokes their interfaces to complete the training and evaluation tasks, and then feeds the results back to Optuna (Akiba et al., 2019) for further optimization.

In this work, we address the fusion of multi-source auxiliary data using the pipeline shown in Fig. 2. First, we store the multi-source data in a UDC. Then, the data are transformed into a format suitable for model input through a fusion embedding module. The output of the module is a hidden representation that integrates the multi-source data, ensuring compatibility with the input requirements of the subsequent models. A more detailed explanation of this process will be provided in the following subsections.

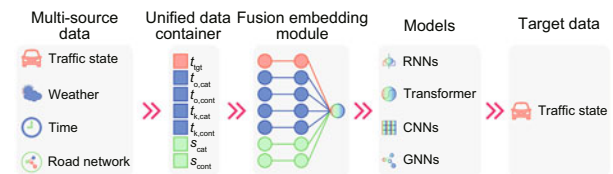


Fig. 2 Multi-source auxiliary data fusion framework for traffic state prediction. tgt: target; k: known; o: observed

**Table 3 Description of the framework modules**

Module	Description
Data	Responsible for preparing the datasets required by the models, including data integration, filtering, cleaning, and constructing UDCs
Model	Provides a series of reproducible or custom models to meet the requirements of prediction tasks
Trainer	Responsible for selecting optimization strategies and customizing relevant decisions to improve model performance
Evaluator	Assesses the accuracy and effectiveness of the model predictions
Executor	Controls the instantiation of modules based on the configuration and executes the training, inference, and evaluation processes using standardized interfaces
Config	Manages all framework-related parameters centrally through YAML files, providing structured and flexible control and adaptability
Tuner	Uses Optuna (Akiba et al., 2019) to select and apply the optimal configuration parameters after each module execution to optimize model performance

### 3.4 UDC

To manage and unify multi-source data, we use a UDC to store traffic state datasets with auxiliary information, which is an extension of the data aspects in TSPP (Bączek et al., 2024). The dataset can be categorized into the seven components in Table 4 based on the following attributes. Notably, the target variable is dynamic and unknown, and can be either continuous or categorical, whereas covariates may include multiple variables.

The data in our study are classified along four dimensions. By source, we distinguish between the target variables to be predicted and covariates. By temporal nature, data are either static (e.g., network topology) or dynamic. With respect to future availability, inputs can be known at prediction time (e.g., timestamps) or unknown (the target value and weather). Finally, by data type, variables are either continuous (real-valued) or categorical (discrete).

### 3.5 Fusion embedding modules

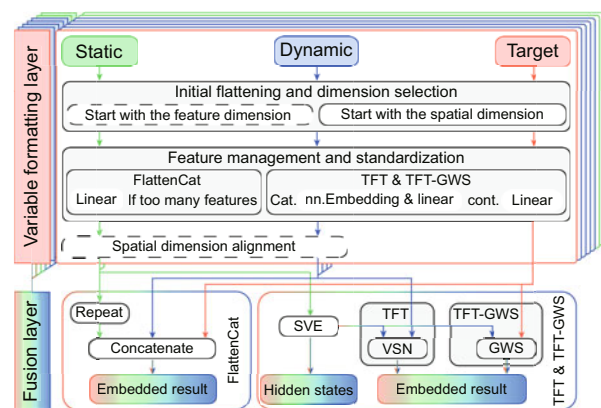
To improve the integration of multi-source data and facilitate model reuse, we propose an abstract

**Table 4 Components of UDC**

Component	Source	Temporal	Observation	Data type
$t_{tgt}$	Target	Dynamic	Unknown	Continuous, categorical
$s_{cont}$	Covariates	Static	Known	Continuous
$s_{cat}$			Known	Categorical
$t_{k,cont}$		Dynamic	Known	Continuous
$t_{k,cat}$			Known	Categorical
$t_{o,cont}$			Unknown	Continuous
$t_{o,cat}$			Unknown	Categorical

fusion embedding module designed specifically for handling such data in Fig. 3. This module acts as an intermediary between the UDC and the model, with the primary objective of projecting multi-source data into a unified feature space for input into subsequent predictive models. We present three concrete instances of this abstract module: FlattenCat, temporal fusion Transformer (TFT), and TFT-gated weighted summation (TFT-GWS). FlattenCat standardizes variables through linear transformation, TFT (Lim et al., 2021) incorporates nonlinear processing and variable selection, and TFT-GWS is a simplified version of TFT that retains the core concepts while reducing computational overhead.

The abstract module consists of two main components: a variable formatting layer and a fusion layer. The variable formatting layer processes each variable through linear transformations, converting them into standardized variables for subsequent analysis. The fusion layer then integrates these standardized variables along the feature dimension,

**Fig. 3 Fusion embedding modules**

producing a unified hidden representation that aligns with the model input specifications.

Consider the target variable  $\mathbf{V}^{\text{tgt}}$  and a set of  $p$  variables  $V = \{\mathbf{V}^1, \mathbf{V}^2, \dots, \mathbf{V}^p\}$ . The target variable  $\mathbf{V}^{\text{tgt}} \in \mathbb{R}^{L^{\text{tgt}} \times S^{\text{tgt}} \times F^{\text{tgt}}}$  has  $L^{\text{tgt}}$  time steps,  $S^{\text{tgt}}$  spatial nodes, and  $F^{\text{tgt}}$  features. Each variable  $\mathbf{V}^i \in \mathbb{R}^{L^i \times S^i \times F^i \times \dots}$  is multidimensional, where  $L^i$  is the temporal length,  $S^i$  is the number of spatial nodes, and  $F^i$  is the number of features. The ellipsis (...) denotes additional higher dimensions. For simplicity, we define the temporal window length of static variables as 1, dynamic variables in the historical dataset as  $T$ , and dynamic variables in the prediction window as  $h$ . Therefore,  $L^{\text{tgt}}$  and  $L^i$  can take values from  $\{1, T, h\}$ .

The fusion embedding module processes each auxiliary data source  $\mathbf{V}^i$  through a three-stage pipeline to produce a standardized tensor  $\xi^i$  compatible with downstream models. The stages are as follows:

1. Initial flattening and dimension selection. First, each high-dimensional input tensor  $\mathbf{V}^i \in \mathbb{R}^{L^i \times S^i \times F^i \times \dots}$  is vectorized. This can be configured in two ways: (1) preserving the spatial dimension to yield a tensor  $\mathbf{V}'^i \in \mathbb{R}^{L^i \times S^i \times F^i}$ , or (2) collapsing the spatial dimension to yield a flat feature vector  $\mathbf{V}''^i \in \mathbb{R}^{F^i}$ .

2. Feature management and standardization. Next, the resulting feature dimension  $F'^i$  is projected to a target model dimension  $F_m$ . We provide two methods for this projection from  $\mathbf{V}'^i_f$  to  $\tilde{\mathbf{V}}_f^i$ , where  $f$  is the feature dimension for  $\mathbf{V}'^i$ : (1) FlattenCat conditionally applies a linear projection to reduce dimensionality as defined in Eq. (1); (2) TFT processes categorical and continuous features separately (using embeddings or linear layers) before projecting them

to the target dimension  $F_m$ , as detailed in Eqs. (2) and (3):

$$\tilde{\mathbf{V}}_f^i = \begin{cases} \text{FC}(\mathbf{V}'^i_f) \in \mathbb{R}^{F_m}, & \text{if } F'^i \geq F_m, \\ \mathbf{V}'^i_f \in \mathbb{R}^{F'^i}, & \text{if } F'^i < F_m, \end{cases} \quad (1)$$

$$\tilde{\mathbf{V}}_f^i \in \mathbb{R}^{F_m} = \text{FC}(\text{Flatten}(\text{Embed}(\mathbf{V}'^i_{f,\text{cat}}))), \quad (2)$$

$$\tilde{\mathbf{V}}_f^i \in \mathbb{R}^{F_m} = \text{FC}(\mathbf{V}''^i_{f,\text{cont}}), \quad (3)$$

where Flatten is matrix flattening, Embed is nn.Embedding, and FC is fully connected layers for DL.

3. Spatial dimension alignment (optional). Finally, if the spatial dimension  $S^i$  was preserved at the first stage, it is aligned to the target model's required spatial dimension  $S^{\text{tgt}}$  via a linear layer ( $\text{FC} : \mathbb{R}^{S^i} \rightarrow \mathbb{R}^{S^{\text{tgt}}}$ ). This step is essential for downstream spatiotemporal models (e.g., T-GCN) that require spatially consistent inputs, but it is bypassed for purely temporal models (e.g., LSTM).

In the fusion layer, the three modules diverge in their methodologies. Consider the feature dimension of the formatted variable  $\xi_f^i \in \mathbb{R}^{F_f^i}$ . FlattenCat extends the static values along the temporal axis and combines the formatted outputs from all variables to form a unified hidden representation  $\mathbf{H}_f = [(\xi_f^1)^T, (\xi_f^2)^T, \dots, (\xi_f^p)^T]^T \in \mathbb{R}^{F_H}$ , where  $F_H$  is the sum of the individual feature dimensions, calculated as  $F_H = \sum_{i=1}^p F_f^i = \sum_{i=1}^p \min(F^i, F_m)$ .

TFT processes formatted static variables using a static variable encoder (SVE), as shown in Fig. 4c. The SVE incorporates a variable selection network (VSN) and four gated residual networks (GRNs) to generate four static encodings:  $c_s$ ,  $c_e$ ,  $c_h$ , and  $c_c$ , as shown in Figs. 4a and 4d. Subsequently, TFT employs an additional VSN to process dynamic variables, including the target, and combines these

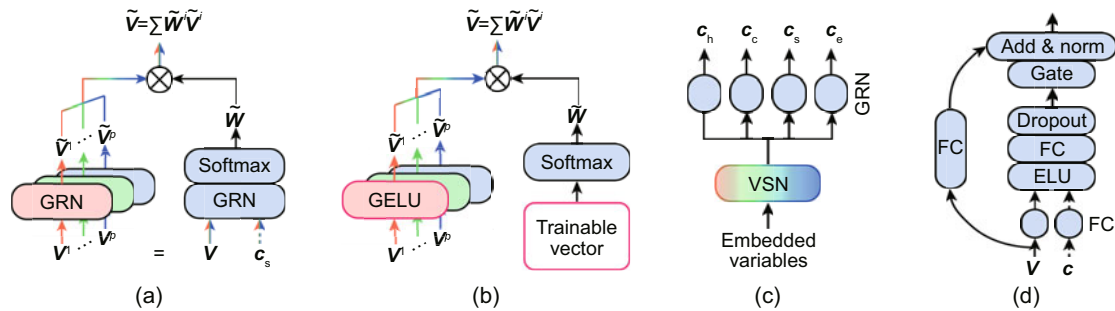


Fig. 4 The structures of VSN (a), GWS (b), SVE (c), and GRN (d). GWS: gated weighted summation; ELU: exponential linear unit

processed dynamic variables with the static encoding  $\mathbf{c}_s$ .

The VSN implementation within TFT, crucial for identifying salient variables, follows three main steps, whose logical relationship is key to its function:

1. Individual feature transformation. Each formatted variable  $\xi_f^i$  is individually processed by a GRN, which applies a gating mechanism to filter out irrelevant variables for the prediction task. That is, each GRN generates a transformed variable  $\tilde{\xi}_f^i \in \mathbb{R}^{F_m}$ , as shown in Eq. (4):

$$\tilde{\xi}_f^i = \text{GRN}(\xi_f^i, \mathbf{c}_s) \in \mathbb{R}^{F_m}. \quad (4)$$

2. Variable weight generation. To determine the relevance of each variable, the VSN computes dynamic attention weights. As described in Eq. (5), all original formatted input variables  $\{\xi_f^1, \xi_f^2, \dots, \xi_f^p\}$  are first stacked and flattened into a single vector  $\Xi_f \in \mathbb{R}^{p \times F_m}$ . This vector is then passed through another GRN with a context vector  $\mathbf{c}_s$  if providing static variables, and the output is fed into a Softmax function. This yields a set of dynamic weights  $\tilde{\mathbf{w}}^1, \tilde{\mathbf{w}}^2, \dots, \tilde{\mathbf{w}}^p$ , where each  $\tilde{\mathbf{w}}^i \in \mathbb{R}^p$  corresponds to the importance of the  $i^{\text{th}}$  variable. These weights are data-dependent, meaning that they are recalculated for each input instance.

$$\begin{aligned} \Xi_f &= \text{Flatten}([\xi_f^1, \xi_f^2, \dots, \xi_f^p]) \in \mathbb{R}^{p \times F_m}, \\ \tilde{\mathbf{W}} &= \text{Softmax}(\text{GRN}(\Xi_f)) \in \mathbb{R}^p. \end{aligned} \quad (5)$$

3. Weighted summation. Finally, the VSN produces its output  $\mathbf{H}_f$  by performing a weighted summation of the transformed features  $\tilde{\xi}_f^i$  (from Eq. (4)), using the generated variable weights  $\tilde{\mathbf{w}}^i$  (from Eq. (5)). This process, illustrated in Eq. (6), effectively selects and combines the most relevant features:

$$\mathbf{H}_f = \sum_{i=1}^p \tilde{\mathbf{w}}^i \tilde{\xi}_f^i. \quad (6)$$

In practice, we observe that TFT's VSN, particularly the GRNs within it, significantly contributes to memory usage, especially with a large number of variables. To address this inefficiency, we introduce the GWS module, as depicted in Fig. 4b. The GWS module modifies the VSN architecture for improved efficiency. The key distinctions between the GWS module and TFT's VSN are as follows:

1. Modification to individual feature transformation in Eq. (4). The computationally intensive

GRN is replaced with a simple gated unit, i.e., Gaussian error linear unit (GELU) activation function. The transformation becomes  $\tilde{\xi}_f^i = \text{GELU}(\xi_f^i)$ . This significantly reduces parameters and computational load compared to a full GRN.

2. Modification to variable weight generation in Eq. (5). Weight generation is simplified by making weights data-independent. It uses a learnable (trainable) parameter vector  $\mathbf{W} \in \mathbb{R}^p$ . The variable selection weights  $\tilde{\mathbf{W}}$  are then obtained by applying Softmax directly to this vector:  $\tilde{\mathbf{W}} = \text{Softmax}(\mathbf{W})$ . These weights are learned during the training process, and remain fixed thereafter for inference, rather than being re-computed for each input.

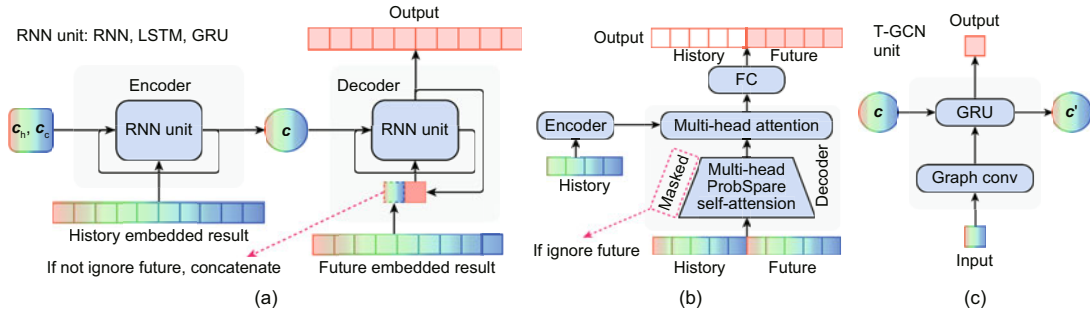
3. No change in weighted summation in Eq. (6).

### 3.6 Models with fusion embedding results

We investigate the integration of embedding modules in canonical time series and spatiotemporal prediction models, focusing on three key models. We select sequence to sequence (Seq2Seq) (Sutskever et al., 2014) as a representative of RNNs and Informer (Zhou et al., 2021) as an example of Transformer-based approaches. Given the graph-structured nature of our data, we choose temporal graph convolutional network (T-GCN) (Zhao et al., 2020) as the main spatiotemporal model and adapt its architecture like Seq2Seq for multi-step forecasting, as shown in Fig. 5c. We also choose Graph WaveNet (Wu ZH et al., 2019).

In our model design, we observe that the original Seq2Seq and Informer models intentionally disregard or occlude variables within the future window. This design choice stems from their primary use in forecasting a single unknown variable, for which the true value in the future window is unavailable. However, in traffic-state prediction tasks that involve multiple auxiliary variables, covariates from various data sources within the future window are accessible and can potentially improve predictive accuracy. To address this, we refine the model to better use information from the future window as shown in Fig. 5.

Specifically, we use both the hidden state and the output from the embedding module as the input to the decoder of Seq2Seq instead of only the hidden state. For Informer, we remove the mask in the decoder's first layer, allowing it to receive information from the future window. Additionally, when extracting samples from the dataset for the future window,



**Fig. 5 Fusion embedding modules and models: (a) Seq2Seq; (b) Informer; (c) T-GCN unit. Conv: convolution**

we assign a value of 0 to unknown variables to simulate the uncertainty encountered in the real-world scenarios.

## 4 Experiments

### 4.1 Experimental setup

We conduct experiments on the augmented PEMS-BAY and METR-LA datasets detailed in Section 3.2. All input variables are normalized using  $Z$ -score, and each dataset is divided into training, validation, and test sets using a 7:2:1 ratio. Model performance is evaluated using mean absolute error (MAE) and mean squared error (MSE).

All experiments are conducted on a server equipped with an NVIDIA RTX 4090 GPU (24 GB) and an Intel Xeon Gold 6430 CPU, running Ubuntu 22.04. Our framework is built with Python 3.12.3 and PyTorch 2.3.0 (CUDA 12.1). To enhance training efficiency, we use automatic mixed precision (AMP) via the torch.amp package.

Our experimental process consists of two main stages: parameter tuning and final evaluation. To ensure a fair comparison, we first tune hyperparameters for baseline models (Seq2Seq, Informer, and T-GCN). At this stage, models are trained for up to three epochs using only traffic variables as input, consistent with their original designs. We use the Adam optimizer with a learning rate of 0.001 and a batch size of 512, and employ early stopping to prevent overfitting. For the main evaluation, we assess model performance under various conditions summarized in Table 5. To ensure robustness, each experiment is performed 10 times with different random seeds. Unless otherwise specified, all models are trained for up to 15 epochs using the Adam op-

timizer (with learning rate 0.001 and batch size 512) with early stopping. Minor adjustments to batch size and AMP usage are made on a case-by-case basis to prevent out-of-memory errors.

### 4.2 Performance and resource utilization of fusion embedding modules

In Section 3.5, we introduce three fusion embedding modules: FlattenCat, TFT, and TFT-GWS. To evaluate their performance and resource utilization, we measure the time and memory costs of each module during both inference and training, as detailed in Table 6. Our analysis focuses on varying sequence lengths, spatial dimensions, and the number of external features to provide a comprehensive understanding of their performance characteristics. We do not reiterate the computational overhead of subsequent models (e.g., Informer and T-GCN) because this has been extensively covered in previous research (Shao et al., 2025). When running the experiments, each module receives multiple time-series variables as the input, spanning historical and future windows of length  $L_w$ . Specifically, each variable has a shape of  $(16, 2L_w, N_s, 1)$ , encompassing one target component and  $N_{Co}$  variables for each covariate component as detailed in Table 4. When the retaining spatial dimension (RSD) is enabled, the module outputs a variable of shape  $(16, 2L_w, N_s, 512)$ ; otherwise, the output shape is  $(16, 2L_w, 512)$ .

### 4.3 Basic performance of forecasting

The baseline model performance on the PEMS-BAY and METR-LR datasets using only traffic data is presented in Table 7 (more in the supplementary materials). The analysis highlights four primary findings: (1) The prediction error increases as

**Table 5 Experimental conditions for model evaluation**

Condition	Setting
Number of nodes ( $N_s$ )	1-1, 4-1, 16- $x$ , all ( $n-i$ denotes $n$ nodes with index intervals $i$ , and $x \in \{1, 2, 4, 8\}$ )
Time window ( $L_w$ )	$L_w \in \{12, 24, 48, 96\}$ (input and output windows are identical)
Auxiliary variables (Co)	Traffic-only (Co = No); All auxiliary variables (Co = All)
Embedding modules	FlattenCat, TFT, TFT-GWS
Model perception of future window	With or without future window information
Model after fusion embedding	Seq2Seq, Informer, T-GCN, graph WaveNet

the length of the forecast horizon increases; (2) The model performance is highly sensitive to the number of nodes used, with different architectures responding differently to the network scale; (3) The choice of embedding module is critical because nonlinear embeddings provide a significant advantage for long-term forecasting over simpler methods; (4) As the number of nodes increases, the complex spatial dependencies can act as noise, and the resulting modeling challenges can overwhelm the informational benefits, ultimately harming predictive accuracy.

#### 4.4 Improvement with covariates

We evaluate the impact of covariates by comparing prediction performance with and without them under identical experimental configurations (Section 4.3). As shown in Table 8, covariates significantly improve the accuracy for long-term forecasts and in low-complexity, single-node settings, highlighting their value for extended prediction horizons.

The impact of covariates varies with prediction horizon and node complexity. For short horizons (e.g., 12 and 24 steps), covariates offer minimal or even detrimental effects. This is likely due to the sampling frequency mismatch between high-frequency traffic data and low-frequency covariates like weather data. Conversely, their benefits are more pronounced for longer horizons (e.g., 48 and 96 steps), where their inclusion consistently enhances the model performance. The positive impact of covariates is most significant in single-node tasks. In multi-node scenarios, however, models can still effectively leverage covariate information by capturing interactions across different nodes, leading to performance gains.

The FlattenCat method is less stable and effective than TFT and TFT-GWS, especially in short-term and full-node scenarios. FlattenCat exhibits inconsistent performance; although effective in some long-term, low-node-count tasks, it shows marked

instability in short-term and full-node scenarios. Specifically, it causes significant performance degradation in short-term predictions (12 and 24 steps), and leads to training anomalies in full-node configurations. In contrast, TFT and TFT-GWS consistently outperform FlattenCat and remain stable across most scenarios. Their advantage is particularly notable in single-node configurations. Although their performance gains are less pronounced in full-node tasks, they avoid the anomalous behavior seen with FlattenCat. It is worth noting that TFT-GWS occasionally underperforms TFT, particularly in short-term predictions, which we attribute to the potential inefficiencies in its gating mechanism.

#### 4.5 Impact of covariates in the future window

We conduct a comparative analysis of prediction tasks that consider future window information with those that do not. The analysis in Table 9 shows that incorporating future window information, particularly with known covariates, significantly improves prediction performance.

TFT and TFT-GWS benefit from future covariates, whereas FlattenCat struggles, often showing adverse effects due to weak covariate-target correlations. When using the TFT and TFT-GWS embedding modules, known covariates in the future window generally improve prediction accuracy, enabling the model to more effectively capture the evolving patterns of the target data. However, for the FlattenCat module, the inclusion of known covariates in the future window exhibits minimal impact on the prediction performance, and in some instances, it even adversely affects the performance, especially in the Informer model. This can be attributed to covariates with low correlation to the target data, which tend to interfere more with the model than provide valuable information. Future window information offers minimal benefits without covariates, and may introduce biases that hinder prediction performance.

**Table 6 Performance and resource utilization of fusion embedding modules**

Embedding module	$N_s$	RSD	Number of parameters							
			$N_{Co}=0$		1		4			
			$L_w=6$	48	6	48	6	48		
FlattenCat	8	False	0	0	0	0	0	0	0	0
		True	0	0	432	432	$1.7 \times 10^3$	$1.7 \times 10^3$		
	16	False	0	0	0	0	0	0		
		True	0	0	$1.6 \times 10^3$	$1.6 \times 10^3$	$6.5 \times 10^3$	$6.5 \times 10^3$		
TFT	8	False	$2.1 \times 10^6$	$2.1 \times 10^6$	$17.3 \times 10^6$	$17.3 \times 10^6$	$50.7 \times 10^6$	$50.7 \times 10^6$		
		True	$2.1 \times 10^6$	$2.1 \times 10^6$	$16.9 \times 10^6$	$16.9 \times 10^6$	$49.3 \times 10^6$	$49.3 \times 10^6$		
	16	False	$2.1 \times 10^6$	$2.1 \times 10^6$	$17.7 \times 10^6$	$17.7 \times 10^6$	$52.3 \times 10^6$	$52.3 \times 10^6$		
		True	$2.1 \times 10^6$	$2.1 \times 10^6$	$16.9 \times 10^6$	$16.9 \times 10^6$	$49.3 \times 10^6$	$49.3 \times 10^6$		
TFT-GWS	8	False	$4.6 \times 10^3$	$4.6 \times 10^3$	$6.7 \times 10^6$	$6.7 \times 10^6$	$14.3 \times 10^6$	$14.3 \times 10^6$		
		True	$1.0 \times 10^3$	$1.0 \times 10^3$	$6.4 \times 10^6$	$6.4 \times 10^6$	$12.9 \times 10^6$	$12.9 \times 10^6$		
	16	False	$8.7 \times 10^3$	$8.7 \times 10^3$	$7.1 \times 10^6$	$7.1 \times 10^6$	$16.0 \times 10^6$	$16.0 \times 10^6$		
		True	$1.0 \times 10^3$	$1.0 \times 10^3$	$6.4 \times 10^6$	$6.4 \times 10^6$	$12.9 \times 10^6$	$12.9 \times 10^6$		

Embedding module	$N_s$	RSD	Inference											
			Time (ms)						Memory (MB)					
			$N_{Co}=0$		1		4		$N_{Co}=0$		1		4	
			$L_w=6$	48	6	48	6	48	6	48	6	48	6	48
FlattenCat	8	False	0.1	0.1	0.4	0.4	0.7	0.7	0.0	0.2	0.2	1.3	0.6	4.6
		True	0.1	0.1	0.6	0.6	1.6	1.5	0.0	0.2	1.1	1.9	1.4	4.6
	16	False	0.1	0.1	0.4	0.4	0.7	0.7	0.1	0.5	0.3	2.6	1.1	9.9
		True	0.1	0.1	0.7	0.6	1.5	1.4	0.1	0.5	1.2	2.9	1.8	9.9
TFT	8	False	1.2	1.2	4.4	5.3	13.3	10.4	2.7	16.6	8.9	62.4	26.6	208.8
		True	1.2	2.0	4.8	11.6	14.6	40.9	16.5	132.1	63.3	482.1	209.5	1637.0
	16	False	1.0	1.0	5.4	5.4	13.4	11.9	2.7	16.7	9.0	62.9	26.8	210.4
		True	1.1	3.7	6.1	21.8	14.3	81.3	33.0	264.3	126.9	964.2	418.2	3274.0
TFT-GWS	8	False	0.3	0.4	2.2	2.8	4.8	4.7	1.8	12.1	9.2	67.9	29.9	232.7
		True	0.4	0.6	3.3	5.6	5.9	22.1	12.0	96.1	69.3	530.0	233.4	1828.6
	16	False	0.4	0.4	2.8	2.8	4.7	5.6	1.8	12.2	9.3	68.4	30.1	234.3
		True	0.4	1.4	4.0	11.3	7.4	46.9	24.0	192.2	138.9	1059.9	466.1	3657.2

Embedding module	$N_s$	RSD	Training											
			Time (ms)						Memory (MB)					
			$N_{Co}=0$		1		4		$N_{Co}=0$		1		4	
			$L_w=6$	48	6	48	6	48	6	48	6	48	6	48
FlattenCat	8	False												
		True			2.5	2.7	5.5	5.7			1.1	2.7	1.5	10.2
	16	False												
		True			2.6	2.4	6.0	5.3			1.3	5.4	2.5	20.9
TFT	8	False	5.7	6.4	21.3	20.0	53.7	48.3	9.2	32.8	55.9	121.7	183.4	409.1
		True	6.1	6.4	20.9	35.4	60.2	146.0	32.3	324.4	151.8	1462.6	510.5	5049.4
	16	False	6.6	5.5	20.0	23.1	54.8	51.4	9.2	32.9	57.5	126.1	189.8	428.6
		True	5.7	11.6	24.2	67.4	56.2	283.3	76.6	654.9	351.4	2976.4	1179.2	1 0273.7
TFT-GWS	8	False	1.9	2.1	6.2	7.8	15.3	17.8	3.4	27.1	11.3	87.3	34.3	293.8
		True	1.8	2.1	7.8	17.8	20.4	97.5	27.0	216.1	117.4	840.8	396.2	2961.1
	16	False	1.9	2.1	7.2	7.5	16.9	15.8	3.4	27.2	10.5	91.2	31.3	313.3
		True	2.2	4.4	9.9	35.6	28.8	197.3	54.0	432.2	235.2	1680.3	788.5	5921.0

$N_{Co}$ : number of variables per covariate component in Table 4 (0 for no covariate, 1 for 6 covariates, and 4 for 24 covariates). RSD: retaining spatial dimension (false if not retained, true if retained)

Table 7 Forecast performance without covariates or future information in the PEMS-BAY dataset

Model	Embedding module	$N_s$	Forecast performance								
			$L_w=12$		24		48		96		
			MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
Seq2Seq	FlattenCat	1-1	0.2967	0.2702	0.3750	0.3148	0.5221	0.3949	0.6479	0.4398	
		16-1	0.2864	0.2681	0.4540	0.3360	0.5574	0.3812	0.5541	0.3880	
		16-4	0.1875	0.2050	0.2624	0.2421	0.3250	0.2810	0.3782	0.3083	
		All	<b>0.2841</b>	<b>0.2762</b>	<b>0.3178</b>	<b>0.2905</b>	<b>0.3810</b>	<b>0.3281</b>	0.3867	0.3271	
	TFT	1-1	0.3065	0.2673	0.3663	0.3105	0.5422	0.4024	0.6146	0.4245	
		16-1	0.3067	0.2766	0.4266	0.3234	0.5638	0.3879	0.5541	0.3851	
		16-4	0.1964	0.2065	0.2645	0.2430	0.3378	0.2882	0.3873	0.3132	
		All	0.2890	0.2780	0.3194	0.2917	0.3817	0.3296	<b>0.3803</b>	<b>0.3224</b>	
	TFT-GWS	1-1	0.2910	0.2629	0.3645	0.3078	<b>0.5069</b>	0.3850	<b>0.5767</b>	<b>0.3964</b>	
		16-1	0.3096	0.2772	0.4507	0.3295	0.5516	0.3833	<b>0.5266</b>	<b>0.3687</b>	
		16-4	0.1974	0.2094	0.2602	0.2418	<b>0.3217</b>	<b>0.2778</b>	<b>0.3698</b>	<b>0.3055</b>	
		All	0.2916	0.2767	0.3216	0.2912	0.3834	0.3302	0.3857	0.3247	
Informer	FlattenCat	1-1	0.3140	0.2770	0.4252	0.3352	0.5931	0.4249	0.6496	0.4424	
		16-1	0.3215	0.2936	0.4859	0.3569	0.5775	0.3990	0.5835	0.4030	
		16-4	0.2136	0.2234	0.2885	0.2605	0.3921	0.3184	0.4073	0.3359	
		All	0.3061	0.2908	0.3481	0.3118	0.3901	0.3342	0.3847	0.3300	
	TFT	1-1	0.3267	0.2879	0.4361	0.3421	0.5885	0.4160	0.6416	0.4393	
		16-1	0.3509	0.3066	0.4943	0.3543	0.5993	0.4080	0.5940	0.4074	
		16-4	0.2205	0.2275	0.2880	0.2587	0.4117	0.3247	0.4179	0.3412	
		All	0.3106	0.2914	0.3386	0.3043	0.4220	0.3494	0.3841	0.3318	
	TFT-GWS	1-1	0.3120	0.2775	0.4245	0.3369	0.5708	0.4161	0.6294	0.4319	
		16-1	0.4076	0.3147	0.5071	0.3612	0.5663	0.3951	0.5714	0.3950	
		16-4	0.2242	0.2321	0.2902	0.2618	0.3488	0.2950	0.4020	0.3283	
		All	0.3175	0.2964	0.3464	0.3120	0.3912	0.3357	0.3814	0.3291	
T-GCN	FlattenCat	1-1	0.3035	0.2777	0.3730	0.2959	0.5411	0.3951	0.7115	0.4678	
		16-1	0.2880	0.2677	0.4683	0.3677	0.7011	0.4772	0.9316	0.5682	
		16-4	0.2454	0.2265	0.4009	0.3071	0.6237	0.4003	0.8116	0.4885	
	TFT	1-1	0.2892	0.2667	<b>0.3531</b>	<b>0.2927</b>	0.5205	0.3748	0.6151	0.3995	
		16-1	0.2614	0.2534	0.4213	0.3456	0.6394	0.4411	0.8651	0.5424	
		16-4	0.2242	0.2153	0.3771	0.2894	0.5850	0.3710	0.7847	0.4747	
	TFT-GWS	1-1	<b>0.2839</b>	<b>0.2586</b>	0.3554	0.2946	0.5209	<b>0.3731</b>	0.6579	0.4447	
		16-1	0.2636	0.2512	0.4308	0.3482	0.6521	0.4644	0.8797	0.5495	
		16-4	0.2313	0.2159	0.3867	0.2943	0.5994	0.3922	0.7958	0.4853	
	Graph WaveNet	FlattenCat	1-1	0.3002	0.2633	0.4449	0.3372	0.5529	0.4001	0.6642	0.4382
			16-1	<b>0.2299</b>	<b>0.2283</b>	<b>0.3449</b>	<b>0.2892</b>	<b>0.4979</b>	<b>0.3468</b>	0.5751	0.4011
			16-4	<b>0.1794</b>	<b>0.1897</b>	<b>0.2383</b>	0.2242	0.3592	0.2845	0.4367	0.3344
TFT		1-1	0.3278	0.2889	0.4142	0.3178	0.5237	0.3938	0.6281	0.4163	
		16-1	0.2386	0.2331	0.3579	0.2897	0.5262	0.3614	0.5566	0.3867	
		16-4	0.1933	0.1943	0.2422	<b>0.2222</b>	0.3740	0.2920	0.4395	0.3277	
TFT-GWS		1-1	0.3274	0.3075	0.4355	0.3341	0.5874	0.4215	0.6572	0.4244	
		16-1	0.2493	0.2564	0.3869	0.3133	0.5337	0.3673	0.5710	0.3948	
		16-4	0.1998	0.2134	0.2604	0.2391	0.3676	0.2919	0.4366	0.3292	

The best results within the same window and node configuration are in bold

**Table 8** Relative forecast performance gain with covariates in the PEMS-BAY dataset

Model	Embedding module	$N_s$	Forecast performance gain							
			$L_w=12$		24		48		96	
			MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Seq2Seq	FlattenCat	1-1	<b>2.28</b>	-5.56	<b>0.08</b>	-3.05	<b>4.96</b>	<b>0.71</b>	<b>15.99</b>	-0.55
		16-1	-46.26	-16.85	-7.69	-3.56	<b>2.56</b>	<b>0.72</b>	<b>5.88</b>	<b>4.98</b>
		All	-275.81	-110.98	-237.09	-100.99	-183.22	-78.76	-176.47	-78.77
	TFT	1-1	<b>13.21</b>	<b>5.88</b>	<b>5.79</b>	<b>3.34</b>	<b>31.39</b>	<b>17.95</b>	<b>43.00</b>	<b>22.74</b>
		16-1	<b>0.19</b>	-0.00	-3.23	<b>0.09</b>	<b>13.05</b>	<b>10.14</b>	<b>13.26</b>	<b>11.88</b>
		All	<b>0.42</b>	<b>0.66</b>	<b>3.15</b>	<b>1.95</b>	<b>10.96</b>	<b>6.96</b>	<b>10.13</b>	<b>5.29</b>
	TFT-GWS	1-1	<b>5.59</b>	<b>0.46</b>	<b>5.41</b>	-1.08	<b>24.07</b>	<b>9.69</b>	<b>36.15</b>	<b>13.33</b>
		16-1	-41.74	-14.29	-6.95	-2.22	<b>9.05</b>	<b>8.04</b>	<b>7.49</b>	<b>5.98</b>
		All	-1.80	-1.88	<b>0.61</b>	-0.63	<b>9.75</b>	<b>6.63</b>	<b>10.18</b>	<b>3.89</b>
Informer	FlattenCat	1-1	-3.83	-10.12	-0.50	-4.46	<b>14.35</b>	<b>6.69</b>	<b>22.63</b>	<b>3.90</b>
		16-1	-38.03	-11.90	-2.84	-0.91	<b>7.38</b>	<b>5.04</b>	<b>7.43</b>	<b>6.44</b>
		All	-40.14	-23.43	-50.60	-29.83	-139.86	-61.33	-46.46	-24.30
	TFT	1-1	<b>10.85</b>	<b>4.36</b>	<b>10.34</b>	<b>2.25</b>	<b>32.53</b>	<b>16.46</b>	<b>37.18</b>	<b>18.11</b>
		16-1	-10.36	-1.20	<b>6.61</b>	<b>4.74</b>	<b>19.64</b>	<b>15.00</b>	<b>16.31</b>	<b>12.79</b>
		All	<b>4.56</b>	<b>2.50</b>	<b>3.90</b>	<b>1.64</b>	<b>19.23</b>	<b>11.53</b>	<b>8.26</b>	<b>4.43</b>
	TFT-GWS	1-1	<b>5.68</b>	-1.34	<b>17.03</b>	<b>4.96</b>	<b>35.30</b>	<b>20.55</b>	<b>35.88</b>	<b>15.16</b>
		16-1	-10.08	-3.52	<b>7.18</b>	<b>7.18</b>	<b>16.55</b>	<b>14.01</b>	<b>14.10</b>	<b>11.91</b>
		All	<b>3.35</b>	<b>2.00</b>	<b>3.34</b>	<b>2.17</b>	<b>9.08</b>	<b>5.95</b>	<b>4.01</b>	<b>1.61</b>
T-GCN	FlattenCat	1-1	<b>7.69</b>	<b>1.15</b>	<b>7.95</b>	-4.44	<b>26.67</b>	<b>13.73</b>	<b>39.14</b>	<b>22.71</b>
		16-1	<b>5.79</b>	-3.43	<b>8.48</b>	<b>0.57</b>	<b>17.87</b>	<b>10.78</b>	<b>23.99</b>	<b>16.52</b>
		All	<b>18.83</b>	<b>9.13</b>	<b>10.08</b>	<b>1.13</b>	<b>30.54</b>	<b>14.73</b>	<b>35.22</b>	<b>14.92</b>
	TFT	1-1	<b>11.60</b>	<b>4.64</b>	<b>21.70</b>	<b>13.58</b>	<b>30.42</b>	<b>19.21</b>	<b>40.49</b>	<b>30.26</b>
		16-1	<b>11.60</b>	<b>4.64</b>	<b>21.70</b>	<b>13.58</b>	<b>30.42</b>	<b>19.21</b>	<b>40.49</b>	<b>30.26</b>
		All	<b>11.60</b>	<b>4.64</b>	<b>21.70</b>	<b>13.58</b>	<b>30.42</b>	<b>19.21</b>	<b>40.49</b>	<b>30.26</b>
	TFT-GWS	1-1	-1.76	-9.90	<b>3.40</b>	-8.58	<b>23.22</b>	<b>7.08</b>	<b>39.22</b>	<b>23.18</b>
		16-1	<b>0.63</b>	-8.72	<b>5.89</b>	-0.07	<b>14.80</b>	<b>9.95</b>	<b>19.05</b>	<b>12.80</b>
		All	<b>0.63</b>	-8.72	<b>5.89</b>	-0.07	<b>14.80</b>	<b>9.95</b>	<b>19.05</b>	<b>12.80</b>
Graph WaveNet	FlattenCat	1-1	<b>5.46</b>	-3.11	<b>13.89</b>	<b>2.16</b>	<b>12.18</b>	<b>8.38</b>	<b>15.27</b>	<b>5.07</b>
		16-1	-0.39	-2.02	<b>1.09</b>	<b>1.73</b>	<b>12.60</b>	<b>6.67</b>	<b>15.59</b>	<b>11.83</b>
		All	<b>16.16</b>	<b>4.43</b>	-14.48	-16.44	-6.30	-3.58	-27.52	-19.43
	TFT	1-1	<b>16.16</b>	<b>4.43</b>	-14.48	-16.44	-6.30	-3.58	-27.52	-19.43
		16-1	<b>2.12</b>	<b>0.07</b>	<b>6.36</b>	<b>0.58</b>	<b>12.75</b>	<b>5.66</b>	<b>2.13</b>	<b>1.67</b>
		All	<b>2.12</b>	<b>0.07</b>	<b>6.36</b>	<b>0.58</b>	<b>12.75</b>	<b>5.66</b>	<b>2.13</b>	<b>1.67</b>
	TFT-GWS	1-1	<b>0.20</b>	-2.42	<b>11.75</b>	<b>4.75</b>	<b>7.58</b>	<b>10.61</b>	-40.02	-18.22
		16-1	-19.56	-7.75	-3.10	-1.76	<b>14.42</b>	<b>5.92</b>	<b>3.77</b>	<b>3.17</b>
		All	-19.56	-7.75	-3.10	-1.76	<b>14.42</b>	<b>5.92</b>	<b>3.77</b>	<b>3.17</b>

The percentage improvement is calculated as  $v = -(m_w/ - m_{w/o})/m_{w/o} \times 100\%$ , where  $m_w/$  and  $m_{w/o}$  are metrics with and without future covariates, respectively. Positive improvements are in bold

Moreover, we observe that the benefits of future window information are marginal in the absence of covariates, and it even results in a decline in the prediction performance in some cases. This can be attributed to biases inherent in the embedding module, which interfere with the ability of the model to learn the target data.

#### 4.6 Variable importance

To understand variable importance, we first note that simple statistical measures, such as Pear-

son correlation and maximum information coefficient (MIC) (Reshef et al., 2011; Albanese et al., 2013), reveal weak relationships between traffic speed and most auxiliary variables (Fig. 6), except for the MIC between traffic speed and HoD or DoW. However, analyzing the learned variable selection weights from our fusion modules reveals a more nuanced and interpretable picture, as shown in Fig. 7.

The model correctly assigns the highest importance to the historical target variable. More revealingly, it identifies the key predictive features that

**Table 9** Relative forecast performance gain with future information in the PEMS-BAY dataset

Co	Model	Embedding module	Forecast performance gain							
			$L_w=12$		24		48		96	
			MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
All	Seq2Seq	FlattenCat	<b>2.22</b>	<b>1.21</b>	<b>2.79</b>	<b>1.49</b>	<b>3.41</b>	<b>1.81</b>	<b>9.16</b>	<b>6.74</b>
		TFT	<b>9.19</b>	<b>5.52</b>	<b>12.85</b>	<b>8.18</b>	<b>14.27</b>	<b>11.23</b>	<b>10.29</b>	<b>9.50</b>
		TFT-GWS	<b>4.87</b>	<b>3.09</b>	<b>10.46</b>	<b>6.61</b>	<b>12.93</b>	<b>10.17</b>	<b>13.41</b>	<b>11.30</b>
	Informer	FlattenCat	-0.34	-0.17	<b>3.57</b>	<b>1.75</b>	<b>4.50</b>	<b>2.21</b>	-11.24	-4.55
		TFT	<b>7.95</b>	<b>3.84</b>	<b>4.85</b>	<b>3.74</b>	<b>6.42</b>	<b>5.43</b>	<b>12.00</b>	<b>10.14</b>
		TFT-GWS	<b>0.46</b>	<b>0.21</b>	<b>1.50</b>	<b>1.82</b>	<b>2.53</b>	<b>3.17</b>	<b>10.40</b>	<b>8.56</b>
	T-GCN	FlattenCat	<b>4.78</b>	<b>3.51</b>	<b>5.26</b>	<b>3.64</b>	<b>8.93</b>	<b>6.46</b>	<b>19.95</b>	<b>12.49</b>
		TFT	<b>10.72</b>	<b>6.73</b>	<b>15.04</b>	<b>10.73</b>	<b>20.22</b>	<b>15.94</b>	<b>21.89</b>	<b>16.90</b>
		TFT-GWS	<b>8.76</b>	<b>5.79</b>	<b>15.02</b>	<b>10.58</b>	<b>21.64</b>	<b>16.38</b>	<b>31.37</b>	<b>22.21</b>
	Graph WaveNet	FlattenCat	<b>5.26</b>	<b>0.07</b>	<b>3.19</b>	-0.70	<b>0.26</b>	-3.09	-9.11	-3.62
		TFT	-3.44	-1.91	<b>4.91</b>	-0.36	<b>7.41</b>	<b>1.71</b>	<b>15.75</b>	<b>10.76</b>
		TFT-GWS	<b>7.14</b>	<b>4.26</b>	<b>1.56</b>	-1.10	<b>5.67</b>	<b>1.02</b>	<b>11.21</b>	<b>6.31</b>
No	Seq2Seq	FlattenCat	-0.61	-0.18	-0.10	-0.09	<b>0.78</b>	<b>0.18</b>	<b>1.32</b>	<b>1.22</b>
		TFT	-1.28	-1.87	-0.98	-1.00	<b>3.35</b>	<b>2.20</b>	<b>1.58</b>	<b>2.01</b>
		TFT-GWS	-6.10	-3.16	-2.23	-1.58	<b>0.44</b>	<b>0.30</b>	-2.64	-1.46
	Informer	FlattenCat	<b>1.19</b>	<b>0.65</b>	<b>0.17</b>	-0.24	-0.54	-0.62	-0.64	-0.53
		TFT	<b>3.38</b>	<b>1.43</b>	-2.02	-1.76	-5.73	-4.23	-5.53	-4.37
		TFT-GWS	<b>4.80</b>	<b>1.99</b>	<b>0.52</b>	<b>0.03</b>	-2.16	-1.29	-1.00	-1.32
	T-GCN	FlattenCat	<b>0.02</b>	<b>0.07</b>	-0.71	-0.79	-0.07	<b>0.12</b>	<b>0.89</b>	<b>0.14</b>
		TFT	-0.69	-0.39	<b>1.43</b>	<b>1.01</b>	<b>0.82</b>	-0.57	<b>4.60</b>	<b>1.07</b>
		TFT-GWS	<b>0.02</b>	<b>0.86</b>	<b>0.03</b>	<b>0.66</b>	-0.53	-0.76	<b>4.99</b>	<b>3.03</b>
	Graph WaveNet	FlattenCat	-4.25	-5.32	-1.61	-1.86	<b>0.22</b>	-0.07	<b>0.53</b>	<b>0.71</b>
		TFT	<b>2.11</b>	-0.93	<b>1.40</b>	-3.88	-2.56	-1.57	-0.64	-0.74
		TFT-GWS	-0.84	<b>2.83</b>	<b>3.00</b>	<b>0.44</b>	<b>1.54</b>	<b>0.04</b>	-2.64	-3.20

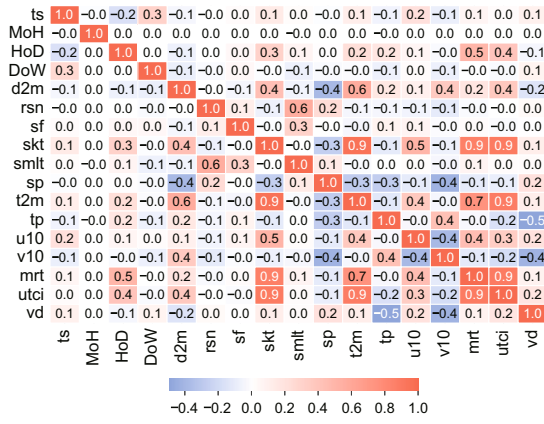
The percentage improvement is calculated as  $v = -(m_{w/} - m_{w/o})/m_{w/o} \times 100\%$ , where  $m_{w/}$  and  $m_{w/o}$  are metrics with and without future information, respectively. Positive improvements are in bold

simple correlations miss. For instance, it gives high weight to future HoD and DoW and to historical precipitation (tp) and visibility (vd). The importance of these variables is empirically supported: Traffic patterns are strongly tied to DoW (Fig. 8), and adverse weather is known to impact driving speed (FHWA, 2024). This demonstrates that the model learns meaningful, nonlinear relationships beyond simple statistical correlation. Furthermore, these weights allow for a comparative analysis of the embedding modules. The weights in the TFT-GWS module exhibit a smaller dynamic range than those in the TFT module. This suggests that TFT is more selective in focusing on critical variables, whereas TFT-GWS gives more even consideration to all inputs, which could make it more susceptible to noise from irrelevant features and negatively impact performance.

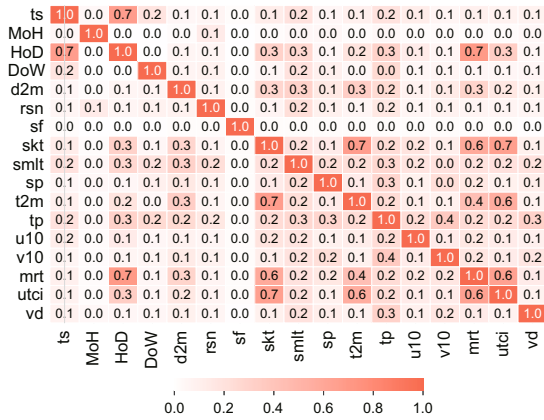
## 5 Conclusions

In this study, we introduce MltAuxTSPP, an innovative benchmarking framework designed to en-

hance traffic state prediction by integrating multi-source auxiliary data. Through a UDC and adaptable fusion embedding modules, MltAuxTSPP effectively addresses key challenges in data utilization and scalability. Our experimental results demonstrate that incorporating auxiliary data, particularly weather and temporal features, significantly improves long-term forecast accuracy. This work validates the efficacy of fusion embedding modules (e.g., FlattenCat, TFT, and TFT-GWS) in transforming heterogeneous data into a unified, predictive feature space. The comparative analysis of these methods offers valuable insights for model selection. Ultimately, MltAuxTSPP provides a robust and extensible framework that not only advances prediction accuracy but also serves as a critical tool for developing more resilient and efficient strategies for ITS. Although this study shows promising results, it has several limitations. First, computational overhead increases significantly with large-scale, multi-source data. Second, the model's sensitivity to the quality of auxiliary data poses a challenge. Although



(a)

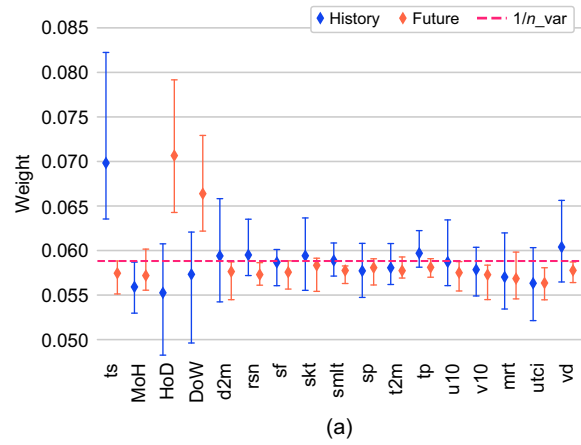


(b)

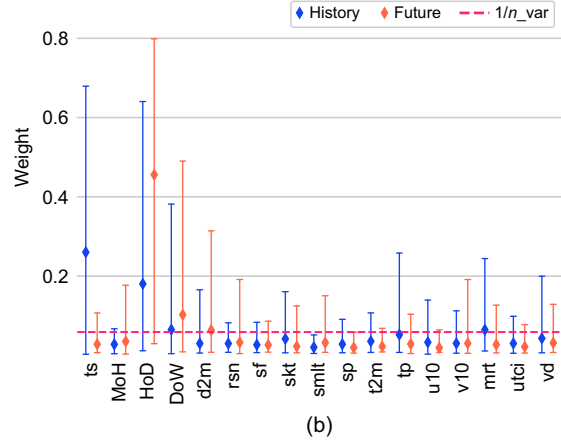
Fig. 6 Correlation between traffic state and covariates: (a) Pearson; (b) MIC. ts: traffic speed; d2m: 2-m dewpoint temperature; rsn: mass of snow per cubic meter in the snow layer; sf: accumulated total snow that has fallen to the Earth’s surface; skt: skin temperature; smlt: melting of snow averaged over the grid box; sp: surface pressure; t2m: 2-m temperature; tp: total precipitation; u10: 10-m u-component of wind; v10: 10-m v-component of wind; mrt: mean radiant temperature; utci: universal thermal climate index; vd: visibility distance

our experimental results demonstrate an improved performance, longer prediction horizons are needed to fully assess the influence of external factors such as weather and special events. Furthermore, the current dataset lacks examples of extreme weather conditions, and the model architecture is limited to classical deep learning approaches.

Future work will address these limitations by focusing on more efficient data processing techniques and adaptive modeling strategies. We also plan to evaluate our method on larger datasets, such as LargeST (Liu X et al., 2023), with extended time horizons and an increased number of traffic



(a)



(b)

Fig. 7 Variable selection weights in Informer: (a) with TFT-GWS; (b) with TFT

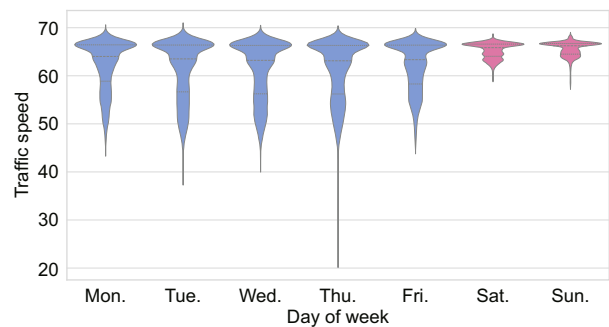


Fig. 8 Average traffic speed distribution by DoW on PEMS-BAY

nodes, building upon related research in traffic patterns (Shao et al., 2025; Zheng et al., 2025). Additionally, we recognize the importance of collecting data on extreme weather events to ensure the algorithm robustness in real-world scenarios.

### Contributors

Yusong ZHOU conceptualized the research, performed the data collection and analysis, conducted the experiments,

and drafted the paper. Xiaoyu JIANG and Shu SUN provided valuable suggestions during the conceptualization phase. Xinmin ZHANG, Yuanqiu MO, and Zhihuan SONG performed critical review and editing of the paper, and helped with its finalization.

### Conflict of interest

All the authors declare that they have no conflict of interest.

### Data availability

All the data that support the findings of this study are available through the open-source LibCity project. Datasets are already uploaded in network disks BaiduDisk with code 1231 or Google Drive.

### References

- Akiba T, Sano S, Yanase T, et al., 2019. Optuna: a next-generation hyperparameter optimization framework. *Proc 25<sup>th</sup> ACM SIGKDD Int Conf on Knowledge Discovery & Data Mining*, p.2623-2631. <https://doi.org/10.1145/3292500.3330701>
- Albanese D, Filosi M, Visintainer R, et al., 2013. Minerva and minepy: a C engine for the MINE suite and its R, Python and MATLAB wrappers. *Bioinformatics*, 29(3):407-408. <https://doi.org/10.1093/bioinformatics/bts707>
- Bączek J, Zhylyk D, Titericz G, et al., 2024. TSPP: a unified benchmarking tool for time-series forecasting. <https://doi.org/10.48550/arXiv.2312.17100>
- Bai HY, Liu X, 2025. T-Graphormer: using Transformers for spatiotemporal forecasting. <https://doi.org/10.48550/arXiv.2501.13274>
- Do LNN, Taherifar N, Vu HL, 2019. Survey of neural network-based models for short-term traffic state prediction. *WIREs Data Min Knowl Disc*, 9(1):e1285. <https://doi.org/10.1002/widm.1285>
- Essien A, Petrounias I, Sampaio P, et al., 2018. The impact of rainfall and temperature on peak and off-peak urban traffic. *Int Conf on Database and Expert Systems Applications*, p.399-407. [https://doi.org/10.1007/978-3-319-98812-2\\_36](https://doi.org/10.1007/978-3-319-98812-2_36)
- Essien A, Petrounias I, Sampaio P, et al., 2021. A deep-learning model for urban traffic flow prediction with traffic events mined from twitter. *World Wide Web*, 24(4):1345-1368. <https://doi.org/10.1007/s11280-020-00800-3>
- Fang Z, Long QQ, Song GJ, et al., 2021. Spatial-temporal graph ODE networks for traffic flow forecasting. *Proc 27<sup>th</sup> ACM SIGKDD Conf on Knowledge Discovery & Data Mining*, p.364-373. <https://doi.org/10.1145/3447548.3467430>
- FHWA, 2024. How Do Weather Events Affect Roads? [https://ops.fhwa.dot.gov/weather/q1\\_roadimpact.htm](https://ops.fhwa.dot.gov/weather/q1_roadimpact.htm) [Accessed on May 20, 2025].
- Guo SN, Lin YF, Feng N, et al., 2019. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. *Proc 39<sup>th</sup> AAAI Conf on Artificial Intelligence*, p.922-929. <https://doi.org/10.1609/aaai.v33i01.3301922>
- Hochreiter S, Schmidhuber J, 1997. Long short-term memory. *Neur Comput*, 9(8):1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Lee K, Eo M, Jung E, et al., 2021. Short-term traffic prediction with deep neural networks: a survey. *IEEE Access*, 9:54739-54756. <https://doi.org/10.1109/ACCESS.2021.3071174>
- Li Y, Yu R, Shahabi C, et al., 2018. Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting. <https://openreview.net/forum?id=SJIHXGWAZ> [Accessed on May 20, 2025].
- Lim B, Arik SO, Loeff N, et al., 2021. Temporal fusion Transformers for interpretable multi-horizon time series forecasting. *Int J Forecast*, 37(4):1748-1764. <https://doi.org/10.1016/j.ijforecast.2021.03.012>
- Liu J, Guan W, 2004. A summary of traffic flow forecasting methods. *J Highway Transport Res Dev*, 21(3):82-85 (in Chinese). <https://doi.org/10.3969/j.issn.1002-0268.2004.03.022>
- Liu X, Xia YT, Liang YX, et al., 2023. LargeST: a benchmark dataset for large-scale traffic forecasting. *Proc 37<sup>th</sup> Int Conf on Neural Information Processing Systems*, p.75354-75371.
- Ma XL, Yu HY, Wang YP, et al., 2015. Large-scale transportation network congestion evolution prediction using deep learning theory. *PLoS ONE*, 10(3):e0119044. <https://doi.org/10.1371/journal.pone.0119044>
- Meta Research, 2019. Hydra—a Framework for Elegantly Configuring Complex Applications. <https://github.com/facebookresearch/hydra> [Accessed on May 20, 2025].
- Nagy AM, Simon V, 2018. Survey on traffic prediction in smart cities. *Pervasive Mob Comput*, 50:148-163. <https://doi.org/10.1016/j.pmcj.2018.07.004>
- Reshef DN, Reshef YA, Finucane HK, et al., 2011. Detecting novel associations in large data sets. *Science*, 334(6062):1518-1524. <https://doi.org/10.1126/science.1205438>
- Shao ZZ, Wang F, Xu YJ, et al., 2025. Exploring progress in multivariate time series forecasting: comprehensive benchmarking and heterogeneity analysis. *IEEE Trans Knowl Data Eng*, 37(1):291-305. <https://doi.org/10.1109/TKDE.2024.3484454>
- Shaygan M, Meese C, Li WX, et al., 2022. Traffic prediction using artificial intelligence: review of recent advances and emerging opportunities. *Transport Res Part C Emerg Technol*, 145:103921. <https://doi.org/10.1016/j.trc.2022.103921>
- Song C, Lin YF, Guo SN, et al., 2020. Spatial-temporal synchronous graph convolutional networks: a new framework for spatial-temporal network data forecasting. *Proc 34<sup>th</sup> AAAI Conf on Artificial Intelligence*, p.914-921. <https://doi.org/10.1609/aaai.v34i01.5438>
- Sutskever I, Vinyals O, Le QV, 2014. Sequence to sequence learning with neural networks. *Proc 27<sup>th</sup> Int Conf on Neural Information Processing Systems*, p.3104-3112.
- Taghipour H, Parsa AB, Mohammadian AK, 2020. A dynamic approach to predict travel time in real time using data driven techniques and comprehensive data sources.

- Transport Eng*, 2:100025.  
<https://doi.org/10.1016/j.treng.2020.100025>
- Vaswani A, Shazeer N, Parmar N, et al., 2017. Attention is all you need. Proc 31<sup>st</sup> Int Conf on Neural Information Processing Systems, p.6000-6010.
- Wang JY, Jiang JW, Jiang WJ, et al., 2021. LibCity: an open library for traffic prediction. Proc 29<sup>th</sup> Int Conf on Advances in Geographic Information Systems, p.145-148. <https://doi.org/10.1145/3474717.3483923>
- Wu HX, Hu TG, Liu Y, et al., 2023. TimesNet: temporal 2D-variation modeling for general time series analysis. <https://doi.org/10.48550/arXiv.2210.02186>
- Wu ZH, Pan SR, Long GD, et al., 2019. Graph WaveNet for deep spatial-temporal graph modeling. Proc 28<sup>th</sup> Int Joint Conf on Artificial Intelligence, p.1907-1913. <https://doi.org/10.24963/IJCAI.2019/264>
- Yang XX, Zou YJ, Tang JJ, et al., 2020. Evaluation of short-term freeway speed prediction based on periodic analysis using statistical models and machine learning models. *J Adv Transport*, 2020(1):9628957. <https://doi.org/10.1155/2020/9628957>
- Yu B, Yin HT, Zhu ZX, 2018. Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting. Proc 27<sup>th</sup> Int Joint Conf on Artificial Intelligence, p.3634-3640. <https://doi.org/10.24963/ijcai.2018/505>
- Zhao L, Song YJ, Zhang C, et al., 2020. T-GCN: a temporal graph convolutional network for traffic prediction. *IEEE Trans Intell Transport Syst*, 21(9):3848-3858. <https://doi.org/10.1109/TITS.2019.2935152>
- Zheng X, Bagloee SA, Sarvi M, 2025. A comparative study of deep learning and ensemble learning to extend the horizon of traffic forecasting. <https://doi.org/10.48550/arXiv.2504.21358>
- Zhou HY, Zhang SH, Peng JQ, et al., 2021. Informer: beyond efficient Transformer for long sequence time-series forecasting. Proc 39<sup>th</sup> AAAI Conf on Artificial Intelligence, p.11106-11115. <https://doi.org/10.1609/aaai.v35i12.17325>
- Zhu L, Yu FR, Wang YG, et al., 2019. Big data analytics in intelligent transportation systems: a survey. *IEEE*

*Trans Intell Trans Syst*, 20(1):383-398.  
<https://doi.org/10.1109/TITS.2018.2815678>

## List of supplementary materials

- 1 Related works
  - 2 Variables in the experiments
  - 3 Hyperparameter tuning
  - 4 Complete experimental results
- Table S1 Multi-source data fusion in traffic forecasting
- Table S2 Summary of datasets
- Table S3 Summary of variables
- Table S4 Variables and sources for the unified data container components
- Table S5 Normalization for traffic variables in PEMS-BAY
- Table S6 Normalization for traffic variables in METR-LA
- Table S7 Normalization for weather covariates in PEMS-BAY
- Table S8 Normalization for weather covariates in METR-LA
- Table S9 Hyperparameter tuning for models considered in this work
- Table S10 Experimental conditions for model evaluation
- Table S11 Forecast performance without covariates or future information
- Table S12 Relative forecast performance gain with covariates
- Table S13 Relative forecast performance gain with future information
- Fig. S1 Spatial distribution of data sampling points for traffic and weather datasets
- Fig. S2 Correlation between traffic state and covariates
- Fig. S3 Variable selection weights
- Fig. S4 Average traffic speed distribution by DoW