

A GENERAL SPATIO-TEMPORAL BACKBONE WITH SCALABLE CONTEXTUAL PATTERN BANK FOR URBAN CONTINUAL FORECASTING

Aoyu Liu, Yaying Zhang*

The Key Laboratory of Embedded System and Service Computing, Ministry of Education,
Tongji University, Shanghai, China
{liuaoyu, yaying.zhang}@tongji.edu.cn

ABSTRACT

With the rapid growth of spatio-temporal data fueled by IoT deployments and urban infrastructure expansion, accurate and efficient continual forecasting has become a critical challenge. Most existing Spatio-Temporal Graph Neural Networks rely on static graph structures and offline training, rendering them inadequate for real-world streaming scenarios characterized by graph expansion and distribution shifts. Although Continual Spatio-Temporal Forecasting methods have been proposed to tackle these issues, they often adopt backbones with limited modeling capacity and lack effective mechanisms to balance stability and adaptability. To overcome these limitations, we propose *STBP*, a novel framework that integrates a general spatio-temporal backbone with a scalable contextual pattern bank. The backbone extracts stable representations in the frequency domain and captures dynamic spatial correlations through lightweight linear graph attention. To support continual adaptation and mitigate catastrophic forgetting, the contextual pattern bank is updated incrementally via parameter expansion, enabling the capture of evolving node-level heterogeneity and relevance. During incremental training, the backbone remains fixed to preserve general knowledge, while the pattern bank adapts to new scenarios and distributions. Extensive experiments demonstrate that *STBP* outperforms state-of-the-art baselines in both forecasting accuracy and scalability, validating its effectiveness for continual spatio-temporal forecasting. Code is available at <https://github.com/Aoyu-Liu/STBP>.

1 INTRODUCTION

With the rapid development of urban IoT sensing systems, spatio-temporal data such as traffic flow (Shao et al., 2022b) and air quality (Tian et al., 2025) observations continue to surge (Kumar et al., 2024; Hu et al., 2023; Fang et al., 2026). Conducting efficient and accurate forecasting on data streams has become a core task in the development of smart cities (Jin et al., 2024; Yang et al., 2025). Unlike traditional offline learning based on static assumptions, real-world urban environments are in a state of continuous evolution—dynamic changes in urban structure and behavioral patterns constantly drive the evolution of graph structures and data distributions.

Spatio-Temporal Graph Neural Networks (STGNNs) (Kong et al., 2024; Gao et al., 2024; Liu & Zhang, 2025) have been widely used to model complex spatio-temporal dependencies. However, most existing models still adhere to the paradigm of “fixed topology + offline training”: the graph structure is predefined and fixed during the training phase, and the model is deployed directly after training. Yet, as shown in Figure 1, this static assumption becomes difficult to sustain when the node set continuously expands or the connectivity dynamically reconstructs over time. If one relies solely on structural modifications and continuous fine-tuning to handle node increments, model performance often degrades significantly. Therefore, Continual Spatio-Temporal Forecasting (CSTF) (Miao et al., 2024b; Chen & Liang, 2025; Ma et al., 2025b) has garnered increasing attention. Its goal is to achieve

*Corresponding author.

incremental learning and efficient inference on new data without repeatedly relying on retraining with historical data. As shown in Figure 1, typical CSTF approaches employ a general spatio-temporal backbone integrated with strategies such as regularization, replay, or dynamic architectures to adapt to graph structural expansion and mitigate catastrophic forgetting.

However, two key issues in existing CSTF methods have not yet been adequately addressed. First, the general backbone adopted by most current methods is relatively simple (e.g., stacks of graph and temporal convolutions), making it difficult to effectively handle incremental scenarios characterized by dynamically changing spatio-temporal correlations and long-term distribution drift. Forcibly adapting existing STGNNs for continual learning often leads to performance degradation (Shao et al., 2024; Ma et al., 2025a). Second, continual optimization strategies based on dynamic structural expansion are often weakly coupled with the backbone—such as direct parameter expansion or prompt concatenation—making it challenging to achieve a good balance among model stability, adaptability, and interpretability. Based on the above issues, we argue that an ideal CSTF framework should simultaneously address the following four key challenges: ❶ *handling distributional drift*; ❷ *modeling dynamic spatio-temporal correlations*; ❸ *alleviating catastrophic forgetting*; and ❹ *designing an incremental strategy that efficiently collaborates with the backbone*.

To this end, we bridge the gap between STGNNs and continual learning by introducing a general-purpose spatio-temporal backbone with scalable contextual pattern bank (STBP). Specifically, the backbone in STBP leverages frequency-domain modules to extract stable spatio-temporal components, mitigating distributional drift. Simultaneously, a lightweight, scene-agnostic linear graph attention mechanism is introduced to model dynamic spatial correlations with low computational overhead. To mitigate catastrophic forgetting and support continuous graph structure expansion, we design a contextual pattern bank composed of trainable parameters. It incrementally updates knowledge via parameter expansion and interacts with the backbone through gating and attention mechanisms, thereby uncovering node relevance and heterogeneity, and gradually adapting to scenario expansion at low cost. Within this framework, the backbone is responsible for modeling general and stable patterns, while the contextual pattern bank captures node-related heterogeneous contexts, working collaboratively to adapt to continuously evolving environments.

Our main contributions are summarized as follows: ❶ We propose an efficient and general backbone tailored for continual forecasting tasks, capable of modeling dynamic spatial correlations and mitigating distribution shift; ❷ We design a prompt-based guidance mechanism using contextual pattern bank, supporting dynamic model adaptation and alleviating catastrophic forgetting; ❸ Extensive experiments on multiple real-world datasets demonstrate that STBP significantly outperforms state-of-the-art baselines in terms of forecasting accuracy, adaptability, and scalability.

2 RELATED WORK

Spatio-Temporal Forecasting. Early studies in spatio-temporal forecasting, including methods like STGCN (Yu et al., 2018) and DCRNN (Li et al., 2018), primarily focused on combining basic temporal and spatial elements for prediction tasks. These models typically depended on predefined geographic adjacency matrices, which limited their ability to capture the evolving nature of spatial correlations. In contrast, later advancements, such as GWNet (Wu et al., 2019), DGCRN (Li et al., 2023), and MegaCRN (Jiang et al., 2023b), addressed this limitation by incorporating adaptive adjacency matrices or learning spatial correlations directly from the data. This shift led to a notable improvement in forecasting accuracy. More recently, models like STID (Shao et al., 2022a), STAEformer (Liu et al., 2023a), and HimNet (Dong et al., 2024) have emphasized the significance of distinguishing spatial patterns to further enhance forecasting performance. These methods incorporate trainable components, including spatial embeddings, parameter pools, and contextual pattern bank, to more accurately capture spatial variations, boosting both prediction precision and model adaptability.

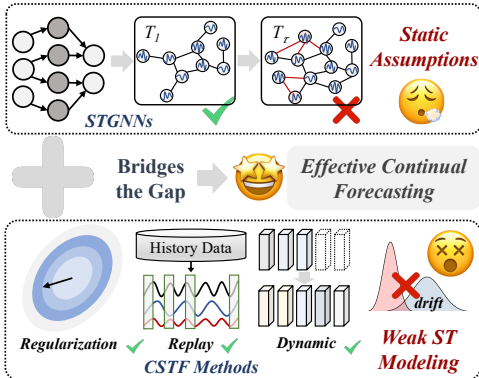


Figure 1: Limitations of existing studies.

Continual Spatio-Temporal Forecasting. TrafficStream (Chen et al., 2021), one of the pioneering frameworks in CSTF, integrates spatio-temporal modeling with continual learning by employing historical data replay and parameter smoothing to manage long-term streaming traffic data and achieve accurate traffic flow prediction. Building on this line of work, STKEC (Wang et al., 2023a) proposes an influence-based knowledge expansion strategy together with a memory-augmented knowledge consolidation mechanism, which better supports the scaling of transportation networks while alleviating catastrophic forgetting. PECPM (Wang et al., 2023b) leverages pattern matching to dynamically maintain a traffic pattern bank, enabling efficient, historical-data-free continual learning with improved accuracy. STRAP (Zhang et al., 2025) adopts retrieval-augmented learning, constructing multi-dimensional pattern libraries and using plug-and-play prompting to fuse retrieved patterns, thereby enhancing out-of-distribution (OOD) generalization and mitigating catastrophic forgetting. EAC (Chen & Liang, 2025) introduces prompt tuning via a dynamic prompt pool that expands and compresses over time, balancing adaptation to new nodes with knowledge preservation in a parameter-efficient manner. Additionally, UFCL (Miao et al., 2025) leverages federated learning to protect data privacy and employs a global replay buffer of synthetic spatio-temporal data, addressing the challenges of distributed streaming environments.

3 PRELIMINARY

Definition 1 (Streaming Spatio-Temporal Graph). We define a streaming spatio-temporal graph as a sequence of evolving graphs $\mathbb{G} = \{G_\tau\}_{\tau=1}^T$, where each graph $G_\tau = (V_\tau, E_\tau, A_\tau)$ represents the graph at incremental period τ . Here, V_τ denotes the node set, E_τ the edge set, and adjacency matrix $A_\tau \in \mathbb{R}^{N_\tau \times N_\tau}$ connections between nodes. The number of nodes at period τ is denoted by $N_\tau = |V_\tau|$. The graph evolves incrementally as $G_\tau = G_{\tau-1} + \Delta G_\tau$, where ΔG_τ captures structural or feature modifications between periods.

Definition 2 (Continual Spatio-Temporal Forecasting). Continual spatio-temporal forecasting aims to develop an optimal predictive model at each stage based on dynamic, streaming spatio-temporal graph data. At each incremental period τ , given the current graph G_τ and historical observations $\mathbf{X}_\tau \in \mathbb{R}^{N_\tau \times T_h}$, the goal is to predict future signals $\mathbf{Y}_\tau \in \mathbb{R}^{N_\tau \times T_f}$ as follows:

$$\hat{\mathbf{Y}}_\tau = f_\theta(G_\tau, \mathbf{X}_\tau), \quad (1)$$

where T_h is the length of the historical observation window, and T_f is the forecasting horizon. The model f_θ is parameterized by θ , and continually updated by minimizing:

$$\theta_\tau^* = \arg \min_{\theta} \mathbb{E}_{(G_\tau, \mathbf{X}_\tau, \mathbf{Y}_\tau) \sim \mathcal{D}_\tau} [\mathcal{L}(f_\theta(G_\tau, \mathbf{X}_\tau), \mathbf{Y}_\tau)], \quad (2)$$

where $\mathcal{L}(\cdot, \cdot)$ is a loss function, and \mathcal{D}_τ denotes the data distribution at period τ .

4 METHODOLOGY

4.1 OVERVIEW OF STBP

The workflow and architecture of STBP are shown in Figure 2. It consists of two core components: a general spatio-temporal backbone and a contextual pattern bank. The backbone, comprising temporal and spatial modules with a prediction layer, captures spatio-temporal correlations in evolving networks. The contextual pattern bank, made of trainable parameters, is dynamically expanded and fine-tuned as data evolves. While the backbone captures general, stable patterns, the contextual pattern bank adapts to environmental changes, focusing on context-specific patterns. Guided by prompts, both components collaborate to form an efficient and robust continual learning system.

In terms of workflow, streaming spatio-temporal data is sequentially fed into the STBP. During the initial incremental training phase, the backbone and contextual pattern bank are jointly trained to capture spatio-temporal correlations from current data. In later stages, the backbone is frozen (denoted by a snowflake) to retain knowledge learned from historical data, while the contextual pattern bank is updated (denoted by a flame) through expansion and fine-tuning. These updates serve as prompts, guiding the frozen backbone to adapt to new data distributions. This continual learning process, driven by the interplay between backbone and contextual pattern bank, enables the model to progressively enhance its representation power and adaptability while preserving core functionality. For detailed workflow steps, refer to Algorithm 1 in Appendix A.3.2.

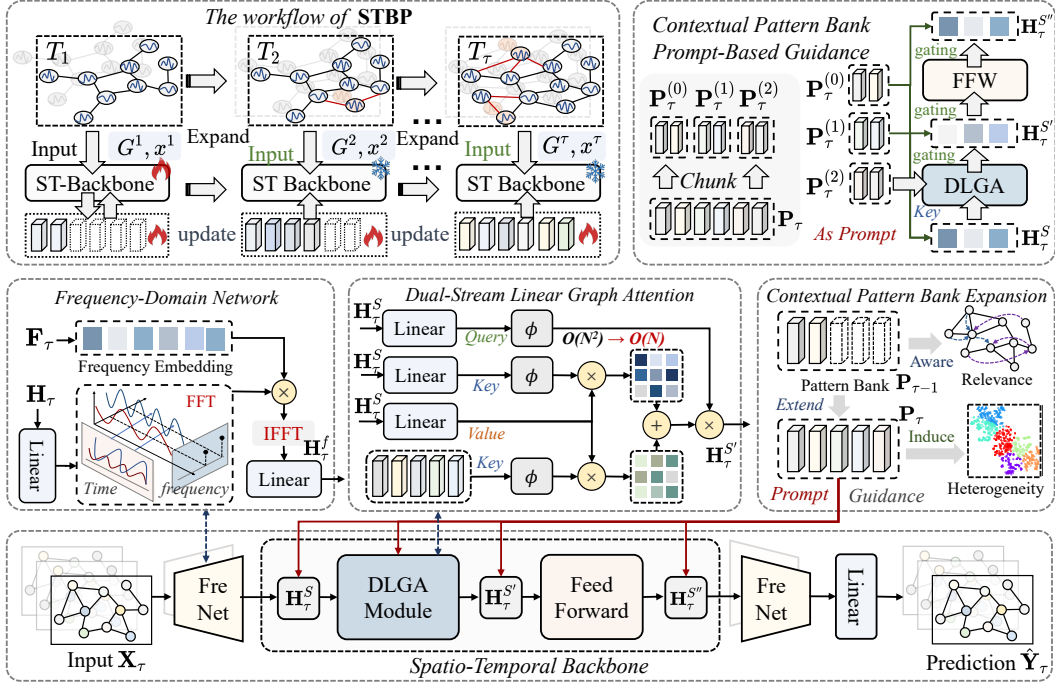


Figure 2: The overall workflow and architecture of STBP.

4.2 CONTEXTUAL PATTERN BANK

Recent studies (Shao et al., 2022a; Dong et al., 2024; Chen & Liang, 2025) have shown that incorporating node-specific trainable parameters into STGNNs can significantly enhance forecasting performance. Following this insight, we propose an expandable contextual pattern bank $\mathbf{P}_\tau \in \mathbb{R}^{N_\tau \times d}$, composed of trainable parameters, to consolidate historical spatio-temporal patterns and generalize to new ones, thereby mitigating *catastrophic forgetting* and continuously adapting to new incremental scenarios, where d denotes the feature dimension.

We posit that the model can utilize \mathbf{P}_τ to effectively distinguish both the *relevance* and *heterogeneity* of nodes, enabling a more nuanced understanding of the underlying data structures. Here, *relevance* refers to shared behavioral patterns among nodes—such as similar trends or periodic fluctuations—while *heterogeneity* captures differences arising from distinct node functions or external factors such as geography, policy, or events. To validate this hypothesis, we conduct a t-SNE-based analysis on \mathbf{P}_τ trained on spatio-temporal datasets (see Figure 3), which reveals meaningful clustering patterns. Each cluster exhibits distinct characteristics, corresponding to *heterogeneity*, while nodes within the same cluster display similar temporal dynamics, reflecting *relevance*.

As shown in Figure 2, given a streaming spatio-temporal input $\mathbf{X}_\tau \in \mathbb{R}^{N_\tau \times T_h}$, the backbone model \mathcal{M}_θ , and contextual pattern bank $\mathbf{P}_\tau \in \mathbb{R}^{N_\tau \times d}$, the incremental learning process is formulated as:

$$\hat{\mathbf{Y}}_\tau = \mathcal{M}_\theta(\mathbf{X}_\tau, \mathbf{P}_\tau). \quad (3)$$

At the initial training stage ($\tau = 1$), both the backbone and contextual pattern bank are jointly trained (denoted with flame). For subsequent stages ($\tau > 1$), the backbone is frozen (denoted with snowflake), and only the contextual pattern bank is updated through expansion:

$$\mathbf{P}'_\tau = \mathbf{P}_{\tau-1} \parallel \Delta \mathbf{P}_\tau, \quad (4)$$

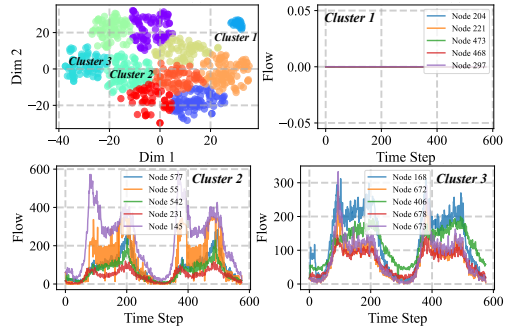


Figure 3: Contextual pattern bank visualization.

where $\Delta \mathbf{P}_\tau \in \mathbb{R}^{(N_\tau - N_{\tau-1}) \times d}$ represents newly introduced parameters for the current incremental period. Only the expanded contextual pattern bank $\mathbf{P}'_\tau \in \mathbb{R}^{N_\tau \times d}$ is fine-tuned during training. Notably, even without explicit clustering constraints, the contextual pattern bank autonomously distinguishes heterogeneous and relevant nodes through data-driven parameter learning and prompt-based interactions with the backbone, driven by the prediction task. This strategy ensures that the backbone retains previously acquired knowledge, while the contextual pattern bank continually adapts to evolving distributions. It incrementally expands to represent an increasingly diverse set of environmental patterns, thereby avoiding the inadequacy exhibited by fixed models in novel scenarios.

Distinct from existing work (Wang et al., 2023a; Chen & Liang, 2025; Wang et al., 2023b), we introduce a *Prompt-Based Guidance* (Peebles & Xie, 2023; Zhang et al., 2023) mechanism to enhance \mathbf{P}_τ 's capacity to model both node-level relevance and heterogeneity. Specifically, the contextual pattern bank comprises three groups of trainable parameters: $\mathbf{P}_\tau^{(i)} \in \mathbb{R}^{N_\tau \times d}$ for $i \in 0, 1, 2$. As illustrated in Figure 2, these components interact with the backbone's hidden representation \mathbf{H}_τ via the following prompt-based gating function:

$$\mathbf{H}'_\tau = \mathbf{P}_\tau^{(1)} \cdot h_\theta(\mathbf{H}_\tau \cdot (1 + \mathbf{P}_\tau^{(0)})), \quad (5)$$

where h_θ denotes an arbitrary submodule within the backbone. This gating mechanism enables adaptive modeling of node heterogeneity. Additionally, $\mathbf{P}_\tau^{(2)}$ acts as a key embedding in the attention module, guiding the backbone to generalize correlation-aware information under task constraints. Importantly, since the contextual pattern bank encodes high-level abstractions rather than raw historical data, our method supports knowledge retention without revisiting prior data—offering advantages in *privacy protection* and *storage efficiency*.

4.3 GENERAL SPATIO-TEMPORAL BACKBONE

While the contextual pattern bank mitigates catastrophic forgetting in continual learning, it lacks the ability to model dynamic spatio-temporal correlations and handle distributional drift. To address this, we design a **general spatio-temporal backbone** aimed at handling distributional drift, spatio-temporal correlation modeling, and graph scalability during continual learning. The term *general* implies that the backbone is independent of the number of nodes and does not rely on any predefined adjacency matrix, making it adaptable to arbitrary spatio-temporal data structures.

As shown in Figure 2, the backbone operates as follows: the input spatio-temporal data first passes through a **frequency-domain network** (FreNet), which maps it into high-dimensional temporal representations and extracts stable components via frequency domain analysis. A **dual-stream linear graph attention** (DLGA) module then captures dynamic spatial correlations, followed by a feedforward layer with a multilayer perceptron for enhanced nonlinear expressivity. Finally, the features are reconstructed to their original shape by another FreNet and passed through a prediction layer. We detail the FreNet and DLGA modules below.

Frequency-Domain Network. Spatio-temporal data in evolving environments often suffer from distributional drift (Wang et al., 2024; Ji et al., 2025; Zhou et al., 2023). Although the contextual pattern bank helps retain stable knowledge, we further address this issue through a dedicated frequency-domain analysis (Xia et al., 2023). FreNet is designed to capture temporal correlations while emphasizing stable components in the data, such as periodicity and trends, which are more resilient to distributional changes (Liu & Zhang, 2025). Specifically, STBP employs two FreNets—one at the beginning and one at the end of the backbone (Figure 2). The first maps input data $\mathbf{X}_\tau \in \mathbb{R}^{N_\tau \times T_h}$ through a linear layer into a high-dimensional representation $\mathbf{H}_\tau \in \mathbb{R}^{N_\tau \times d}$, which is then transformed to the frequency domain using a Fast Fourier Transform (FFT). A learnable frequency-domain embedding $\mathbf{F}_\tau \in \mathbb{C}^{\frac{d}{2}+1}$ adaptively highlights stable features. This process is formalized as:

$$\mathbf{H}_\tau^f = \text{IFFT}(\text{FFT}(\mathbf{H}_\tau) \odot \mathbf{F}_\tau), \quad (6)$$

where $\mathbf{H}_\tau^f \in \mathbb{R}^{N_\tau \times d}$ is further processed by a linear layer. The resulting representation \mathbf{H}_τ^f then interacts with the contextual pattern bank component $\mathbf{P}_\tau^{(0)}$ via gating-based prompt guidance (Eq. 5) to produce $\mathbf{H}_\tau^s \in \mathbb{R}^{N_\tau \times d}$, which serves as input to the subsequent DLGA module. The second FreNet performs an inverse operation, restoring the feature shape to $\mathbb{R}^{N_\tau \times T_h}$. Compared to traditional temporal modules like RNNs (Li et al., 2018; Bai et al., 2020) or TCNs (Zheng et al., 2023; Fang et al., 2023), FreNet offers higher computational efficiency and enhanced ability to extract

stable low-frequency components (e.g., periodicity and trends) while suppressing high-frequency noise, thereby obtaining more robust temporal representations that are resilient to distributional drift across periods and scenarios.

Dual-Stream Linear Graph Attention. After obtaining stable components, it remains essential to capture complex spatial interactions and time-varying node correlations. An effective spatial module must adaptively learn node correlations in a data-driven manner, maintain computational efficiency, and scale to growing graphs. Graph attention mechanisms (Veličković et al., 2018) have emerged as promising solutions, enabling dynamic correlation modeling without relying on fixed adjacency matrices. However, conventional graph attention (Zheng et al., 2020; Jiang et al., 2023a; Liu et al., 2023a) incurs $O(N^2)$ complexity, limiting its scalability. To overcome this, we propose DLGA (Figure 2), which improves efficiency using a *random feature mapping*-based linear attention mechanism (Katharopoulos et al., 2020). Moreover, DLGA introduces a **dual-stream structure** by incorporating the contextual pattern bank $\mathbf{P}_\tau^{(2)} \in \mathbb{R}^{N_\tau \times d}$ as an additional *key*. This enables the model to assess the relationship between evolving input patterns and stored knowledge. Formally:

$$\mathbf{Q} = \mathbf{W}_q \mathbf{H}_\tau^s, \quad \mathbf{K} = \mathbf{W}_k \mathbf{H}_\tau^s, \quad \mathbf{V} = \mathbf{W}_v \mathbf{H}_\tau^s, \quad (7)$$

$$\begin{aligned} \mathbf{H}_\tau^{s'} &= \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{P}_\tau^{(2)}) \\ &= \text{Softmax}(\mathbf{Q}\mathbf{K}^\top + \mathbf{Q}(\mathbf{P}_\tau^{(2)})^\top) \mathbf{V}, \end{aligned} \quad (8)$$

$$\begin{aligned} \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{P}_\tau^{(2)}) &\approx (\phi(\mathbf{Q})\phi(\mathbf{K})^\top + \phi(\mathbf{Q})\phi(\mathbf{P}_\tau^{(2)})^\top) \mathbf{V} \\ &= \phi(\mathbf{Q}) \left(\phi(\mathbf{K})^\top \mathbf{V} + \phi(\mathbf{P}_\tau^{(2)})^\top \mathbf{V} \right). \end{aligned} \quad (9)$$

Here, \mathbf{W}_q , \mathbf{W}_k , and \mathbf{W}_v are trainable projection matrices. \mathbf{H}_τ^s and $\mathbf{H}_\tau^{s'} \in \mathbb{R}^{N_\tau \times d}$ denote the input and the spatially enriched representation passed to the feedforward layer of the DLGA module, respectively. The function $\phi(\cdot)$ denotes a random feature mapping, with Softmax used for approximation in our implementation. For further details on the approximation derivation, see Appendix A.3.1. Notably, the linear attention approximation does not explicitly construct an adjacency matrix. Instead, it implicitly models dynamic correlations by reordering operations in the attention computation. DLGA reduces computational complexity from quadratic to linear, while preserving dynamic spatial modeling and seamlessly integrating prompt-based knowledge from the contextual pattern bank.

5 EXPERIMENT

5.1 EXPERIMENTAL SETTINGS

Datasets. We evaluate our model on three real-world streaming spatio-temporal datasets from the traffic and meteorology domains. The traffic datasets, **PEMS-Stream** (Chen et al., 2001) and **CA-Stream** (Liu et al., 2023b), consist of traffic flow measurements provided by the California Department of Transportation (CalTrans), with a sampling interval of 5 minutes. The meteorological dataset, **AIR-Stream** (Chen & Liang, 2025), is derived from urban air quality platform of the Chinese Environmental Monitoring Center, with hourly sampling intervals. To ensure fair evaluation, all datasets are split into training, validation, and test sets using a fixed ratio of 6:2:2. For each prediction task, the model is trained to forecast the next 12 time steps based on the previous 12 observations. Detailed dataset statistics are provided in Appendix A.4.1.

Baselines and Metrics. We select representative models from two categories as baselines: \triangleright Conventional spatio-temporal forecasting models, including lightweight spatio-temporal architectures such as **GWNet** (Wu et al., 2019), **STID** (Shao et al., 2022a), and **iTransformer** (Liu et al., 2024b). These models are adapted specifically for incremental training in our experiments. \triangleright Continual spatio-temporal forecasting models, including **TrafficStream**, **STKEC** (Wang et al., 2023a), **PECPM** (Wang et al., 2023b), **STRAP** (Zhang et al., 2025), and **EAC** (Chen & Liang, 2025). The performance of all models is evaluated using the following metrics: Mean Absolute Error (**MAE**), Root Mean Squared Error (**RMSE**), and Mean Absolute Percentage Error (**MAPE**). More details on this are included in Appendix A.4.2.

5.2 MAIN RESULTS

The main experimental results are summarized in Table 1, which reports the metrics averaged over all incremental periods. We also present the results at specific forecasting horizons (3, 6, and 12 time steps ahead), together with the overall average across horizons. STGNNs, including GWNet and STID,

Table 1: Main experimental results. **Bold**: best, underline: second best.

Dataset	Metric	Horizon	GWNet	STID	iTransformer	TrafficStream	STKEC	PECPM	STRAP	EAC	STBP	
PEMS-Stream	MAE	3	19.64±0.12	24.34±0.13	17.63±0.76	14.23±0.09	14.29±0.12	14.26±0.13	14.30±0.11	<u>13.86</u> ±0.16	11.62 ±0.09	
		6	19.68±0.19	25.45±0.21	20.82±0.76	16.43±0.03	16.44±0.11	16.35±0.12	16.34±0.10	<u>15.40</u> ±0.19	12.26 ±0.10	
		12	20.63±0.09	29.42±0.38	28.33±0.86	21.76±0.07	21.66±0.11	21.46±0.19	21.52±0.15	<u>18.90</u> ±0.28	13.47 ±0.08	
		Avg.	<u>19.87</u> ±0.10	<u>26.07</u> ±0.23	21.60±0.79	<u>16.95</u> ±0.03	16.96±0.09	16.86±0.12	16.88±0.10	<u>15.67</u> ±0.20	12.31 ±0.07	
	RMSE	3	32.20±0.17	39.37±0.13	28.20±1.15	23.00±0.09	23.08±0.14	23.07±0.15	23.06±0.13	<u>22.26</u> ±0.23	19.20 ±0.13	
		6	32.34±0.32	40.86±0.19	33.80±1.13	26.87±0.04	26.93±0.15	26.76±0.20	26.71±0.14	<u>24.99</u> ±0.28	20.51 ±0.15	
		12	33.73±0.09	46.20±0.43	45.98±1.25	35.29±0.11	35.19±0.11	34.77±0.37	34.80±0.19	<u>30.56</u> ±0.45	22.67 ±0.13	
		Avg.	<u>32.59</u> ±0.18	<u>41.67</u> ±0.21	34.88±1.17	<u>27.52</u> ±0.05	27.56±0.11	27.37±0.20	27.35±0.13	<u>25.30</u> ±0.29	20.52 ±0.11	
	MAPE (%)	3	27.47±0.69	37.79±2.23	32.46±3.04	18.34±0.67	18.54±0.61	<u>18.19</u> ±0.66	18.69±0.52	18.35±0.31	<u>15.00</u> ±0.24	15.00 ±0.24
		6	27.22±0.58	39.70±2.43	36.73±3.84	20.77±0.71	20.64±0.48	20.79±0.57	21.33±0.41	<u>20.11</u> ±0.36	<u>15.55</u> ±0.26	15.55 ±0.26
		12	29.38±1.18	47.94±2.91	54.31±4.66	27.88±0.26	27.05±0.62	28.33±0.52	28.20±1.10	<u>24.30</u> ±0.57	<u>16.75</u> ±0.23	16.75 ±0.23
		Avg.	<u>27.79</u> ±0.76	<u>41.09</u> ±2.49	39.63±3.81	<u>21.66</u> ±0.54	21.50±0.52	21.73±0.45	22.17±0.46	<u>20.42</u> ±0.41	<u>15.65</u> ±0.21	15.65 ±0.21
CA-Stream	MAE	3	23.49±0.80	27.71±0.23	20.16±0.06	17.82±0.26	17.69±0.19	17.93±0.12	23.59±0.61	<u>17.66</u> ±0.37	15.01 ±0.18	
		6	23.31±0.69	28.93±0.26	24.37±0.06	20.38±0.17	20.41±0.04	20.33±0.09	25.38±0.68	<u>19.68</u> ±0.54	15.78 ±0.07	
		12	24.78±0.86	33.61±0.45	34.05±0.06	26.92±0.53	27.05±0.17	26.68±0.19	31.10±0.89	<u>24.86</u> ±1.33	17.19 ±0.09	
		Avg.	<u>23.73</u> ±0.75	<u>29.71</u> ±0.28	25.34±0.05	<u>21.09</u> ±0.29	21.09±0.13	21.04±0.11	26.25±0.62	<u>20.20</u> ±0.69	15.77 ±0.09	
	RMSE	3	35.87±0.98	41.53±0.31	31.58±0.09	28.01±0.22	28.02±0.19	28.00±0.16	34.73±0.74	<u>27.46</u> ±0.46	24.37 ±0.27	
		6	35.68±0.88	43.14±0.35	37.76±0.10	32.19±0.22	32.43±0.05	31.94±0.09	37.97±0.86	<u>30.64</u> ±0.83	25.71 ±0.22	
		12	37.57±1.11	49.18±0.58	51.24±0.10	41.59±0.64	42.08±0.21	41.14±0.30	46.74±1.36	<u>37.71</u> ±1.94	28.08 ±0.14	
		Avg.	<u>36.20</u> ±0.96	<u>44.12</u> ±0.37	38.94±0.09	<u>33.01</u> ±0.35	33.24±0.13	32.77±0.17	39.05±0.80	<u>31.18</u> ±0.99	25.70 ±0.16	
	MAPE (%)	3	24.61±0.95	29.24±0.65	21.76±0.17	17.05±0.41	<u>16.60</u> ±0.19	17.63±0.91	19.11±0.49	18.26±1.88	<u>14.22</u> ±0.03	14.22 ±0.03
		6	24.44±0.80	30.66±0.78	26.76±0.22	19.22±0.30	<u>18.98</u> ±0.17	19.74±0.92	20.48±0.39	19.45±1.16	<u>14.85</u> ±0.07	14.85 ±0.07
		12	25.71±0.81	36.88±1.29	39.81±0.38	25.47±0.46	24.99±0.29	25.94±1.09	24.97±0.59	<u>24.52</u> ±1.10	<u>16.20</u> ±0.08	16.20 ±0.08
		Avg.	<u>24.79</u> ±0.85	<u>31.73</u> ±0.86	28.34±0.20	<u>19.98</u> ±0.30	<u>19.61</u> ±0.19	20.49±0.91	21.15±0.47	<u>20.17</u> ±1.25	<u>14.94</u> ±0.05	14.94 ±0.05
AIR-Stream	MAE	3	28.48±1.43	32.85±0.21	22.37±0.76	20.73±0.40	20.95±0.17	20.82±0.35	21.41±0.33	<u>20.41</u> ±0.36	20.00 ±0.14	
		6	29.79±0.89	33.15±0.22	26.22±0.48	25.64±0.34	25.54±0.08	25.54±0.19	26.12±0.34	<u>25.20</u> ±0.29	24.70 ±0.30	
		12	31.30±0.52	33.88±0.25	29.45±0.31	29.04±0.23	28.94±0.12	28.95±0.11	29.38±0.31	<u>28.57</u> ±0.42	28.28 ±0.63	
		Avg.	<u>29.66</u> ±1.01	<u>33.23</u> ±0.22	25.53±0.56	<u>24.58</u> ±0.34	24.63±0.11	24.60±0.21	25.16±0.32	<u>24.21</u> ±0.43	23.64 ±0.23	
	RMSE	3	44.38±2.04	51.24±0.28	34.98±1.18	32.80±0.57	33.13±0.28	33.07±0.52	33.72±0.41	<u>32.19</u> ±0.57	32.15 ±0.24	
		6	46.22±1.28	51.61±0.31	40.95±0.73	40.41±0.53	40.38±0.20	40.48±0.39	41.13±0.40	39.63 ±0.43	39.81 ±0.26	
		12	48.34±0.85	52.55±0.39	45.70±0.55	45.54±0.47	45.53±0.27	45.63±0.27	46.07±0.34	<u>44.65</u> ±0.63	44.97 ±0.97	
		Avg.	<u>46.01</u> ±1.46	<u>51.72</u> ±0.33	39.67±0.91	<u>38.58</u> ±0.53	38.70±0.26	38.76±0.41	39.37±0.38	<u>37.83</u> ±0.60	37.76 ±0.30	
	MAPE (%)	3	38.02±2.60	43.52±0.64	28.64±1.28	26.33±0.30	26.24±0.30	<u>25.79</u> ±0.50	26.80±0.36	26.06±0.71	<u>24.64</u> ±0.16	24.64 ±0.16
		6	39.98±1.70	44.12±0.57	34.91±0.68	33.33±0.21	33.10±0.28	32.97±0.18	33.30±0.19	<u>32.88</u> ±0.64	<u>30.66</u> ±0.42	30.66 ±0.42
		12	42.37±1.14	45.06±0.62	40.79±0.39	39.27±0.24	39.02±0.18	<u>38.67</u> ±0.02	38.87±0.24	38.85±0.67	<u>36.23</u> ±0.52	36.23 ±0.52
		Avg.	<u>39.87</u> ±1.87	<u>44.16</u> ±0.60	34.15±0.76	<u>32.29</u> ±0.29	32.12±0.21	31.82±0.16	32.37±0.28	<u>31.77</u> ±0.53	29.70 ±0.35	

rely on static graph assumptions and are not designed for continual learning. Following prior work (Chen & Liang, 2025), we therefore **re-train** the backbone from scratch at each incremental stage using only data from the current period. In contrast, iTransformer is scenario-agnostic, so we adopt an **online** training regime: at each stage it is trained on the complete node set of the current spatio-temporal graph, initialized from the previous period’s weights, enabling end-to-end fine-tuning. More detailed experimental results are provided in Appendix A.4.4.

Results of conventional methods. As shown in Table 1, STGNNs trained from scratch achieve only poor performance on all datasets. Although these methods work well under static assumptions, they fail to exploit past spatio-temporal knowledge, resulting in unsatisfactory performance. In contrast, iTransformer performs better by leveraging historical spatio-temporal information through online training, but it still suffers from catastrophic forgetting and is therefore not an ideal solution.

Results of CSTF methods. The best-performing models are those that explicitly mitigate catastrophic forgetting, including CSTF methods such as PECPM, STRAP, and EAC. Compared with full-parameter fine-tuning strategies (e.g., PECPM, STKEC, TrafficStream), lightweight prompt-based adaptation on a frozen backbone (e.g., EAC, STRAP, STBP) yields higher average accuracy, highlighting the benefits of dynamically tuning only a small set of parameters. Nevertheless, STRAP performs notably poorly on CA-Stream, indicating that retrieval-based pattern matching struggles

Table 2: Comparison of few-shot forecasting performance.

Model	PEMS-Stream 10%			CA-Stream 10%		
	MAE	RMSE	MAPE (%)	MAE	RMSE	MAPE (%)
GWNet	30.15±1.06	45.30±1.59	48.80±3.85	33.73±0.89	50.80±1.43	36.52±0.86
STID	33.42±2.90	50.63±3.73	63.96±12.60	37.09±0.52	55.10±0.69	39.18±1.12
iTransformer	20.99±0.19	32.67±0.25	49.11±1.62	25.43±0.08	39.01±0.10	28.39±0.54
TrafficStream	17.23±0.08	27.49±0.17	27.63±0.43	21.28±0.19	33.25±0.22	20.45±0.45
STKEC	17.75±0.12	28.23±0.13	27.80±0.88	21.20±0.13	33.20±0.08	<u>20.23</u> ±0.46
PECPM	17.05±0.02	27.20±0.07	29.08±1.90	21.48±0.15	33.33±0.13	21.25±0.86
STRAP	17.68±0.10	27.98±0.14	31.67±2.88	26.34±0.79	39.39±1.09	21.34±0.45
EAC	<u>16.13</u> ±0.05	<u>25.57</u> ±0.06	<u>24.02</u> ±1.23	<u>20.94</u> ±0.70	<u>32.19</u> ±1.00	21.37±1.53
STBP	13.58 ±0.05	22.24 ±0.13	17.89 ±0.29	17.11 ±0.03	27.48 ±0.16	17.60 ±0.30

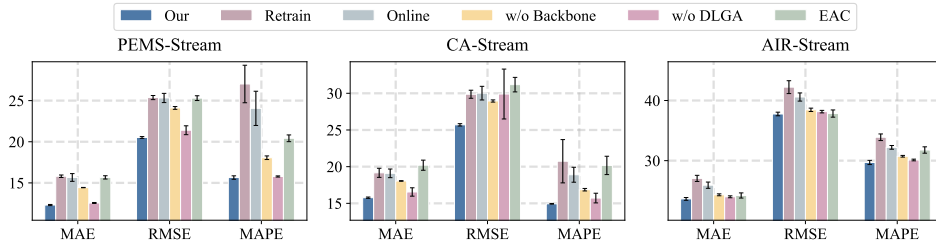


Figure 4: Results of ablation experiments.

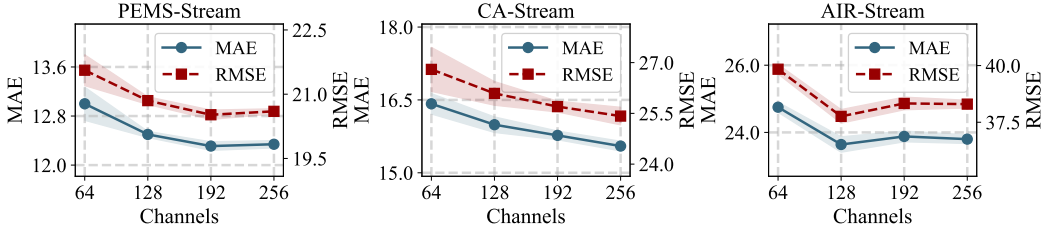


Figure 5: Results of parameter experiments.

in extreme incremental scenarios with rapid, large-scale topology expansion. Overall, our proposed STBP outperforms all competing models. Compared with the best baseline, STBP reduces the average MAE by **21.44%**, **21.93%**, and **2.35%** on the PEMS-Stream, CA-Stream, and AIR-Stream datasets, respectively. This gain stems from the bridge it establishes between STGNNs and CSTF methods: the carefully designed general spatio-temporal backbone and contextual pattern bank jointly capture dynamic spatio-temporal correlations, thereby mitigating catastrophic forgetting and alleviating distributional drift.

Results of few-shot forecasting task. To further evaluate the robustness of the proposed model under low-resource scenarios, we construct a few-shot training setting and compare it against existing baselines. Specifically, we simulate a few-shot setting in which the sample size of the first incremental period is kept unchanged, while the training set size for subsequent periods is reduced to only 10% of the original. The test set size remains fixed throughout. As shown in Table 2, STBP consistently outperforms all other methods, highlighting its strong ability to extract meaningful patterns from limited data. CSTF baselines are more resilient to low-resource conditions than conventional STGNNs (e.g., GWNet, STID). This demonstrates that when data is extremely scarce, conventional models struggle to capture stable spatio-temporal patterns, whereas CSTF methods can leverage knowledge accumulated from historical stages to adapt more quickly to new nodes. The continual learning mechanism effectively mitigates catastrophic forgetting, allowing the model to continuously utilize previously learned general features during incremental learning.

5.3 ABLATION STUDY & PARAMETER SENSITIVITY ANALYSIS

Ablation Study Settings. To validate the core contributions of STBP, we design the following variants for ablation experiments: **❶ Retrain:** The contextual pattern bank is removed. Similar to GWNet and STID, a new backbone is trained for each incremental period using the spatio-temporal graph data of that period, with the corresponding model predicting the results for the current test set. **❷ Online:** The contextual pattern bank is removed. Similar to iTransformer, the model is trained on the complete node data of the current spatio-temporal graph and initialized with the model from the previous period, allowing for adjustments across the entire model. **❸ w/o Backbone:** The contextual pattern bank is retained, but the spatio-temporal backbone is replaced with the ones used in TrafficStream, STKEC, and EAC—i.e., replacing F_{reNet} and $DLGA$ with CNN and GCN. **❹ w/o DLGA:** The DLGA module in the spatio-temporal backbone is ablated. **❺ EAC:** We also include EAC, which follows a similar approach, for comparison in the ablation study.

Ablation findings. As shown in Figure 4, the ablation results demonstrate that parameter expansion in the contextual pattern bank, together with spatio-temporal pattern distinction and prompt guidance, is essential for alleviating catastrophic forgetting in continual learning. The performance of the **Retrain** and **Online** variants supports this conclusion. Notably, even without the contextual pattern

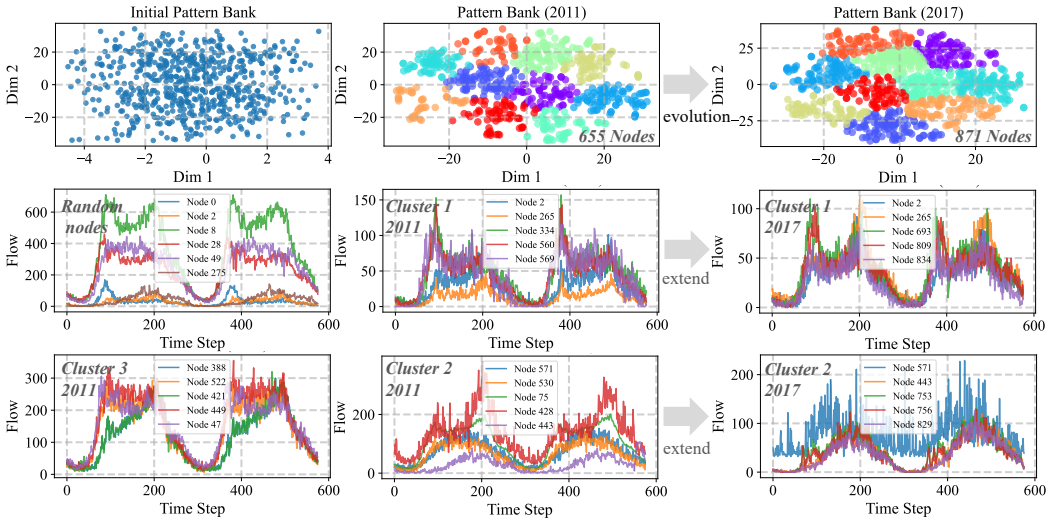


Figure 6: Case study on PEMS-Stream.

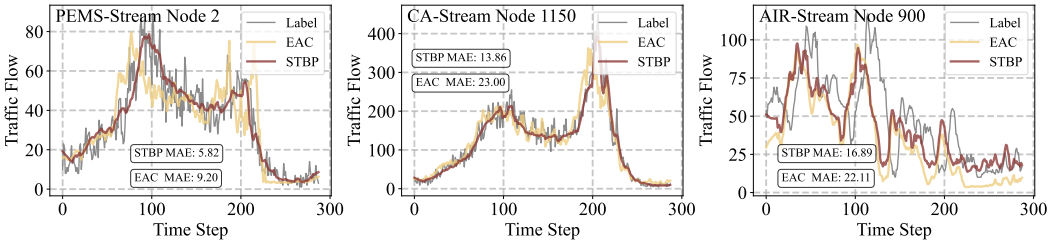


Figure 7: Visualization of real forecasting results.

bank on traffic datasets, the spatio-temporal backbone alone attains performance comparable to **EAC** under online training, highlighting the critical role of real-time dynamic correlation modeling and temporal distribution-drift mitigation in adapting to new incremental tasks. The performance drop observed in the **w/o Backbone** variant further confirms the indispensability of the general backbone and highlights the portability and adaptability of the pattern bank across different backbone architectures. Moreover, removing the **DLGA** module leads to significant performance degradation, validating its role in capturing dynamic spatial correlations and integrating prompt-based knowledge. The **FreNet** module also makes a notable contribution by improving computational efficiency and enhancing the extraction of stable temporal components.

Parameter Sensitivity Analysis. Additionally, we perform a sensitivity analysis on the adjustable hyperparameter d in **STBP**. In **STBP**, d represents the feature dimension for each module’s feature mapping, as well as the feature dimension of parameters in the contextual pattern bank. The analysis results are shown in Figure 5. Increasing d enhances the model’s overall parameter count and improves its expressive power. However, the performance gains from increasing d do not grow indefinitely; after reaching a certain threshold, the performance gain stabilizes. Further increases in d not only fail to improve performance but may also lead to negative effects, causing parameter redundancy. More parameter sensitivity analysis can be found in Appendix A.4.5.

5.4 CASE STUDY

To illustrate the distinction and expandability of the contextual pattern bank in **STBP**, we apply t-SNE to reduce the dimensionality of $\mathbf{P}_\tau \in \mathbb{R}^{N_\tau \times d}$ on the PEMS-Stream dataset. As shown in Figure 6, each point represents a graph node. Initially untrained, the pattern bank shows a chaotic distribution. After incremental training, clear clusters emerge. Nodes within the same cluster exhibit similar periodic and trend patterns in their traffic data, while those in different clusters (e.g., Clusters 1–3) show distinct behaviors. New nodes from later stages (e.g., Nodes 693, 809, 834 in 2017) are

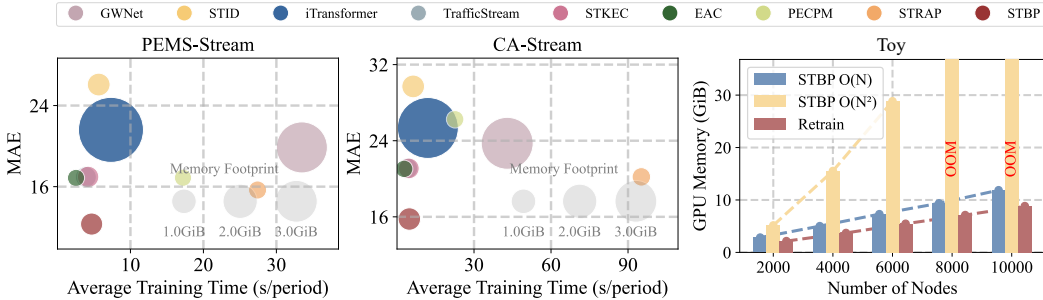


Figure 8: Efficiency comparison. $STBP O(N^2)$ denotes the version without linear attention, and Retrain refers to removing the contextual pattern bank.

correctly grouped into existing clusters, demonstrating that the pattern bank effectively distinguishes and generalizes spatio-temporal patterns through parameter fine-tuning, enabling continual adaptation.

In addition, we conduct an intuitive comparison of the forecasting performance between $STBP$ and EAC in real-world application scenarios. As shown in Figure 7, we select representative nodes from three datasets for visualization. Compared to EAC , $STBP$ more accurately captures dynamic trends, and its predictions demonstrate higher practical relevance in real-world continual learning environments. Additional case studies on other datasets can be found in Appendix A.4.6.

5.5 EFFICIENCY STUDY

An effective CSTF method must balance scalability, computational cost, and performance. We evaluate the efficiency of $STBP$ against baselines under the same settings. As shown in Figure 8, the average computational cost per period on PEMS-Stream and AIR-Stream is reported, with scatter size indicating GPU memory usage. We further analyze the impact of linear attention, full attention, and removal of the contextual pattern bank using a toy dataset. Results indicate that non-continual methods—such as $GWNet$, $STID$, and $iTransformer$ —require global parameter adjustments at each phase, impairing efficiency. $iTransformer$, in particular, incurs high memory overhead due to quadratic attention complexity. Even lightweight non-continual models exhibit limited efficiency in incremental training.

In contrast, CSTF methods such as EAC , $TrafficStream$, and $STKEC$ achieve higher efficiency through lightweight backbones and localized parameter tuning. While $PECPM$ and $STRAP$ maintain low memory usage, their training speeds remain modest. Despite its more complex backbone, $STBP$ incurs only minimal overhead compared to models like EAC , thanks to optimizations including frequency-domain processing and linear attention. This enables $STBP$ to deliver substantial performance gains with negligible cost increase. Results on the toy dataset confirm that linear attention reduces computational load effectively. As node count grows, the contextual pattern bank introduces only linear additional cost through its lightweight interaction with the backbone, avoiding exponential overhead. Furthermore, on CA-Stream, $STBP$ maintains state-of-the-art performance even under drastic graph expansion, demonstrating strong scalability.

6 CONCLUSION

In this work, we propose $STBP$, a novel framework for continual spatio-temporal forecasting. By combining a general-purpose backbone with a scalable contextual pattern bank, $STBP$ efficiently mitigates catastrophic forgetting while capturing dynamic spatio-temporal correlations. It adapts to evolving urban data without retraining from scratch, making it suitable for real-time applications. Validated on multiple datasets, $STBP$ demonstrates strong continual learning capabilities. Nevertheless, $STBP$ currently supports continual learning in a single-task setting. In the future, we plan to extend its application to cross-domain continual spatio-temporal forecasting, which will be a crucial step towards developing a foundational spatio-temporal model.

ACKNOWLEDGMENTS

This work was partly supported by the National Key Research and Development Program of China under Grant 2022YFB4501704, the National Natural Science Foundation of China under Grant 72342026, and Fundamental Research Funds for the Central Universities under Grant 2024-6-ZD-02.

REFERENCES

- Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. volume 33, pp. 17804–17815, 2020.
- Lucas Caccia, Eugene Belilovsky, Massimo Caccia, and Joelle Pineau. Online learned continual compression with adaptive quantization modules. In *International conference on machine learning*, pp. 1240–1250. PMLR, 2020.
- Chao Chen, Karl Petty, Alexander Skabardonis, Pravin Varaiya, and Zhanfeng Jia. Freeway performance measurement system: mining loop detector data. *Transportation research record*, 1748(1): 96–102, 2001.
- Wei Chen and Yuxuan Liang. Expand and compress: Exploring tuning principles for continual spatio-temporal graph forecasting. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Xu Chen, Junshan Wang, and Kunqing Xie. Trafficstream: A streaming traffic flow forecasting framework based on graph neural networks and continual learning. In Zhi-Hua Zhou (ed.), *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pp. 3620–3626, 8 2021. doi: 10.24963/ijcai.2021/498. URL <https://doi.org/10.24963/ijcai.2021/498>.
- Zheng Dong, Renhe Jiang, Haotian Gao, Hangchen Liu, Jinliang Deng, Qingsong Wen, and Xuan Song. Heterogeneity-informed meta-parameter learning for spatiotemporal time series forecasting. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 631–641, 2024.
- Yuchen Fang, Yanjun Qin, Haiyong Luo, Fang Zhao, and Kai Zheng. Stwave+: A multi-scale efficient spectral graph attention network with long-term trends for disentangled traffic flow forecasting. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- Yuchen Fang, Hao Miao, Yuxuan Liang, Liwei Deng, Yue Cui, Ximu Zeng, Yuyang Xia, Yan Zhao, Torben Bach Pedersen, Christian S. Jensen, Xiaofang Zhou, and Kai Zheng. Unraveling Spatio-Temporal Foundation Models Via the Pipeline Lens: A Comprehensive Review. *IEEE Transactions on Knowledge & Data Engineering*, (01):1–24, January 2026. ISSN 1558-2191. doi: 10.1109/TKDE.2026.3651536. URL <https://doi.ieeecomputersociety.org/10.1109/TKDE.2026.3651536>.
- Haotian Gao, Renhe Jiang, Zheng Dong, Jinliang Deng, Yuxin Ma, and Xuan Song. Spatial-temporal-decoupled masked pre-training for spatiotemporal forecasting. In Kate Larson (ed.), *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pp. 3998–4006. International Joint Conferences on Artificial Intelligence Organization, 8 2024. doi: 10.24963/ijcai.2024/442. URL <https://doi.org/10.24963/ijcai.2024/442>. Main Track.
- Danlei Hu, Lu Chen, Hanxi Fang, Ziquan Fang, Tianyi Li, and Yunjun Gao. Spatio-temporal trajectory similarity measures: A comprehensive survey and quantitative study. *IEEE Transactions on Knowledge and Data Engineering*, 36(5):2191–2212, 2023.
- Jiahao Ji, Wentao Zhang, Jingyuan Wang, and Chao Huang. Seeing the unseen: Learning basis confounder representations for robust traffic prediction. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 1*, pp. 577–588, 2025.
- Jiawei Jiang, Chengkai Han, Wayne Xin Zhao, and Jingyuan Wang. Pdformer: Propagation delay-aware dynamic long-range transformer for traffic flow prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023a.

- Renhe Jiang, Zhaonan Wang, Jiawei Yong, Puneet Jeph, Quanjun Chen, Yasumasa Kobayashi, Xuan Song, Shintaro Fukushima, and Toyotaro Suzumura. Spatio-temporal meta-graph learning for traffic forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 8078–8086, 2023b.
- Ming Jin, Huan Yee Koh, Qingsong Wen, Daniele Zambon, Cesare Alippi, Geoffrey I Webb, Irwin King, and Shirui Pan. A survey on graph neural networks for time series: Forecasting, classification, imputation, and anomaly detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning (ICML)*, pp. 5156–5165. PMLR, 2020.
- Weiyang Kong, Ziyu Guo, and Yubao Liu. Spatio-temporal pivotal graph neural networks for traffic flow forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(8):8627–8635, Mar. 2024. doi: 10.1609/aaai.v38i8.28707. URL <https://ojs.aaai.org/index.php/AAAI/article/view/28707>.
- Rahul Kumar, Manish Bhanu, João Mendes-Moreira, and Joydeep Chandra. Spatio-temporal predictive modeling techniques for different domains: a survey. *ACM Computing Surveys*, 57(2):1–42, 2024.
- Sanghyun Lee and Chanyoung Park. Continual traffic forecasting via mixture of experts. *arXiv preprint arXiv:2406.03140*, 2024.
- Fuxian Li, Jie Feng, Huan Yan, Guangyin Jin, Fan Yang, Funing Sun, Depeng Jin, and Yong Li. Dynamic graph convolutional recurrent network for traffic prediction: Benchmark and solution. *ACM Transactions on Knowledge Discovery from Data*, 17(1):1–21, 2023.
- Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations (ICLR)*, 2018.
- Aoyu Liu and Yaying Zhang. An efficient spatial-temporal transformer with temporal aggregation and spatial memory for traffic forecasting. *Expert Systems with Applications*, 250:123884, 2024a.
- Aoyu Liu and Yaying Zhang. Spatial-temporal dynamic graph convolutional network with interactive learning for traffic forecasting. *IEEE Transactions on Intelligent Transportation Systems*, 2024b.
- Aoyu Liu and Yaying Zhang. Crossst: An efficient pre-training framework for cross-district pattern generalization in urban spatio-temporal forecasting. In *41th IEEE International Conference on Data Engineering*, 2025.
- Chenxi Liu, Sun Yang, Qianxiong Xu, Zhishuai Li, Cheng Long, Ziyue Li, and Rui Zhao. Spatial-temporal large language model for traffic prediction. In *25th IEEE International Conference on Mobile Data Management (MDM)*, pp. 31–40, 2024a.
- Chenxi Liu, Kethmi Hirushini Hettige, Qianxiong Xu, Cheng Long, Shili Xiang, Gao Cong, Ziyue Li, and Rui Zhao. ST-LLM+: Graph enhanced spatio-temporal large language models for traffic prediction. *IEEE Transactions on Knowledge and Data Engineering*, 37(8):4846–4859, 2025a.
- Chenxi Liu, Qianxiong Xu, Hao Miao, Sun Yang, Lingzheng Zhang, Cheng Long, Ziyue Li, and Rui Zhao. TimeCMA: Towards llm-empowered multivariate time series forecasting via cross-modality alignment. In *AAAI*, volume 39, pp. 18780–18788, 2025b.
- Hangchen Liu, Zheng Dong, Renhe Jiang, Jiewen Deng, Q Chen, and X Song. Staformer: Spatio-temporal adaptive embedding makes vanilla transformers sota for traffic forecasting. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM)*, pp. 21–25, 2023a.
- Xu Liu, Yutong Xia, Yuxuan Liang, Junfeng Hu, Yiwei Wang, Lei Bai, Chao Huang, Zhenguang Liu, Bryan Hooi, and Roger Zimmermann. Largest: A benchmark dataset for large-scale traffic forecasting. In *Advances in Neural Information Processing Systems*, 2023b.

- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *International Conference on Learning Representations (ICLR)*, 2024b.
- Jiaming Ma, Bingwu Wang, Pengkun Wang, Zhengyang Zhou, Xu Wang, and Yang Wang. Robust spatio-temporal centralized interaction for ood learning. In *Proceedings of the Forty-Second International Conference on Machine Learning (ICML)*, 2025a.
- Minbo Ma, Kai Tang, Huan Li, Fei Teng, Dalin Zhang, and Tianrui Li. Beyond fixed variables: Expanding-variate time series forecasting via flat scheme and spatio-temporal focal learning. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 2054–2065, 2025b.
- Hao Miao, Ziqiao Liu, Yan Zhao, Chenjuan Guo, Bin Yang, Kai Zheng, and Christian S Jensen. Less is more: Efficient time series dataset condensation via two-fold modal matching. *Proceedings of the VLDB Endowment*, 18(2):226–238, 2024a.
- Hao Miao, Yan Zhao, Chenjuan Guo, Bin Yang, Kai Zheng, Feiteng Huang, Jiandong Xie, and Christian S Jensen. A unified replay-based continuous learning framework for spatio-temporal prediction on streaming data. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, pp. 1050–1062. IEEE, 2024b.
- Hao Miao, Yan Zhao, Chenjuan Guo, Bin Yang, Kai Zheng, and Christian S Jensen. Spatio-temporal prediction on streaming data: A unified federated continuous learning framework. *IEEE Transactions on Knowledge and Data Engineering*, 2025.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- Zezhi Shao, Zhao Zhang, Fei Wang, Wei Wei, and Yongjun Xu. Spatial-temporal identity: A simple yet effective baseline for multivariate time series forecasting. In *Proceedings of the 31st ACM international conference on information & knowledge management*, pp. 4454–4458, 2022a.
- Zezhi Shao, Zhao Zhang, Wei Wei, Fei Wang, Yongjun Xu, Xin Cao, and Christian S Jensen. Decoupled dynamic spatial-temporal graph neural network for traffic forecasting. *arXiv preprint arXiv:2206.09112*, 2022b.
- Zezhi Shao, Fei Wang, Yongjun Xu, Wei Wei, Chengqing Yu, Zhao Zhang, Di Yao, Tao Sun, Guangyin Jin, Xin Cao, et al. Exploring progress in multivariate time series forecasting: Comprehensive benchmarking and heterogeneity analysis. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- Chao Song, Youfang Lin, Shengnan Guo, and Huaiyu Wan. Spatial-temporal synchronous graph convolutional networks: A new framework for spatial-temporal network data forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 914–921, 2020.
- Jindong Tian, Yuxuan Liang, Ronghui Xu, Peng Chen, Chenjuan Guo, Aoying Zhou, Lujia Pan, Zhongwen Rao, and Bin Yang. Air quality prediction with physics-guided dual neural odes in open systems. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- Binwu Wang, Yudong Zhang, Jiahao Shi, Pengkun Wang, Xu Wang, Lei Bai, and Yang Wang. Knowledge expansion and consolidation for continual traffic prediction with expanding graphs. *IEEE Transactions on Intelligent Transportation Systems*, 24(7):7190–7201, 2023a.
- Binwu Wang, Yudong Zhang, Xu Wang, Pengkun Wang, Zhengyang Zhou, Lei Bai, and Yang Wang. Pattern expansion and consolidation on evolving graphs for continual traffic prediction. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, pp. 2223–2232, New York, NY, USA, 2023b. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599463. URL <https://doi.org/10.1145/3580305.3599463>.

- Binwu Wang, Jiaming Ma, Pengkun Wang, Xu Wang, Yudong Zhang, Zhengyang Zhou, and Yang Wang. Stone: A spatio-temporal ood learning framework kills both spatial and temporal shifts. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2948–2959, 2024.
- Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-temporal graph modeling. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1907–1913, 2019.
- Mingyuan Xia, Chunxu Zhang, Zijian Zhang, Hao Miao, Qidong Liu, Yuanshao Zhu, and Bo Yang. Timeemb: A lightweight static-dynamic disentanglement framework for time series forecasting. In *NeurIPS*, 2025.
- Yutong Xia, Yuxuan Liang, Haomin Wen, Xu Liu, Kun Wang, Zhengyang Zhou, and Roger Zimmermann. Deciphering spatio-temporal graph forecasting: A causal lens and treatment. *Advances in Neural Information Processing Systems*, 36:37068–37088, 2023.
- Bin Yang, Yuxuan Liang, Chenjuan Guo, and Christian S Jensen. Data driven decision making with time series and spatio-temporal data. In *2025 IEEE 41st International Conference on Data Engineering (ICDE)*, pp. 4517–4522. IEEE Computer Society, 2025.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. In *Proceedings of the 33rd International Joint Conference on Artificial Intelligence (IJCAI)*, 2018.
- Haoyu Zhang, Wentao Zhang, Hao Miao, Xinke Jiang, Yuchen Fang, and Yifan Zhang. Strap: Spatio-temporal pattern retrieval for out-of-distribution generalization. *The Thirty-Ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Zijian Zhang, Xiangyu Zhao, Qidong Liu, Chunxu Zhang, Qian Ma, Wanyu Wang, Hongwei Zhao, Yiqi Wang, and Zitao Liu. Promptst: Prompt-enhanced spatio-temporal multi-attribute prediction. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pp. 3195–3205, 2023.
- Chuanpan Zheng, Xiaoliang Fan, Cheng Wang, and Jianzhong Qi. Gman: A graph multi-attention network for traffic prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 1234–1241, 2020.
- Chuanpan Zheng, Xiaoliang Fan, Shirui Pan, Haibing Jin, Zhaopeng Peng, Zonghan Wu, Cheng Wang, and S Yu Philip. Spatio-temporal joint graph convolutional networks for traffic forecasting. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- Zhengyang Zhou, Qihe Huang, Kuo Yang, Kun Wang, Xu Wang, Yudong Zhang, Yuxuan Liang, and Yang Wang. Maintaining the status quo: Capturing invariant relations for ood spatiotemporal learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, pp. 3603–3614, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599421. URL <https://doi.org/10.1145/3580305.3599421>.

A APPENDIX

A.1 NOTATIONS

Table 3 summarizes the notations frequently used throughout this manuscript.

Table 3: The notations that are commonly used in the manuscript.

Notation	Definition
$\mathbb{G} = \{G_\tau\}_{\tau=1}^T$	Streaming spatio-temporal graph
\mathbf{X}_τ	Inputs for the τ period
\mathbf{Y}_τ	Prediction for the τ period
\mathbf{P}_τ	contextual pattern bank for the τ period
\mathbf{P}'_τ	Expanded contextual pattern bank
\mathbf{H}_τ	Hidden representation for the τ period
\mathcal{M}_θ	Spatio-Temporal backbone
\mathbf{F}_τ	Frequency domain embedding
\mathbf{H}_τ^f	Representation after frequency-domain processing
\mathbf{W}_q	Trainable parameter weights
\mathbf{W}_k	Trainable parameter weights
\mathbf{W}_v	Trainable parameter weights
$\phi(\cdot)$	Random mapping function
\mathbf{H}_τ^s	Input of the DLGA module
$\mathbf{H}_\tau^{s'}$	Output of the DLGA module

A.2 RELATED WORK DETAILS

A.2.1 SPATIO-TEMPORAL FORECASTING

Spatio-temporal forecasting aims to support decision-making in critical domains such as intelligent transportation and smart cities by uncovering dynamic correlations embedded in spatio-temporal data. These data typically exhibit strong spatial-temporal correlations and pronounced heterogeneity. In recent years, deep learning-based STGNNs have emerged as effective tools for such forecasting tasks. STGNNs generally employ temporal modules (e.g., recurrent neural networks (RNNs) (Li et al., 2018; Jiang et al., 2023b; Shao et al., 2022b) and convolutional neural networks (CNNs)) (Yu et al., 2018; Liu & Zhang, 2024b;a) to capture temporal correlations, while leveraging spatial modules (e.g., graph neural networks (GNNs)) (Veličković et al., 2018; Song et al., 2020) to model spatial relationships.

Early STGNNs, such as STGCN (Yu et al., 2018) and DCRNN (Li et al., 2018), combined basic temporal and spatial components for forecasting tasks, often relying on predefined geographic adjacency matrices. However, these static assumptions hinder their ability to model dynamically changing spatial correlations in a data-driven manner. Subsequent works—such as GWNet (Wu et al., 2019), DGCRN (Li et al., 2023), and MegaCRN (Jiang et al., 2023b)—introduced adaptive adjacency matrices or learned spatial correlations directly from data, significantly improving prediction accuracy. More recent advances, including STID (Shao et al., 2022a), STAEformer (Liu et al., 2023a), and HimNet (Dong et al., 2024), have highlighted the importance of spatial pattern distinction in enhancing forecasting performance. These models incorporate trainable mechanisms such as spatial embeddings, parameter pools, and contextual pattern banks to distinguish spatial patterns more precisely, thereby improving both accuracy and adaptability.

Despite these advancements, most existing STGNNs are built on static assumptions and are not designed to operate in dynamic, continually evolving spatio-temporal environments—limiting their applicability in continual learning scenarios.

A.2.2 CONTINUAL SPATIO-TEMPORAL LEARNING

Early research in continual learning primarily focused on computer vision (Lee & Park, 2024; Miao et al., 2024a) and natural language processing (Caccia et al., 2020; Xia et al., 2025). With the rapid development of IoT and intelligent transportation systems, attention has increasingly shifted toward CSTF (Chen et al., 2021; Wang et al., 2023a; Chen & Liang, 2025; Wang et al., 2023b; Miao et al., 2025), which addresses the challenges of dynamically evolving and expanding spatio-temporal data. CSTF aims to enable models to continuously learn and adapt to new patterns and knowledge in changing environments, while minimizing forgetting of previously acquired information or performance degradation.

One of the earliest frameworks in this domain, TrafficStream (Chen et al., 2021), pioneered the integration of spatio-temporal modeling with continual learning. It employed strategies such as historical data replay and parameter smoothing to handle long-term streaming traffic data, achieving accurate traffic flow forecasting. Subsequently, the STKEC (Wang et al., 2023a) introduced an influence-based knowledge expansion strategy and a memory-augmented knowledge consolidation mechanism to better accommodate the growth of transportation networks while mitigating catastrophic forgetting. The PECPM (Wang et al., 2023b) framework employs a pattern matching mechanism to maintain and dynamically update a bank of representative traffic patterns from evolving road networks, enabling efficient continual learning without historical data and improving both prediction accuracy and training efficiency. Meanwhile, STRAP (Zhang et al., 2025) introduces a retrieval-augmented approach that builds multi-dimensional pattern libraries. During inference, it retrieves and fuses relevant historical patterns with current inputs via a plug-and-play prompting mechanism, effectively boosting generalization in OOD scenarios while mitigating catastrophic forgetting.

The EAC (Chen & Liang, 2025) further advanced the field by incorporating prompt tuning, enabling CSTF with a small number of trainable parameters. Its dynamic prompt pool, which supports both “expansion” and “compression,” enhances adaptability to new nodes while preserving historical knowledge, improving both generalization and computational efficiency. In addition, the UFCL (Miao et al., 2025) leveraged federated learning to preserve data privacy and introduced a global replay buffer for synthetic spatio-temporal data, addressing the challenges of distributed streaming environments. Despite these advancements, most existing methods primarily focus on alleviating knowledge forgetting, while overlooking the critical role of the spatio-temporal backbone in continual learning scenarios.

Advancements Beyond Existing Prompt Methods. Unlike EAC’s “expand-and-compress” prompt pool that may lead to historical information loss during compression, our contextual pattern bank adopts pure parametric incremental expansion without compression, more completely preserving historical knowledge. Additionally, while EAC’s prompt interaction is relatively simple (e.g., feature addition), our pattern bank employs structured multi-component design ($\mathbf{P}_\tau^{(0)}$, $\mathbf{P}_\tau^{(1)}$, $\mathbf{P}_\tau^{(2)}$) that jointly models node relevance and heterogeneity through gating and attention mechanisms, enabling more comprehensive spatio-temporal representation learning.

A.3 FURTHER METHODS DETAILS

A.3.1 APPROXIMATION DERIVATION OF EQ. 9

An approximate derivation of the attention mechanism in the dual-stream linear graph attention is presented below:

$$\begin{aligned}
 \text{Attention}(\mathbf{q}_u, \mathbf{k}_v, \mathbf{v}_v, \mathbf{p}_v) &= \sum_{v=1}^N \frac{\exp(\mathbf{q}_u^\top \mathbf{k}_v) \mathbf{v}_v}{\sum_{w=1}^N \exp(\mathbf{q}_u^\top \mathbf{k}_w)} + \sum_{v=1}^N \frac{\exp(\mathbf{q}_u^\top \mathbf{p}_v) \mathbf{v}_v}{\sum_{w=1}^N \exp(\mathbf{q}_u^\top \mathbf{p}_w)} \\
 &\approx \frac{\sum_{v=1}^N \phi(\mathbf{q}_u)^\top \phi(\mathbf{k}_v) \mathbf{v}_v}{\sum_{w=1}^N \phi(\mathbf{q}_u)^\top \phi(\mathbf{k}_w)} + \frac{\sum_{v=1}^N \phi(\mathbf{q}_u)^\top \phi(\mathbf{p}_v) \mathbf{v}_v}{\sum_{w=1}^N \phi(\mathbf{q}_u)^\top \phi(\mathbf{p}_w)} \quad (10) \\
 &= \underbrace{\begin{bmatrix} \phi(\mathbf{q}_u)^\top \sum_{v=1}^N \phi(\mathbf{k}_v) \mathbf{v}_v^\top \\ \phi(\mathbf{q}_u)^\top \sum_{w=1}^N \phi(\mathbf{k}_w) \end{bmatrix}}_{\text{Term 1: Representation-based aggregation}} + \underbrace{\begin{bmatrix} \phi(\mathbf{q}_u)^\top \sum_{v=1}^N \phi(\mathbf{p}_v) \mathbf{v}_v^\top \\ \phi(\mathbf{q}_u)^\top \sum_{w=1}^N \phi(\mathbf{p}_w) \end{bmatrix}}_{\text{Term 2: Prompt-based aggregation}}
 \end{aligned}$$

Algorithm 1 The workflow of STBP for continual spatio-temporal forecasting

Require: Spatio-temporal backbone \mathcal{M}_θ , contextual pattern bank $\{\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_\tau\}$, streaming train data $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_\tau\}$.

Ensure: Optimized backbone \mathcal{M}_{θ^*} and contextual pattern bank $\{\mathbf{P}_1^*, \dots, \mathbf{P}_\tau^*\}$

Initialize: $\mathcal{M}_\theta \leftarrow \{\}$, $\mathbf{P}_1 \leftarrow \{\}$

for each period i in $\{1, 2, 3, \dots, \tau\}$ **do**

if $i == 1$ **then**

\triangleright *Initial training phase* \triangleleft

 Construct the initial contextual pattern bank \mathbf{P}_1

 Optimize backbone and contextual pattern bank with initial data \mathbf{X}_1 :

$(\mathcal{M}_{\theta^*}, \mathbf{P}_1^*) \leftarrow \operatorname{argmin}_\theta \mathcal{M}_\theta(\mathbf{X}_1, \mathbf{P}_1)$

else

\triangleright *Streaming learning phase* \triangleleft

 Expand contextual pattern bank \mathbf{P}_i : $\mathbf{P}_i \leftarrow \mathbf{P}_{i-1} \parallel \Delta \mathbf{P}_i$

 Inherit parameters: $(\mathcal{M}_\theta, \mathbf{P}_i) \leftarrow (\mathcal{M}_{\theta^*}, \mathbf{P}_{i-1}^*)$

 Freeze backbone parameters \mathcal{M}_θ : $\theta \leftarrow \operatorname{freeze}(\theta)$

 Fine-tune \mathbf{P}_i with backbone \mathcal{M}_θ on \mathbf{X}_i :

$\mathbf{P}_i^* \leftarrow \operatorname{argmin}_\theta \mathcal{M}_\theta(\mathbf{X}_i, \mathbf{P}_i)$

Table 4: Overview of continual spatio-temporal forecasting datasets.

Dataset	Domain	Time Range	Period	Node Expansion	Frequency
PEMS-Stream	Traffic	07/10/2011 - 09/08/2017	7	655 \rightarrow 715 \rightarrow 786 \rightarrow 822 \rightarrow 834 \rightarrow 850 \rightarrow 871	5 min
CA-Stream	Traffic	01/01/2019 - 04/30/2019	4	480 \rightarrow 691 \rightarrow 1175 \rightarrow 1698	5 min
AIR-Stream	Air Quality	01/01/2016 - 12/31/2019	4	1087 \rightarrow 1154 \rightarrow 1193 \rightarrow 1202	1 hour

where q_u is the query tensor of node u ; k_v and v_v are the key and value tensors of node v , respectively; and p_v represents the prompt information for node v .

A.3.2 ALGORITHM WORKFLOW

The overall workflow of STBP for continual spatio-temporal forecasting is presented in a more intuitive manner in Algorithm 1.

A.4 ADDITIONAL EXPERIMENT DETAILS

A.4.1 DATASET DETAILS

Table 4 and Table 5 jointly summarize the characteristics of the three continual spatio-temporal datasets used in this study: **PEMS-Stream**, **CA-Stream**, and **AIR-Stream**. These datasets differ in domain (traffic¹ vs. air quality²), temporal span, and topological evolution, collectively covering a broad spectrum of real-world non-stationary scenarios suitable for evaluating continual learning models. PEMS-Stream contains highway traffic sensor readings collected across California from July 2011 to September 2017. It spans seven periods with a gradual increase in the number of sensor nodes—from 655 to 871—resulting in a +33% relative growth. This dataset simulates realistic, long-term infrastructure expansion and serves as a benchmark for evaluating model adaptability under progressive and stable topological changes. CA-Stream, also in the traffic domain, covers a much

¹<https://dot.ca.gov/programs/traffic-operations/mpr/pems-source>

²<https://air.cnemc.cn:18007/>

Table 5: Topological dynamics and evaluation purposes of the datasets.

Dataset	Topology Change	Δ Nodes	Relative Change	Primary Purpose
PEMS-Stream	Gradual expansion	216	+33%	Realistic progressive growth
CA-Stream	Explosive expansion	1,218	+254%	Extreme incremental stress test
AIR-Stream	Stable expansion	115	+10%	Cross-domain validation

shorter period (January to April 2019) but features a sharp and sudden node expansion—from 480 to 1,698—corresponding to a +254% relative increase.

This explosive growth introduces significant distributional shifts, making CA-Stream a challenging testbed for assessing model robustness under rapidly evolving conditions. AIR-Stream focuses on urban air quality and environmental measurements from 2016 to 2019. It exhibits modest but steady node growth—from 1,087 to 1,202 (+10%)—and represents a relatively stable expansion setting. Its distinct domain and smoother structural changes make it particularly suitable for evaluating cross-domain generalization and robustness to gradual environmental variation.

To further assess non-stationarity, we conduct Maximum Mean Discrepancy (MMD) tests across different periods, separately evaluating *original nodes* (present from the beginning) and *added nodes* (introduced during expansion), as shown in Table 6. A distribution shift is considered significant when $\text{MMD} > 0.1$ or $p < 0.05$. Across all datasets, added nodes consistently exhibit stronger distributional shifts, reflecting the spatial disruptions caused by topological expansion. For instance, CA-Stream shows a substantial shift for added nodes ($\text{MMD} = 0.3361$, $p = 0.0010$), consistent with its rapid growth. Interestingly, AIR-Stream records the highest MMD among original nodes (0.3324, $p = 0.001$), despite minimal structural change—indicating notable temporal drift in environmental data. This highlights AIR-Stream’s importance for evaluating robustness to evolving distributions even under stable topology. By contrast, PEMS-Stream shows only moderate drift among original nodes ($\text{MMD} = 0.0939$), aligning with its smoother expansion. CA-Stream presents weaker drift in original nodes ($\text{MMD} = 0.0792$, $p = 0.1119$), likely due to its limited temporal span. These results underscore the dual challenge in continual spatio-temporal learning: managing both spatial shifts induced by node expansion and temporal non-stationarity inherent to dynamic environments, with their nature and intensity varying across domains.

Table 6: Distribution shift analysis based on MMD tests.

Type	PEMS-Stream		AIR-Stream		CA-Stream	
	MMD	p	MMD	p	MMD	p
Original Node	0.0939	0.008	0.3324	0.001	0.0792	0.1119
Added Node	0.2958	0.001	0.2679	0.001	0.3361	0.0010

A.4.2 BASELINES AND METRICS DETAILS

In this paper, we provide a detailed comparison with two categories of representative models:

Conventional Spatio-Temporal Forecasting Models. ① **GWNet** (Wu et al., 2019): A STGNN model based on an adaptive adjacency matrix that can adaptively capture latent spatial dependencies. This model combines graph convolutional networks and temporal convolutions to effectively capture spatio-temporal correlations in the data. ② **STID** (Shao et al., 2022a): An efficient multilayer perceptron model that solves the problem of sample non-separability using trainable embeddings, showing outstanding performance in spatio-temporal forecasting tasks. ③ **iTransformer** (Liu et al., 2024b): A time-series model that does not rely on a static graph structure. By modeling the interactions between variables, it captures temporal features and is effectively applied to multivariate time series forecasting tasks.

Continual Spatio-Temporal Forecasting Models. These models are designed to handle time-varying data and are suitable for continual training tasks. Like STBP, they belong to the category of continual learning models. We selected the following three representative models for comparison: ① **TrafficStream** (Chen et al., 2021): The first model for CSTF, it employs a traffic pattern fusion approach, historical data replay, and parameter smoothing strategies to efficiently integrate and learn new spatio-temporal patterns in the continuously expanding and evolving traffic network. ② **STKEC** (Wang et al., 2023a): A traffic forecasting model based on the continual learning paradigm. Through an influence-based knowledge expansion strategy and a memory-augmented knowledge

Table 7: Comparison of prediction performance for each incremental period on PEMS-Stream. **Bold**: best, underline: second best.

Model	Metric	PEMS-Stream Period							
		2011	2012	2013	2014	2015	2016	2017	Avg.
GWNet	MAE	25.90 \pm 1.19	20.51 \pm 2.14	18.59 \pm 1.31	16.17 \pm 0.39	19.37 \pm 1.73	17.58 \pm 0.83	20.95 \pm 1.03	19.87 \pm 0.10
	RMSE	44.79 \pm 1.08	32.43 \pm 2.60	29.32 \pm 1.47	25.37 \pm 0.60	31.37 \pm 2.11	29.89 \pm 1.14	34.96 \pm 1.12	32.59 \pm 0.18
	MAPE (%)	29.79 \pm 1.61	29.18 \pm 2.95	28.89 \pm 2.63	25.45 \pm 2.51	29.61 \pm 4.28	24.79 \pm 1.44	26.86 \pm 2.42	27.79 \pm 0.76
STID	MAE	32.68 \pm 1.34	26.10 \pm 0.97	25.36 \pm 2.28	22.78 \pm 1.25	23.37 \pm 1.38	24.28 \pm 1.09	27.98 \pm 1.08	26.07 \pm 0.23
	RMSE	53.22 \pm 2.16	41.02 \pm 1.64	39.50 \pm 3.50	35.67 \pm 2.07	37.15 \pm 2.34	40.94 \pm 1.78	44.20 \pm 1.39	41.67 \pm 0.21
	MAPE (%)	44.33 \pm 5.66	38.68 \pm 1.13	41.40 \pm 0.58	39.50 \pm 3.34	42.05 \pm 5.09	39.65 \pm 6.86	42.05 \pm 4.57	41.09 \pm 2.49
iTransformer	MAE	25.44 \pm 3.24	20.90 \pm 0.70	20.37 \pm 0.58	20.58 \pm 0.85	20.23 \pm 0.71	20.83 \pm 0.80	22.83 \pm 1.04	21.60 \pm 0.79
	RMSE	39.73 \pm 4.15	33.32 \pm 1.23	32.39 \pm 1.00	33.25 \pm 1.35	33.04 \pm 1.04	35.70 \pm 1.22	36.69 \pm 1.64	34.88 \pm 1.17
	MAPE (%)	35.71 \pm 2.97	38.61 \pm 4.15	37.35 \pm 3.82	42.89 \pm 5.17	40.49 \pm 6.04	40.83 \pm 5.24	41.55 \pm 4.25	39.63 \pm 3.81
TrafficStream	MAE	18.15 \pm 0.19	16.81 \pm 0.36	16.16 \pm 0.12	16.62 \pm 0.14	16.39 \pm 0.01	16.47 \pm 0.15	18.04 \pm 0.09	16.95 \pm 0.03
	RMSE	27.75 \pm 0.23	26.72 \pm 0.58	25.64 \pm 0.13	26.96 \pm 0.06	27.29 \pm 0.02	28.93 \pm 0.05	29.31 \pm 0.10	27.52 \pm 0.05
	MAPE (%)	21.35 \pm 0.76	21.14 \pm 0.46	21.33 \pm 0.85	22.61 \pm 1.31	21.36 \pm 0.53	20.84 \pm 0.86	22.99 \pm 1.47	21.66 \pm 0.54
STKEC	MAE	<u>18.09</u> \pm 0.46	16.83 \pm 0.36	16.26 \pm 0.20	16.48 \pm 0.24	16.38 \pm 0.15	16.31 \pm 0.13	18.41 \pm 0.35	16.96 \pm 0.09
	RMSE	<u>27.47</u> \pm 0.47	26.89 \pm 0.56	25.74 \pm 0.21	26.97 \pm 0.40	27.17 \pm 0.22	28.70 \pm 0.38	30.03 \pm 0.68	27.56 \pm 0.11
	MAPE (%)	21.00 \pm 0.77	21.42 \pm 1.07	20.54 \pm 0.47	21.53 \pm 0.63	21.71 \pm 0.46	20.38 \pm 0.86	23.87 \pm 1.01	21.50 \pm 0.52
PECPM	MAE	18.43 \pm 0.41	16.91 \pm 0.43	16.03 \pm 0.23	16.27 \pm 0.12	16.09 \pm 0.13	16.21 \pm 0.23	18.05 \pm 0.39	16.86 \pm 0.12
	RMSE	28.09 \pm 0.39	26.94 \pm 0.64	25.48 \pm 0.43	26.49 \pm 0.28	26.71 \pm 0.28	28.55 \pm 0.24	29.31 \pm 0.63	27.37 \pm 0.20
	MAPE (%)	21.57 \pm 0.84	21.18 \pm 0.63	20.71 \pm 0.82	22.82 \pm 1.61	22.12 \pm 1.68	20.99 \pm 0.37	22.73 \pm 1.83	21.73 \pm 0.45
STRAP	MAE	18.18 \pm 0.08	17.40 \pm 0.56	16.07 \pm 0.11	16.30 \pm 0.06	16.04 \pm 0.06	16.16 \pm 0.13	18.02 \pm 0.12	16.88 \pm 0.10
	RMSE	27.72 \pm 0.04	27.60 \pm 0.82	25.46 \pm 0.16	26.49 \pm 0.06	26.53 \pm 0.07	28.34 \pm 0.18	29.30 \pm 0.19	27.35 \pm 0.13
	MAPE (%)	22.66 \pm 0.84	21.06 \pm 0.61	22.58 \pm 0.79	22.87 \pm 0.51	21.57 \pm 1.66	21.90 \pm 0.73	22.54 \pm 1.14	22.17 \pm 0.46
EAC	MAE	18.12 \pm 0.26	<u>15.41</u> \pm 0.40	<u>14.67</u> \pm 0.21	<u>15.09</u> \pm 0.15	<u>14.96</u> \pm 0.15	<u>14.78</u> \pm 0.21	<u>16.67</u> \pm 0.19	<u>15.67</u> \pm 0.20
	RMSE	27.72 \pm 0.70	<u>24.23</u> \pm 0.62	<u>23.08</u> \pm 0.32	<u>24.29</u> \pm 0.22	<u>24.58</u> \pm 0.25	<u>26.26</u> \pm 0.30	<u>26.96</u> \pm 0.30	<u>25.30</u> \pm 0.29
	MAPE (%)	<u>19.62</u> \pm 0.16	<u>19.90</u> \pm 0.51	<u>20.40</u> \pm 0.63	<u>20.01</u> \pm 0.32	<u>20.72</u> \pm 0.27	<u>19.98</u> \pm 0.70	<u>22.33</u> \pm 1.19	<u>20.42</u> \pm 0.41
STBP	MAE	14.29 \pm 0.05	12.13 \pm 0.11	11.60 \pm 0.09	11.71 \pm 0.07	11.61 \pm 0.06	11.38 \pm 0.09	13.42 \pm 0.17	12.31 \pm 0.07
	RMSE	22.00 \pm 0.07	19.47 \pm 0.18	18.58 \pm 0.15	19.54 \pm 0.10	19.55 \pm 0.07	21.91 \pm 0.07	22.56 \pm 0.26	20.52 \pm 0.11
	MAPE (%)	16.34 \pm 0.12	15.26 \pm 0.28	15.38 \pm 0.15	15.51 \pm 0.23	15.47 \pm 0.35	14.80 \pm 0.23	16.76 \pm 0.20	15.65 \pm 0.21

consolidation mechanism, STKEC helps the model effectively integrate new spatio-temporal traffic patterns in an ever-expanding road network while retaining previously learned spatio-temporal patterns. ⑤ **PECPM** (Wang et al., 2023b): A continual spatio-temporal forecasting model for evolving traffic networks, relying on a pattern-matching pattern bank to store representative patterns without full historical data. It fine-tunes with new/conflict nodes for knowledge expansion and uses preservation/traceability mechanisms to avoid forgetting. ④ **STRAP** (Zhang et al., 2025): A retrieval-augmented framework for OOD generalization, building a multi-dimensional key-value pattern library (spatial/temporal/spatio-temporal) during training. It retrieves similar patterns to fuse with current data in inference, achieving SOTA without task-specific fine-tuning. ⑥ **EAC** (Chen & Liang, 2025): A CSTF based on prompt tuning. By integrating a base STGNN with a continual prompt pool, it efficiently addresses incremental learning and catastrophic forgetting in streaming data using lightweight trainable parameters.

The Excluded Models. Some baselines that might be considered relevant for comparison were excluded, and we provide explanations for their exclusion below. ① **STAEformer** (Liu et al., 2023a): a widely recognized baseline, was not included in our comparison due to non-convergence observed when applying the same experimental setting as used for GWNet and STID on the selected three datasets. To ensure fair and unambiguous evaluation, we excluded it from the results and have provided the corresponding training logs in the code repository. ② **UFCL** (Miao et al., 2025): The CSTF method UFCL is not included in the comparison due to differences in experimental settings, which prevent a fair evaluation.

Metrics Details. Additionally, the performance metrics used in the experiments to evaluate the model, namely MAE, RMSE, and MAPE, are defined as follows:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (11)$$

Table 8: Comparison of prediction performance for each incremental period on CA-Stream and AIR-Stream. **Bold**: best, underline: second best.

Model	Metric	CA-Stream Period					AIR-Stream Period				
		Jan-19	Feb-19	Mar-19	Apr-19	Avg.	2016	2017	2018	2019	Avg.
GWNet	MAE	26.94±1.39	26.20±1.56	20.98±0.48	20.82±0.24	23.73±0.75	35.54±3.70	29.94±1.11	26.87±1.57	26.29±1.96	29.66±1.01
	RMSE	41.01±1.71	39.79±1.97	32.44±0.55	31.58±0.13	36.20±0.96	55.84±5.75	45.28±1.33	43.63±2.52	39.29±2.53	46.01±1.46
	MAPE (%)	23.13±0.91	30.31±2.77	23.72±1.09	22.00±0.38	24.79±0.85	41.77±4.84	34.94±0.62	41.09±1.99	41.68±3.65	39.87±1.87
STID	MAE	33.83±0.51	30.83±0.23	26.66±0.46	27.50±0.73	29.71±0.28	40.38±0.56	34.92±0.82	28.36±0.43	29.28±0.16	33.23±0.22
	RMSE	49.98±0.53	45.98±0.25	39.84±0.51	40.70±0.97	44.12±0.37	64.10±0.91	52.12±1.11	46.78±0.78	43.90±0.42	51.72±0.33
	MAPE (%)	31.21±2.21	33.63±0.89	31.11±2.31	30.95±1.25	31.73±0.86	44.78±0.79	42.04±0.82	45.92±0.63	43.87±1.16	44.16±0.60
iTransformer	MAE	30.00±0.11	25.99±0.08	23.27±0.08	22.11±0.06	25.34±0.05	32.03±1.25	27.07±0.35	20.74±0.18	22.28±0.68	25.53±0.56
	RMSE	45.21±0.25	40.76±0.15	35.89±0.10	33.90±0.18	38.94±0.09	50.92±2.20	40.22±0.61	34.16±0.39	33.37±0.65	39.67±0.91
	MAPE (%)	31.23±0.90	29.31±0.59	28.31±0.70	24.50±0.92	28.34±0.20	38.05±2.33	32.25±0.33	33.64±0.43	32.66±0.69	34.15±0.76
TrafficStream	MAE	23.63±0.42	21.87±0.28	19.71±0.37	19.15±0.18	21.09±0.29	30.09 ±0.63	26.43±0.33	20.34±0.39	21.48±0.31	24.58±0.34
	RMSE	36.46±0.53	34.85±0.34	31.07±0.43	29.64±0.20	33.01±0.35	48.04 ±1.10	39.24±0.40	34.22±0.68	32.80±0.42	38.58±0.53
	MAPE (%)	21.34±1.65	20.03±0.36	<u>19.65</u> ±0.81	18.88±0.33	19.98±0.30	34.34±0.33	31.34±0.31	33.17±1.08	30.30±0.38	32.29±0.29
STKEC	MAE	23.22±0.54	22.29±0.35	19.79±0.26	19.09±0.13	21.09±0.13	<u>30.12</u> ±0.30	26.74±0.41	20.34±0.27	21.33±0.31	24.63±0.11
	RMSE	36.16±0.62	35.73±0.52	31.36±0.38	29.71±0.23	33.24±0.13	48.48±0.61	39.57±0.57	34.33±0.70	32.45±0.29	38.70±0.26
	MAPE (%)	<u>20.12</u> ±0.46	<u>20.01</u> ±0.13	19.80±0.79	<u>18.48</u> ±0.14	<u>19.61</u> ±0.19	33.22±0.32	31.87±0.39	32.71±0.48	30.69±0.35	32.12±0.21
PECPM	MAE	24.29±0.22	21.46±0.11	19.34±0.10	19.06±0.15	21.04±0.11	30.86±1.38	26.02±0.60	20.43±0.06	21.15±0.45	24.74±0.25
	RMSE	37.05±0.42	34.11±0.10	30.54±0.25	29.39±0.17	32.77±0.17	49.63±2.19	38.92±0.60	34.66±0.05	32.28±0.49	39.00±0.48
	MAPE (%)	23.05±3.21	20.07±0.43	19.77±1.00	19.07±0.21	20.49±0.91	34.29±0.90	<u>30.69</u> ±0.63	<u>31.84</u> ±0.92	<u>30.22</u> ±0.25	31.93±0.22
STRAP	MAE	30.64±1.85	25.93±0.48	23.87±0.44	24.56±0.66	26.25±0.62	32.14±1.35	26.36±0.78	20.22±0.66	21.91±0.30	25.16±0.32
	RMSE	44.86±2.30	38.96±0.60	35.98±0.70	36.41±0.61	39.05±0.80	51.46±2.36	38.99±0.73	34.07±1.51	32.97±0.10	39.37±0.38
	MAPE (%)	21.72±0.97	22.02±0.63	20.41±0.41	20.43±0.49	21.15±0.47	34.25±0.97	31.62±0.94	32.30±1.13	31.33±0.87	32.37±0.28
EAC	MAE	<u>22.70</u> ±0.82	<u>20.76</u> ±0.73	<u>18.77</u> ±0.66	<u>18.55</u> ±0.57	<u>20.20</u> ±0.69	30.36±1.21	<u>25.74</u> ±0.33	<u>19.74</u> ±0.34	<u>21.02</u> ±0.23	<u>24.21</u> ±0.43
	RMSE	<u>34.85</u> ±1.19	<u>32.70</u> ±0.96	<u>29.04</u> ±0.96	<u>28.12</u> ±0.87	<u>31.18</u> ±0.99	<u>48.33</u> ±1.71	<u>38.22</u> ±0.24	<u>32.88</u> ±0.51	<u>31.89</u> ±0.26	<u>37.83</u> ±0.60
	MAPE (%)	20.56±1.51	20.72±1.94	20.26±1.18	19.14±0.90	20.17±1.25	33.43±0.72	30.78±0.52	32.02±1.31	30.86±0.31	<u>31.77</u> ±0.53
STBP	MAE	17.88 ±0.37	15.86 ±0.10	14.62 ±0.20	14.73 ±0.22	15.77 ±0.09	30.95±0.55	24.32 ±0.09	18.82 ±0.20	20.47 ±0.15	23.64 ±0.23
	RMSE	29.56 ±0.84	26.54 ±0.14	23.58 ±0.30	23.13 ±0.34	25.70 ±0.16	49.34±0.79	37.65 ±0.18	32.14 ±0.24	31.92 ±0.34	37.76 ±0.30
	MAPE (%)	15.47 ±0.32	14.86 ±0.05	14.30 ±0.12	15.11 ±0.14	14.94 ±0.05	31.89 ±0.48	27.98 ±0.28	30.22 ±0.42	28.72 ±0.32	29.70 ±0.35

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (12)$$

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100\% \quad (13)$$

where n represents the number of observed samples, y_i denotes the i -th real sample, and \hat{y}_i is the corresponding predicted value.

A.4.3 IMPLEMENTATION DETAILS

All experiments are conducted on a machine with an NVIDIA Tesla V100 GPU and 32 GB of memory. The Adam optimizer, with an initial learning rate of 0.01, is used to optimize the training process. The batch size is set to 64, the number of training epochs is set to 200, and an early stopping mechanism is implemented to ensure efficient convergence. The reported results for all baselines are the average of five repeated runs.

A.4.4 EXPERIMENT RESULTS

Tables 7 and 8 report detailed results for each incremental period, where the metrics of a given period are averaged over 12 forecasting steps. Figure 9 provides a visual summary of the same results to better illustrate the performance trends across periods. Overall, STBP consistently achieves strong performance throughout the entire continual spatio-temporal forecasting process, including both the aggregated performance across all periods and the period-wise results. This advantage is largely attributed to the well-designed spatio-temporal backbone and the contextual pattern bank, which together support effective knowledge reuse and adaptation under evolving spatio-temporal patterns.

A.4.5 PARAMETER SENSITIVITY ANALYSIS

Beyond the feature dimension d , we further investigated the sensitivity of two key architectural hyperparameters: the number of DLGA layers and attention heads. As shown in Figure 10, increasing either parameter yields marginal gains at best, and in some cases, even leads to slight performance

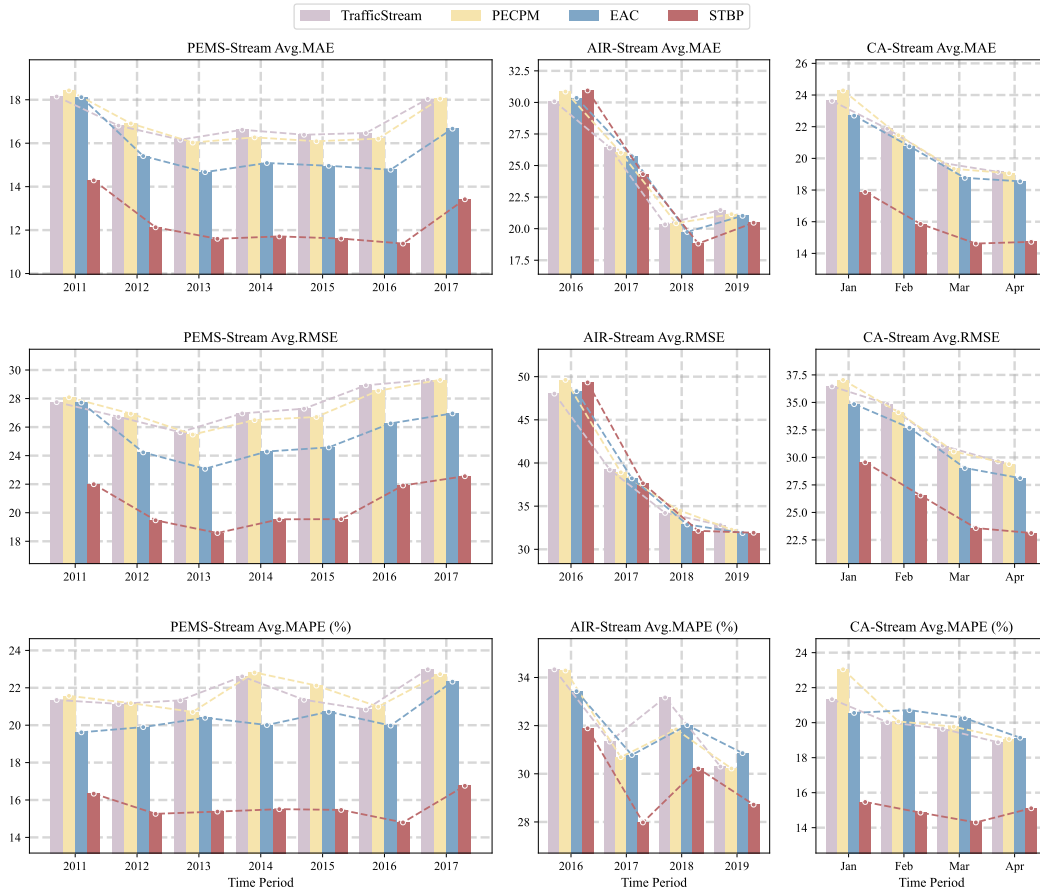


Figure 9: Visualization of period-wise forecasting performance across incremental periods.

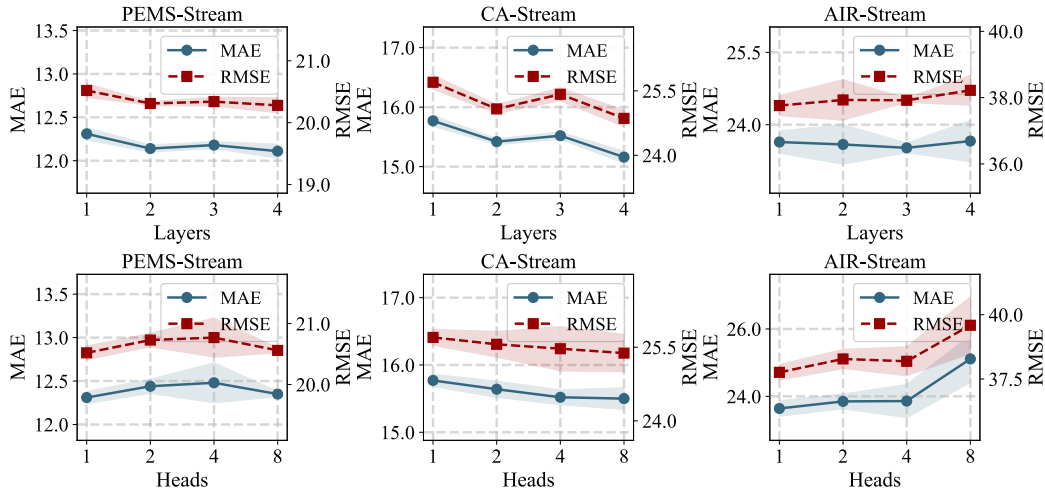


Figure 10: Additional Results of parameter experiments.

degradation. Overall, apart from the feature dimension, model performance remains relatively insensitive to these hyperparameter variations.

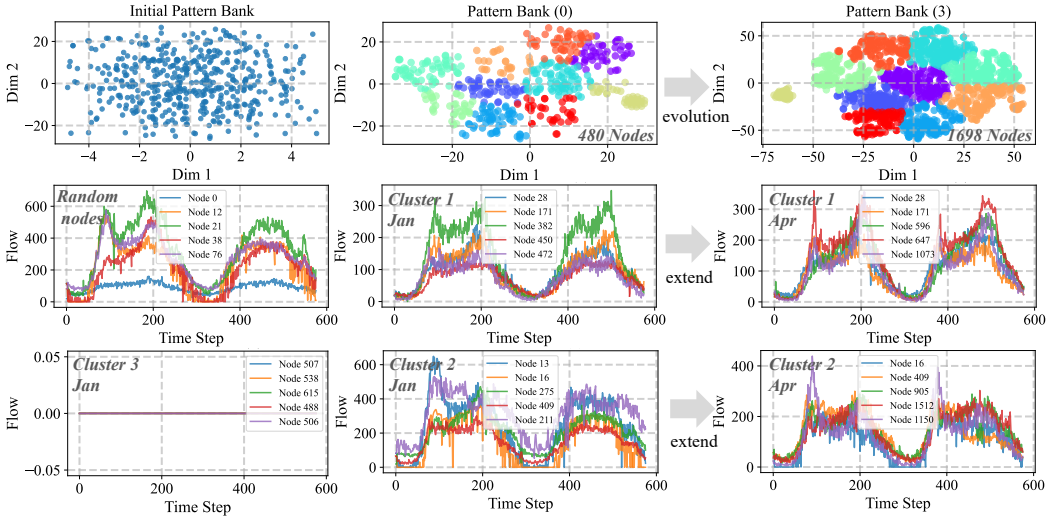


Figure 11: Case Study on CA-Stream.

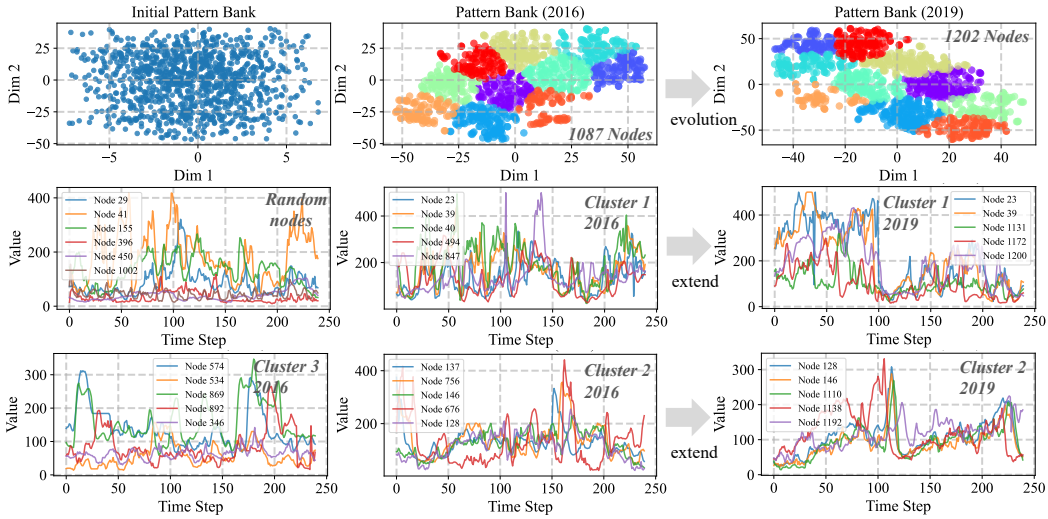


Figure 12: Case Study on AIR-Stream.

A.4.6 ADDITIONAL CASE STUDY

To maintain consistency with the case study on PEMS-Stream, we also conduct case studies on the CA-Stream and AIR-Stream datasets to further validate the expansion and distinction capabilities of the contextual pattern bank in STBP. The experimental results for CA-Stream are shown in Figure 11. Even in the more challenging task of node increment, STBP’s contextual pattern bank effectively distinguishes and consolidates different spatio-temporal patterns, incorporating new patterns introduced by newly added nodes into the existing pattern clusters.

Figure 12 presents the results on AIR-Stream. Compared to traffic flow data, the spatio-temporal patterns in this dataset exhibit more complex periodic and trend changes. Nonetheless, STBP continues to accurately differentiate and consolidate diverse patterns, indicating that its contextual pattern bank has adaptive inductive capabilities for various types of spatio-temporal patterns, independent of the specific dataset type. This mechanism enables STBP to exhibit greater flexibility and adaptability in CSTF tasks.

In addition, Figure 13 provides an intuitive comparison of the predictive performance of STBP and the second-best model, EAC, in real-world application scenarios. These representative cases further substantiate the superior practical utility of STBP in realistic continual learning settings.

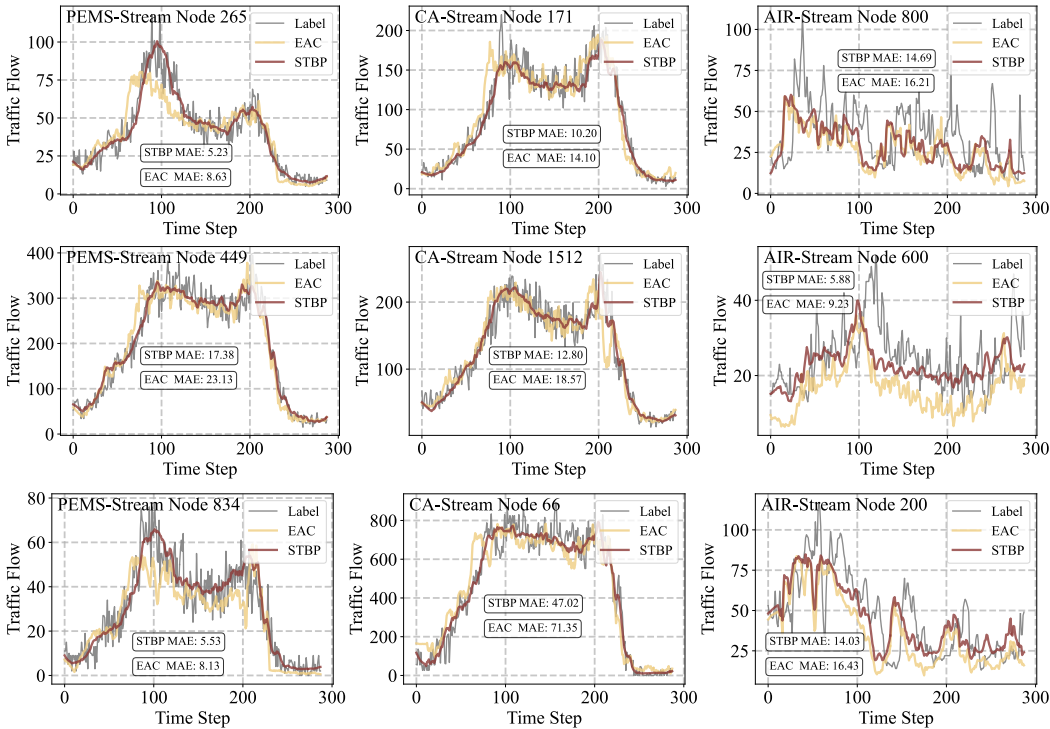


Figure 13: Additional visualization of real forecasting results.

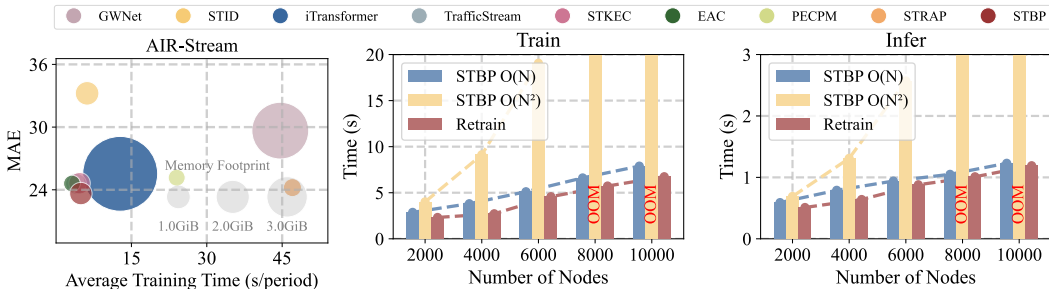


Figure 14: Additional efficiency comparison. STBP $O(N^2)$ denotes the version without linear attention, and Retrain refers to removing the contextual pattern bank.

A.4.7 EFFICIENCY STUDY

Figure 14 provides additional experiments assessing the efficiency and scalability of STBP. Overall, these results confirm that STBP achieves favorable scalability and efficiency, and that its linear-attention design and modular contextual pattern bank structure enable it to handle large-scale spatio-temporal graphs in continual learning settings.

A.5 LIMITATION

Despite STBP’s strong performance on benchmark datasets, its limitations in cross-domain generalization warrant further investigation. Current continual learning approaches, including ours, typically assume incremental tasks originate from similar domains—an idealization that diverges from real-world dynamic and heterogeneous environments. In practice, cross-domain distribution shift introduce dual challenges: feature space misalignment and exacerbated catastrophic forgetting. While STBP’s architecture exhibits inherent adaptability—with DLGA dynamically capturing topological variations and FreNet extracting domain-invariant frequency patterns—its robustness remains unverified under significant structural divergence between source and target domains. Future work should therefore validate the framework’s efficacy in such complex cross-domain scenarios.

A.6 BROADER IMPACT

STBP, with its carefully designed general spatio-temporal backbone and contextual pattern bank expansion mechanism tailored for dynamic scenario changes, effectively achieves continual spatio-temporal forecasting. This approach demonstrates that the spatio-temporal backbone can serve as a stable infrastructure, consistently retaining the ability to model general spatio-temporal dependencies. When facing new or evolving scenarios, there is no need to retrain the backbone. Instead, by introducing scalable parameters relevant to the current scenario, the model can rapidly adapt to new tasks. Building on this concept, we aim to advance the development of a spatio-temporal foundational model with enhanced cross-domain generalization, while concurrently exploring the potential of Large Language Models (LLMs) in spatio-temporal and time-series forecasting tasks (Liu et al., 2025b; 2024a; 2025a).

This involves two key directions: ❶ introducing explicit domain adaptation mechanisms to better distinguish between domain-specific and shared features, and ❷ exploring cross-domain shared contextual pattern banks to enhance adaptability while maintaining efficiency. This approach involves continuously training a unified backbone model with spatio-temporal data from multiple heterogeneous domains, thereby enhancing its spatio-temporal representational capacity. As data from various domains are continuously integrated and trained, the spatio-temporal foundational model will evolve, enabling efficient generalization and adaptation to entirely new scenarios or tasks by incorporating only a small number of additional parameters. Such a model holds the potential to benefit society by improving intelligent transportation through more accurate traffic forecasting and supporting climate resilience via advanced environmental modeling.

A.7 LLM USAGE

In accordance with the ICLR 2026 policy on large language model (LLM) usage, we disclose that we used an LLM (ChatGPT) solely for the purpose of improving the grammar, clarity, and fluency of the manuscript. The content, structure, technical contributions, experiments, analysis, and all scientific writing were entirely conceived, drafted, and validated by the human authors. The LLM was not involved in research ideation, experimental design, data analysis, or any aspect of the scientific content creation. All outputs generated by the LLM were reviewed and edited by the authors to ensure accuracy and correctness. We confirm that no hidden prompts, prompt injections, or LLM-generated falsehoods were introduced in the manuscript, and all use of LLMs complies with the ICLR Code of Ethics.