

# Node Feature Forecasting in Temporal Graphs: an Interpretable Online Algorithm

Anonymous authors  
Paper under double-blind review

## Abstract

1 In this paper, we propose an online algorithm `mSPACE` for forecasting node features in  
 2 temporal graphs, which captures spatial cross-correlation among different nodes as well as the  
 3 temporal auto-correlation within a node. The algorithm can be used for both probabilistic  
 4 and deterministic multi-step forecasting, making it applicable for estimation and generation  
 5 tasks. Comparative evaluations against various baselines, including temporal graph neural  
 6 network (TGNN) models and classical Kalman filters, demonstrate that `mSPACE` performs  
 7 at par with the state-of-the-art and even surpasses them on some datasets. Importantly,  
 8 `mSPACE` demonstrates consistent performance across datasets with varying training sizes, a  
 9 notable advantage over TGNN models that require abundant training samples to effectively  
 10 learn the spatiotemporal trends in the data. Therefore, employing `mSPACE` is advantageous  
 11 in scenarios where the training sample availability is limited. Additionally, we establish  
 12 theoretical bounds on multi-step forecasting error of `mSPACE` and show that it scales linearly  
 13 with the number of forecast steps  $q$  as  $\mathcal{O}(q)$ . For an asymptotically large number of nodes  $n$ ,  
 14 and timesteps  $T$ , the computational complexity of `mSPACE` grows linearly with both  $n$ , and  $T$ ,  
 15 i.e.,  $\mathcal{O}(nT)$ , while its space complexity remains constant  $\mathcal{O}(1)$ . We compare the performance  
 16 of various `mSPACE` variants against ten recent TGNN baselines and two classical baselines,  
 17 ARIMA and the Kalman filter across ten real-world datasets. Lastly, we have investigated the  
 18 interpretability of different `mSPACE` variants by analyzing model parameters alongside dataset  
 19 characteristics to jointly derive model-centric and data-centric insights.

## 20 1 Introduction

21 Temporal graphs are a powerful tool for modelling real-world data that evolves over time. They are increasingly  
 22 being used in diverse fields, such as recommendation systems (Gao et al., 2022), social networks (Deng et al.,  
 23 2019), and transportation systems (Yu et al., 2018), to name a few. Temporal graph learning (TGL) can be  
 24 viewed as the task of learning on a sequence of graphs that form a time series. The changes in the graph can  
 25 be of several types: changes to the number of nodes, the features of existing nodes, the configuration of edges,  
 26 or the features of existing edges. Moreover, a temporal graph can result from a single or a combination of  
 27 these changes. The TGL methods can be applied to various tasks, such as regression, classification, and  
 28 clustering, at three levels: node, edge, and graph (Longa et al., 2023).

29 In this work, we focus on node feature forecasting, also known as node regression, where the previous temporal  
 30 states of a graph are used to predict its future node features. In most temporal graph neural network (TGNN)  
 31 models, the previous states are encoded into a super-state or dynamic graph embedding (Barros et al., 2021),  
 32 guided by the graph structure. This dynamic embedding is then used to forecast the future node features.  
 33 While the TGNN models perform well, they could be more interpretable, as a direct relationship between the  
 34 node features and the embeddings cannot be understood straightforwardly. Furthermore, most embedding  
 35 aggregation mechanisms impose a strong assumption that the neighbours influence a node in proportion to  
 36 their edge weight (Wang et al., 2021).

37 TGNN methods (Li et al., 2018; Micheli & Tortorella, 2022; Wu et al., 2019; Fang et al., 2021; Liu et al.,  
 38 2023) typically involve a training phase where the model learns from training data and is then deployed on

39 test data without further training due to computational costs. If the test data distribution differs from the  
 40 training data, an offline model cannot adapt. Therefore, when dealing with time-series data, it is crucial  
 41 to use a lightweight online algorithm that can adapt to changes in data distribution while also performing  
 42 forecasts. Moreover, TGNN models are typically trained to forecast a predetermined number of future steps.  
 43 If we want to increase the number of forecast steps, even by one, the model needs to be reinitialized and  
 44 retrained.

45 Inspired by the simplicity of Markov models, we define the state of a graph at a given time in an interpretable  
 46 manner and propose a lightweight model that can be deployed without any training. The algorithm is  
 47 designed with a mechanism to prioritize recent trends in the data over historical ones, allowing it to adapt to  
 48 changes in data distribution.

49 **Contributions** The contributions of our work are summarized as follows:

- 50 • We have proposed an online learning algorithm `mspace` for node feature forecasting in temporal  
 51 graphs, which can sequentially predict the node features for  $q \in \mathbb{N}$  future timesteps after observing  
 52 only two past node features.
- 53 • The algorithm `mspace` can produce both probabilistic and deterministic forecasts, making it suitable  
 54 for generative and predictive tasks.
- 55 • The root mean square error (RMSE) of  $q$ -step iterative forecast scales linearly in the number of steps  
 56  $q$ , i.e.  $\text{RMSE}(q) = \mathcal{O}(q)$ .
- 57 • For asymptotically large number of nodes  $n$ , and timesteps  $T$ , the computational complexity of  
 58 `msapce` grows linearly with both  $n$ , and  $T$ , i.e.,  $\mathcal{O}(nT)$ , while the space complexity is constant  $\mathcal{O}(1)$ .
- 59 • We have compared the performance of different variants of `mspace` against ten recent TGNN baselines,  
 60 and two classical baselines ARIMA, and Kalman filter.
- 61 • We have evaluated `mspace` on four datasets for single-step forecasting and six datasets for multi-step  
 62 forecasting.
- 63 • In addition to the evaluation on ten real-world datasets, we have proposed a technique to generate  
 64 synthetic datasets that can aid in a more thorough evaluation of node feature forecasting methods.  
 65 The synthetic datasets have the potential to serve as benchmark for future research.
- 66 • We have investigated the interpretability of different `mspace` variants by analyzing the model  
 67 parameters along with the dataset characteristics to jointly derive model-centric and data-centric  
 68 insights.
- 69 • To facilitate the reproducibility of results, the `code` is made available here.

70 **Notation** We denote vectors with lowercase boldface  $\mathbf{x}$ , and matrices and tensors with uppercase boldface  
 71  $\mathbf{X}$ . Sets are written in calligraphic font such as  $\mathcal{V}, \mathcal{U}, \mathcal{S}, \mathcal{C}$ , with the exception of graphs  $\mathcal{G}$ , and queues  $\mathcal{Q}$ .  
 72 The operator  $\succ$  is used in two contexts:  $\mathbf{x} \succ \mathbf{0}$  is an element wise positivity check on the vector  $\mathbf{x}$ , and  $\mathbf{A} \succ \mathbf{0}$   
 73 indicates that the matrix  $\mathbf{A}$  is positive definite.  $\mathbb{I}(\cdot)$  is the indicator function, and  $[m] \triangleq \{1, 2, \dots, m\}$  for any  
 74  $m \in \mathbb{N}$ . We denote the distributions of continuous variables by  $p(\cdot)$ , and of discrete variables by  $P(\cdot)$ . The  
 75 statement  $\mathbf{x} \sim p$  means that  $\mathbf{x}$  is sampled from  $p$ . The Hadamard product operator is denoted by  $\odot$  while  
 76 the Kronecker product operator is denoted by  $\otimes$ . The trace of a matrix  $\mathbf{A}$  is written as  $\text{tr}(\mathbf{A})$ .

77 We denote the neighbours of a node  $v$  for an arbitrary number of hops as  $\mathcal{U}_v$ . The neighbours of node  $v$  up  
 78 to  $K$  number of hops is defined as follows. Let  $\mathbf{N} = \sum_{k \in [K]} \mathbf{A}^k$ , then  $\mathcal{U}_v = \{u : \mathbf{N}_{v,u} > 0, \forall u \in \mathcal{V}\}$ . Since  
 79  $\mathbf{A}_{v,v} = \mathbf{1}$ ,  $v \in \mathcal{U}_v$ . We introduce the operator  $\langle \cdot \rangle$  to arrange the nodes in a set  $\mathcal{U}$  in ascending numerical  
 80 order of the node IDs. When another set or vector is super-scripted with  $\langle \mathcal{U} \rangle$ , the elements within that set or  
 81 vector are filtered and arranged as per  $\langle \mathcal{U} \rangle$ .

82 A Markov chain is represented using  $\mathfrak{J}$  with different subscripts for identification. The transition kernel of a  
 83 Markov chain is denoted as  $\mathbf{P}$  with  $\mathbf{P}_{\mathbf{a},\mathbf{b}}$  representing the probability of transitioning from state  $\mathbf{a}$  to  $\mathbf{b}$ .

84 **Organization** In Sec. 2 we formulate the problem of node feature forecasting and also a propose a model  
 85 to solve it. In Sec. 3 we expand upon the solution and present it as an algorithm. We discuss the related  
 86 works in Sec. 4 and present the results on single-step and multi-step node feature forecasting in Sec. 5. In  
 87 Sec. 6 we discuss the interpretability of the proposed algorithm and then discuss the limitations in Sec. 7.  
 88 Finally, we conclude in Sec. 8.

## 89 2 Methodology

90 **Problem Formulation** A discrete-time temporal graph is defined as  $\{\mathcal{G}_t = (\mathcal{V}, \mathcal{E}, \mathbf{X}_t) : t \in [T]\}$ , where  
 91  $\mathcal{V} = [n]$  is the set of nodes,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  is the set of edges, and  $\mathbf{X}_t \in \mathbb{R}^{n \times d}$  is the node feature matrix at time  
 92  $t$ . The set of edges  $\mathcal{E}$  can alternatively be represented by the adjacency matrix denoted as  $\mathbf{A} \in \{0, 1\}^{n \times n}$ .  
 93 The node feature vector is denoted by  $\mathbf{x}_t(v) \in \mathbb{R}^d$  such that  $\mathbf{X}_t = [\mathbf{x}_t(v)]_{v \in \mathcal{V}}^\top$ , and we refer to the first-order  
 94 differencing (Shumway & Stoffer, 2017) of a node feature vector as **shock**. For a node  $v \in \mathcal{V}$  we define the  
 95 shock at time  $t$  as  $\varepsilon_t(v) \triangleq \mathbf{x}_t(v) - \mathbf{x}_{t-1}(v)$ . The shock of the nodes in an ordered set  $\mathcal{U}$  at time  $t$  is denoted  
 96 by  $\varepsilon_t^{(\mathcal{U})} \in \mathbb{R}^{|\mathcal{U}|d}$ . The shock at time  $t$  for an arbitrary set of nodes is  $\varepsilon_t$ .

97 **Assumption 2.1.** The shocks  $\{\varepsilon_1, \varepsilon_2, \varepsilon_3 \dots \varepsilon_T\}$  is assumed to be sampled from a continuous-state Markov  
 98 chain defined on  $\mathbb{R}^{nd}$  such that  $p(\varepsilon_{t+1} | \varepsilon_t, \varepsilon_{t-1}, \dots) = p(\varepsilon_{t+1} | \varepsilon_t)$ .

99 This is a weak assumption because a continuous-state Markov chain has infinite number of states. However,  
 100 having infinite number of states makes it impossible to learn the transition kernel from limited samples  
 101 without additional assumptions on the model. To circumvent this, *linear dynamical systems* and *autoregressive*  
 102 *models* are used in the literature (Barber, 2012).

103 Let  $p(\varepsilon' | \varepsilon)$  denote the transition probability  $\varepsilon \rightarrow \varepsilon'$  in a continuous-state Markov chain  $\mathfrak{Z}_0$  defined over  $\mathcal{C}$ .  
 104 A discrete-state Markov chain  $\mathfrak{Z}_1$  defined over finite  $\mathcal{S}$  with transition probability  $\mathbf{P}_{s,s'}$  can be constructed  
 105 from  $p(\varepsilon' | \varepsilon)$  through a mapping  $\Psi : \mathcal{C} \rightarrow \mathcal{S}$  as

$$\mathbf{P}_{s,s'} = \frac{\int_{\mathcal{C}} \int_{\mathcal{C}} p(\varepsilon' | \varepsilon) p(\varepsilon) \mathbb{I}(\Psi(\varepsilon) = s) \mathbb{I}(\Psi(\varepsilon') = s') d\varepsilon d\varepsilon'}{\int_{\mathcal{C}} \int_{\mathcal{C}} p(\varepsilon' | \varepsilon) p(\varepsilon) \mathbb{I}(\Psi(\varepsilon) = s) d\varepsilon d\varepsilon'}. \quad (1)$$

106 For a continuous-state Markov chain sample  $\{\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T\}$ , we can estimate  $\mathbf{P}$  directly from  
 107  $\{\Psi(\varepsilon_1), \Psi(\varepsilon_2), \dots, \Psi(\varepsilon_T)\}$  without the need of  $p(\varepsilon' | \varepsilon)$ . Now, consider a *random function*  $\Omega : \mathcal{S} \rightarrow \mathcal{C}$ ,  
 108 such that: (a)  $\Psi(\varepsilon) = s$ , (b)  $\Psi(\varepsilon') = s'$ , (c)  $\varepsilon' = \Omega(s)$ , from which follows  $p(\Omega(s)) = p(\varepsilon' | s)$ .

109 The approximate transition kernel  $\hat{\mathbf{P}}$  due to  $(\Psi, \Omega)$  can be written as:

$$\hat{\mathbf{P}}_{s,s'} = \int_{\{\varepsilon' \in \mathcal{C} : \Psi(\varepsilon') = s'\}} p(\varepsilon' | s) d\varepsilon' = \int_{\mathcal{C}} p(\Omega(s)) \mathbb{I}(\Psi(\varepsilon') = s') d\varepsilon'. \quad (2)$$

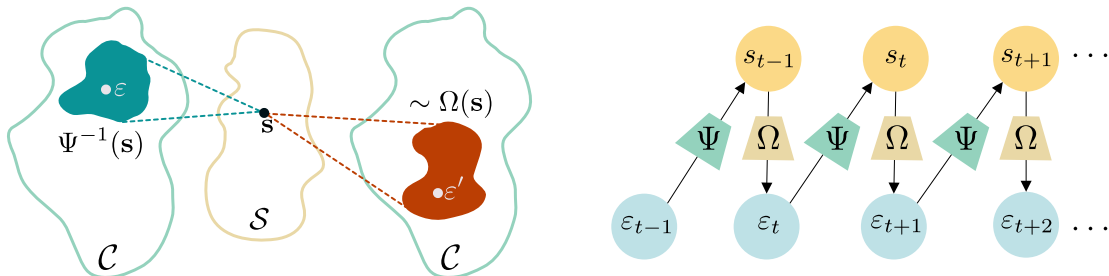


Figure 1: (left) state and sampling functions visualized, (right) Markov approximation.

In Fig. 1 (left) we depict the functions  $\Psi$  mapping from continuous space in  $\mathcal{C}$  to a discrete space  $\mathcal{S}$ . We also depict  $\Omega$  mapping from  $\mathcal{S}$  to  $\mathcal{C}$ . In a red patch we show the range of  $\Omega(\mathbf{s})$ , and in the green patch we show the domain of  $\Psi(\mathbf{s})$ . In Fig. 1 (right), we visualize Assumption 2.1 wherein the shocks evolve as a Markov chain through the functions  $\Psi, \Omega$ .

We refer to  $\Psi$  as the **state function**, and  $\Omega$  as the **sampling function**. The approximated Markov chain defined over  $\mathcal{S}$  resulting from  $(\Psi, \Omega)$  is denoted as  $\hat{\mathfrak{Z}}(\Psi, \Omega)$ , with  $p(\hat{\varepsilon}' | \varepsilon) = p(\Omega \circ \Psi(\varepsilon))$ . Ideally, the goal is to find the pair of functions  $(\Psi, \Omega)$  such that: (a)  $\mathbf{P}_{\mathbf{s}, \mathbf{s}'} = \hat{\mathbf{P}}_{\mathbf{s}, \mathbf{s}'}, \forall \mathbf{s}, \mathbf{s}' \in \mathcal{S}$ , (b)  $p(\varepsilon' | \varepsilon) = p(\Omega \circ \Psi(\varepsilon)) \forall \varepsilon \in \mathcal{C}$ . However, in practice this is quite ambitious as the state and sampling functions will induce some error in the encoding and decoding process. Therefore, we frame the problem as follows.

The sequence of shocks drawn from the original Markov chain  $\mathfrak{Z}_0$  is represented as  $\{\varepsilon_t : t \in [T]\} \sim \mathfrak{Z}_0$ . Then, for each  $\varepsilon_t$  we generate a sequence of  $q$  future shocks using the Markov chain  $\hat{\mathfrak{Z}}(\Psi, \Omega)$  as

$$\hat{\varepsilon}_{t+j} = (\Omega \circ \Psi)^j(\varepsilon_t), \quad \forall t \in [T-q], j \in [q].$$

The problem is to design  $\Psi, \Omega$  such that  $\left\| \sum_{j \in [k]} \varepsilon_{t+j} - \hat{\varepsilon}_{t+j} \right\|^2$  is minimized  $\forall k \in [q], t \in [T-q]$ , which can be written alternatively as:

**Problem 2.1** ( $q$ -step node feature forecasting). Design the state and sampling functions  $\Psi, \Omega$  such that

$$\min \sum_{t \in [T-q]} \sum_{k \in [q]} \left\| \sum_{j \in [k]} \varepsilon_{t+j} - (\Omega \circ \Psi)^j \varepsilon_t \right\|^2. \quad (3)$$

In a deep learning context, both  $\Psi$  and  $\Omega$  would typically be neural networks trained directly using the objective in Problem 2.1. In this work, however, we explicitly define  $\Psi$  and  $\Omega$  and learn their parameters through the same objective.

**Proposed Model** Instead of creating a single model to approximate  $p(\varepsilon_{t+1}^{(\mathcal{V})} | \varepsilon_t^{(\mathcal{V})})$ , we create a model for each node  $v \in \mathcal{V}$  to approximate  $p(\varepsilon_{t+1}^{(\mathcal{U}_v)} | \varepsilon_t^{(\mathcal{U}_v)})$  where  $\mathcal{U}_v$  denotes the neighbours of node  $v$  within a certain number of hops. We present this in the following assumption.

**Assumption 2.2.** The shock of node  $v$  at time  $t+1$  can be estimated from the shock of its neighbouring nodes in  $\mathcal{U}_v$  at time  $t$ .

While  $\varepsilon_t(u')$  for any node  $u' \notin \mathcal{U}_v$  may help in estimating  $\varepsilon_{t+1}(v)$ , we assume that enough information is already conveyed by the nodes in  $\mathcal{U}_v$  that the impact of considering node  $u'$  would be minimal. It must be noted that  $\mathcal{U}_v$  denotes the neighbours of node  $v$  up to an arbitrary number of hops, therefore if we consider  $\mathcal{U}_v$  to mean  $k$  hops, then all the nodes that neighbours  $v$  with  $1, 2, \dots, k$  hops are all in  $\mathcal{U}_v$  and their impact is considered. Assumption 2.2 is important to create a scalable model, because in a connected graph every node will be correlated with every other node which will make the state space prohibitively large.

We propose two variants of the **state function**, one which captures the characteristics of the shock  $\Psi_{\mathcal{S}}$ , and the other which is concerned with the timestamps  $\Psi_{\mathcal{T}}$  and captures seasonality.

- $\Psi_{\mathcal{S}} : \mathbb{R}^{|\mathcal{U}|d} \rightarrow \{-1, 1\}^{|\mathcal{U}|d}, \quad \Psi_{\mathcal{S}}(\varepsilon^{(\mathcal{U})}) = \text{sign}(\varepsilon^{(\mathcal{U})})$ .
- $\Psi_{\mathcal{T}} : \mathbb{N} \rightarrow \{0, 1, \dots, \tau_0 - 1\}, \quad \Psi_{\mathcal{T}}(t) = t \bmod \tau_0$ , where  $\tau_0 \in \mathbb{N}$  is the time period.

We also define two variants of the **sampling function**:

- deterministic  $\Omega_{\mu}(\mathbf{s}) = \boldsymbol{\mu}(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}$ .
- probabilistic  $\Omega_{\mathcal{N}}(\mathbf{s}) \sim \mathcal{N}(\varepsilon'; \boldsymbol{\mu}(\mathbf{s}), \boldsymbol{\Sigma}(\mathbf{s})), \forall \mathbf{s} \in \mathcal{S}$ .

More details on the state functions are provided in Sec. 6, where we offer a comprehensive explanation. The proposed model is presented as an online algorithm and discussed in detail in the following section.

### 148 3 Algorithm

149 We name our algorithm **mspace** with a suffix specifying the state and sampling functions. For example,  
 150 **mspace-S $\mathcal{N}$**  represents the algorithm with state function  $\Psi_{\mathbf{s}}$ , and sampling function  $\Omega_{\mathcal{N}}$ . For each node  $v \in \mathcal{V}$ ,  
 151 we approximate  $p(\varepsilon_{t+1}^{(U_v)} \mid \Psi_{\mathbf{s}}(\varepsilon_t^{(U_v)}) = \mathbf{s})$  as a Gaussian distribution with mean vector  $\boldsymbol{\mu}_v(\mathbf{s}) \in \mathbb{R}^{|\mathcal{U}_v|d}$  and  
 152 covariance matrix  $\boldsymbol{\Sigma}_v(\mathbf{s}) \in \mathbb{R}^{|\mathcal{U}_v|d \times |\mathcal{U}_v|d}$  indexed by the state  $\mathbf{s} \in \{-1, 1\}^{|\mathcal{U}_v|d}$ . The parameters  $\boldsymbol{\mu}_v(\mathbf{s}), \boldsymbol{\Sigma}_v(\mathbf{s})$   
 153 are learnt through maximum likelihood estimation (MLE). For each node  $v \in \mathcal{V}$ , and state  $\mathbf{s}$  we maintain a  
 154 **queue**  $\mathcal{Q}_v(\mathbf{s})$  of maximum size  $M$  in which the shocks succeeding a given state  $\mathbf{s}$  are collected. The MLE  
 155 solution is calculated as  $\boldsymbol{\mu}_v(\mathbf{s}) \leftarrow \text{mean}(\mathcal{Q}_v(\mathbf{s}))$ , and  $\boldsymbol{\Sigma}_v(\mathbf{s}) \leftarrow \text{covariance}(\mathcal{Q}_v(\mathbf{s}))$ .

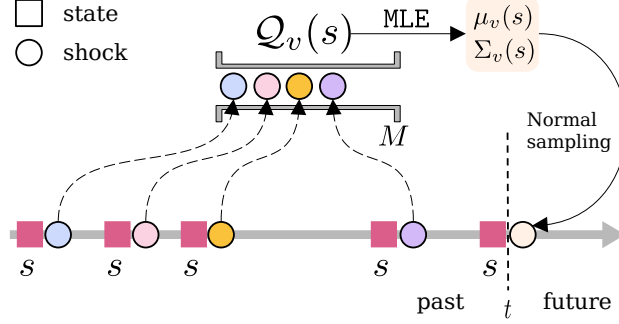


Figure 2: Operation of a queue.

156 The use of a fixed-size queue ensures that the model prioritises recent data over historical data, thereby  
 157 allowing the system to adapt to prevailing trends. It must be noted that obtaining the parameters  $\boldsymbol{\mu}_v(\mathbf{s}), \boldsymbol{\Sigma}_v(\mathbf{s})$   
 158 from historical data *relaxes the Markov assumption* in the original model. The queue  $\mathcal{Q}_v(\mathbf{s}) = \{\varepsilon_{\tau}^{(U_v)} : \Psi_{\mathbf{s}}(\varepsilon_{\tau-1}^{(U_v)}) = \mathbf{s}, \tau < t\}$   
 159 contains shocks from the past (see Fig. 2). Therefore, the estimated sample  $\varepsilon_{t+1}^{(U_v)}$   
 160 depends on certain shocks from the past which violates the Markov property.

161 As **mspace** is an online algorithm, we might encounter unobserved states for which the queue is empty, and  
 162 therefore cannot employ MLE. To facilitate *inductive inference*, as a state  $\mathbf{s}_t$  is encountered, we find the state  
 163  $\mathbf{s}^* \in \mathcal{S}_v$  which is the closest to  $\mathbf{s}_t$ , i.e.,  $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \|\mathbf{s} - \mathbf{s}_t\|$ , where  $\mathcal{S}_v$  is the set of states observed  
 164 before time  $t$ .

---

#### Algorithm 1 **mspace-S $\mathcal{N}$**

---

**Input**  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ ,  $r \in [0, 1)$ ,  $q, M$

**Output**  $\hat{\varepsilon}_t(v), \forall v \in \mathcal{V}, t \in [[r \cdot T], T]$

1:  $\varepsilon_t \leftarrow \mathbf{x}_t - \mathbf{x}_{t-1}, \forall t \in [T]$

*Offline training (A)*

2: **for**  $t \in [[r \cdot T]]$  **do**

3:   **for**  $v \in \mathcal{V}$  **do**

4:      $\mathbf{s}_t \leftarrow \Psi_{\mathbf{s}}(\varepsilon_t^{(U_v)})$

5:      $\mathcal{S}_v \leftarrow \mathcal{S}_v \cup \{\mathbf{s}_t\}$

6:      $\mathcal{Q}_v(\mathbf{s}_t) \leftarrow \text{enqueue } \varepsilon_{t+1}^{(U_v)}$

7:   **end for**

8: **end for**

9: **for**  $v \in \mathcal{V}$  **do**

10:    $\boldsymbol{\mu}_v(\mathbf{s}) \leftarrow \text{mean}(\mathcal{Q}_v(\mathbf{s})), \forall \mathbf{s} \in \mathcal{S}_v$

11:    $\boldsymbol{\Sigma}_v(\mathbf{s}) \leftarrow \text{covariance}(\mathcal{Q}_v(\mathbf{s})), \forall \mathbf{s} \in \mathcal{S}_v$

12: **end for**

*Online learning (B)*

13: **for**  $t \in [[r \cdot T], T - q]$  **do**

14:   **for**  $v \in \mathcal{V}$  **do**

15:      $\mathbf{s}_t \leftarrow \Psi_{\mathbf{s}}(\varepsilon_t^{(U_v)})$

16:      $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \|\mathbf{s} - \mathbf{s}_t\|$

17:      $\hat{\varepsilon}_{t+1}^{(U_v)} \sim \mathcal{N}(\varepsilon; \boldsymbol{\mu}_v(\mathbf{s}^*), \boldsymbol{\Sigma}_v(\mathbf{s}^*))$

18:     **for**  $k \in [q] \setminus \{1\}$  **do**

19:        $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \|\mathbf{s} - \Psi(\hat{\varepsilon}_{t+k-1}^{(U_v)})\|$

20:        $\hat{\varepsilon}_{t+k}^{(U_v)} \sim \mathcal{N}(\varepsilon; \boldsymbol{\mu}_v(\mathbf{s}^*), \boldsymbol{\Sigma}_v(\mathbf{s}^*))$

21:     **end for**

22:      $\hat{\varepsilon}_{t+k}(v) \leftarrow \hat{\varepsilon}_{t+k}^{(U_v)}(v), \forall k \in [q]$

23:     Update  $\mathcal{S}_v, \mathcal{Q}_v; \boldsymbol{\mu}_v(\mathbf{s}), \boldsymbol{\Sigma}_v(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}_v$

24:   **end for**

25: **end for**

---

165 **Example** For the purpose of explaining `mspace-SN` we  
 166 consider an example with two nodes  $n = 2$ , and feature  
 167 dimension  $d = 1$ . In Fig. 3 we first show the shock vector  
 168  $\varepsilon_t \in \mathbb{R}^2$ . The state of shock  $\varepsilon_t$ , denoted by  $\Psi(\varepsilon_t)$  is marked  
 169 in  $\mathcal{S} \in \{-1, 1\}^2$ . Corresponding to this state, we have a  
 170 Gaussian distribution  $\mathcal{N}(\varepsilon; \mu(\Psi(\varepsilon_t)), \Sigma(\Psi(\varepsilon_t)))$  depicted  
 171 as an ellipse. The next shock  $\varepsilon_{t+1}$  is sampled from this  
 172 distribution. This distribution is updated as we gather  
 173 more information over time. The volume of the Gaussian  
 174 density in a quadrant is equal to the probability of the  
 175 next shock’s state being in that quadrant, i.e., the tran-  
 176 sition kernel  $\hat{\mathbf{P}}_{s,s'} = \int_{s' \odot \varepsilon > \mathbf{0}} \mathcal{N}(\varepsilon; \mu(s), \Sigma(s)) d\varepsilon$ . **There-**  
 177 **fore, `mspace-SN` can be viewed as a Markov chain whose**  
 178 **transition function is a multivariate Gaussian.**

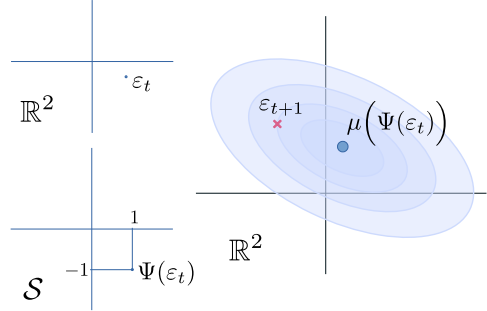


Figure 3: Shock Distribution.

179 **4 Related Works**

180 **Correlated Time Series Forecasting** A set of  $n$  time series data denoted as  $\mathbf{x}_t(v), \forall v \in [n], t \in [T]$  is  
 181 assumed to exhibit spatio-temporal correlation (Wu et al., 2021a; Lai et al., 2023). The correlations can then  
 182 be discerned from the observations to perform forecasting. The correlated time series (CTS) data can be  
 183 viewed as a temporal graph  $\mathcal{G} = (\mathbf{X}_t, \mathbf{A})$ , with  $\mathbf{X}_t \triangleq [\mathbf{x}_t(v)]_{v \in [n]}$  where the *spatial correlation* between  $\mathbf{x}_t(u)$   
 184 and  $\mathbf{x}_t(u)$  is quantified as the edge weight  $\mathbf{A}_{u,v}$ , and  $\mathbf{A}_{u,u}$  signifies the *temporal correlation* within  $\mathbf{x}_t(u)$ .  
 185 The architecture of existing CTS forecasting methods consist of spatial (S) and temporal (T) operators. The  
 186 S-operator can be a graph convolutional network (GCN) (Kipf & Welling, 2017) or a Transformer (Vaswani  
 187 et al., 2017). As for the T-operator, convolutional neural network (CNN), recurrent neural network (RNN)  
 188 (Chung et al., 2014) or Transformer (Zeng et al., 2023) can be used.

189 **Temporal Graph Neural Network** A Graph Neural Network (GNN) is a type of neural network that  
 190 operates on graph-structured data, such as social networks, citation networks, and molecular graphs. GNNs  
 191 aim to learn node and graph representations by aggregating and transforming information from neighbouring  
 192 nodes and edges (Wu et al., 2021b). GNNs have shown promising results in various applications, such as  
 193 node classification, link prediction, and graph classification.

194 **Temporal GNN (TGNN)** (Longa et al., 2023) is an extension of GNNs which operates on temporal graphs  
 195  $\mathcal{G}_t = (\mathbf{X}_t, \mathbf{A}_t)$  where  $\mathbf{X}_t$  denotes the node features, and  $\mathbf{A}_t$  is the evolving adjacency matrix. The TGNN  
 196 architecture can be viewed as a neural network encoder-decoder pair  $(f_\theta, g_\phi)$  (see Fig. 4).

197 A sequence of  $m$  past graph snapshots is first encoded into  
 198 an embedding  $\mathbf{h}_t = f_\theta(\{\mathcal{G}_{t-m+1}, \dots, \mathcal{G}_t\})$ , and then a se-  
 199 quence of  $q$  future graph snapshots is estimated by the de-  
 200 coder as  $\{\hat{\mathcal{G}}_{t+1}, \dots, \hat{\mathcal{G}}_{t+q}\} = g_\phi(\mathbf{h}_t)$ . The parameters  $(\theta, \phi)$  are  
 201 trained to minimize the difference between the true sequence  
 202  $\{\mathcal{G}_{t+1}, \dots, \mathcal{G}_{t+q}\}$  and the predicted sequence  $\{\hat{\mathcal{G}}_{t+1}, \dots, \hat{\mathcal{G}}_{t+q}\}$ .  
 203 In node feature forecasting, the objective is to minimize the  
 204 difference between the node feature matrices  $\{\hat{\mathbf{X}}_{t+1}, \dots, \hat{\mathbf{X}}_{t+q}\}$   
 205 and  $\{\mathbf{X}_{t+1}, \dots, \mathbf{X}_{t+q}\}$ , while in temporal link prediction, the  
 206 goal is to minimize the difference between the graph structures  $\{\hat{\mathbf{A}}_{t+1}, \dots, \hat{\mathbf{A}}_{t+q}\}$  and  $\{\mathbf{A}_{t+1}, \dots, \mathbf{A}_{t+q}\}$ .

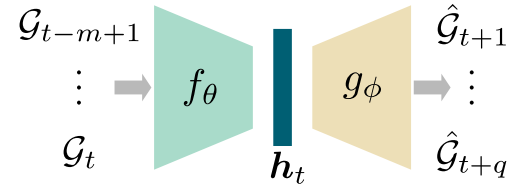


Figure 4: TGNN architecture.

207 There are two main approaches to implementing TGNNs: model evolution and embedding evolution. In  
 208 *model evolution*, the parameters of a static GNN are updated over time to capture the temporal dynamics of  
 209 the graph, e.g., `EvolveGCN` (Pareja et al., 2020). In *embedding evolution*, the GNN parameters remain fixed,  
 210 and the node and edge embeddings are updated over time to learn the evolving graph structure and node  
 211 features (Li et al., 2018; Zhao et al., 2019; Micheli & Tortorella, 2022; Wu et al., 2019; Fang et al., 2021; Liu  
 212 et al., 2023). The TGNN methods are described in Appendix D.3.



251 We notice that `mSPACE-S $\mu$`  achieves a balanced performance between TGNN models and `Kalman-x` across all  
 252 datasets except for `cpox`. The subpar performance of `mSPACE-S*` on the `cpox` dataset may be attributed to  
 253 the seasonal trend, given that it represents the weekly count of chickenpox cases.

254 **Multi-step Forecasting** For the TGNN models, we use the 6 : 2 : 2 train-validation-test chronological split  
 255 in line with the experiments reported by the baselines. For `mSPACE` and `Kalman`, the train-test chronological  
 256 split is 8 : 2, as they do not require a validation set. In Table 3 we report the multi-step  $q = 12^1$  forecasting  
 257 RMSE, and mean absolute error (MAE) on the test set. For `mSPACE`, the queue size  $M = 20^2$ .

Table 3: Multi-step forecasting RMSE and MAE, ( $M = 20$ ).

	PEMS03		PEMS04		PEMS07		PEMS08		PEMSBAY		METRLA	
	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE
GRAM-ODE	26.40	15.72	31.05	19.55	34.42	21.75	25.17	16.05	<b>3.34</b>	<b>1.67</b>	<b>6.64</b>	<b>3.44</b>
STGODE	27.84	16.50	32.82	20.84	37.54	22.99	25.97	16.81	4.89	2.30	7.37	3.75
DCRNN	30.31	18.18	38.12	24.70	38.58	25.30	27.83	17.86	4.74	2.07	7.60	3.60
ARIMA	47.59	33.51	48.80	33.73	59.27	38.17	44.32	31.09	6.50	3.38	13.23	6.90
GWNet	32.94	19.85	39.70	25.45	42.78	26.85	31.05	19.13	4.85	1.95	7.81	3.53
LightCTS	-	-	30.14	18.79	-	-	23.49	14.63	4.32	<u>1.89</u>	<u>7.21</u>	<b>3.42</b>
FOGS	<b>24.09</b>	<b>15.06</b>	31.33	19.35	<b>33.96</b>	<b>20.62</b>	24.09	14.92	-	-	-	-
<code>mSPACE-S<math>\mu</math></code>	36.51	26.43	<u>18.85</u>	<u>13.25</u>	54.39	38.83	<u>14.61</u>	<u>10.36</u>	4.27	2.47	10.24	6.56
<code>mSPACE-T<math>\mu</math></code>	26.53	18.31	<b>13.49</b>	<b>8.70</b>	38.63	24.02	<b>10.35</b>	<b>6.33</b>	<u>3.77</u>	2.19	10.08	6.77
<code>Kalman-x</code>	45.38	33.21	33.75	15.26	64.95	48.01	27.40	12.40	5.71	3.87	13.97	10.7
<code>Kalman-<math>\epsilon</math></code>	749	619	818	709	2313	1988	460	399	50.2	43.1	127.1	109

258 Figure 5 shows the RMSE of the models,  
 259 normalized to the minimum RMSE for the  
 260 dataset, plotted against the number of avail-  
 261 able training samples. We observe that  
 262 `mSPACE-T $\mu$`  performs competitively across  
 263 all datasets with the exception of METRLA.  
 264 Moreover, `mSPACE-T $\mu$`  demonstrates super-  
 265 ior performance compared to `mSPACE-S $\mu$`   
 266 across all the datasets which suggests that  
 267 temporal auto-correlation dominate spatial  
 268 cross-correlation among the nodes.

269 TGNN models, being neural networks, rely  
 270 heavily on the amount of training data avail-  
 271 able. With the relatively small number  
 272 of training samples in PEMS04 and PEMS08,  
 273 these models underperform. In contrast,  
 274 both variants of `mSPACE` significantly surpass  
 275 the state-of-the-art (SoTA), demonstrating their effectiveness  
 276 with smaller datasets <sup>3</sup>. Furthermore, `mSPACE-T $\mu$`  ranks as the second-best model for the largest dataset,  
 277 PEMS08. Therefore, we conclude that `mSPACE` offers consistent performance across datasets with varying  
 sample sizes, and it is particularly advantageous when training data is limited.

278 In Fig. 6, we illustrate how the RMSE scales with the number of forecast steps  $q$  for different variants  
 279 of `mSPACE`. The scaling law for `mSPACE-S*` appears linear, while for `mSPACE-T*`, it appears sublinear. We  
 280 investigate this theoretically in Appendix A.

281 The TGNN baselines perform forecasting for  $q = 12$  future steps, relying on the node features from the  
 282 preceding 12 time steps as input. In contrast, `mSPACE` requires only the node features from the **two previous**

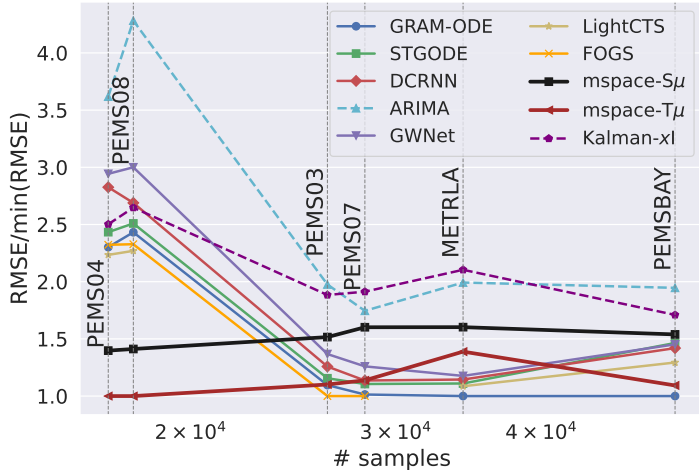


Figure 5: Multi-step forecasting normalised RMSE.

<sup>1</sup> $q = 12$  corresponds to one hour in the traffic datasets used.

<sup>2</sup>a higher value of  $M$  might give better estimates at the cost of higher memory usage and lower adaptability.

<sup>3</sup>single-step forecasting datasets have prohibitively low number of samples ( $< 800$ ), likely limiting `mSPACE`'s performance compared to multi-step forecasting with  $17k+$  samples.



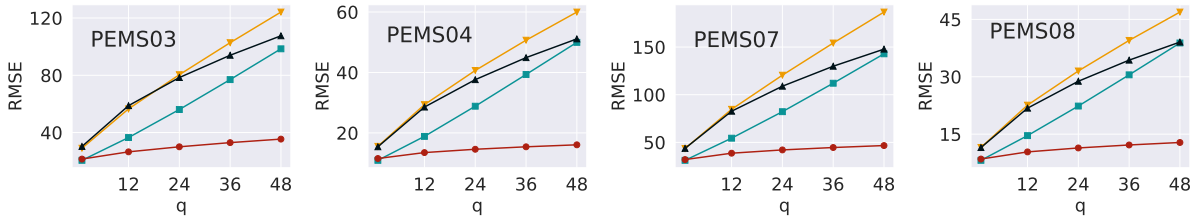


Figure 6: Scaling of error with the number of forecast steps  $q$  using different `mspace` variants:  $\blacktriangledown$  `mspace-S $\mathcal{N}$` ,  $\blacktriangle$  `mspace-T $\mathcal{N}$` ,  $\blacksquare$  `mspace-S $\mu$` ,  $\bullet$  `mspace-T $\mu$` .

time steps. Additionally, `mspace` has the flexibility to forecast for any  $q \in \mathbb{N}$ , whereas TGNN models are limited to forecasting up to the specified number of steps they were trained on. Moreover, `mspace` offers both probabilistic ( $\Omega_{\mathcal{N}}$ ) and deterministic ( $\Omega_{\mu}$ ) forecasts, a capability absent in the baselines. Finally, while TGNN baselines exploit the edge weights information for predictions, `mspace` achieves comparable results using only the graph structure.

## 6 Interpretability

In this section, we examine `mspace` in light of the following definition of Interpretability.

**Definition 6.1.** Consider data  $\mathbf{x} \in \mathcal{D}$  which is processed by a model  $F_{\theta}$  to produce the output  $\hat{\mathbf{y}} \in \mathcal{Y}$ , i.e.,  $\hat{\mathbf{y}} = F_{\theta}(\mathbf{x})$ , where  $\theta$  denotes the model parameters. Moreover, consider a true mapping  $f : \mathbf{x} \mapsto \mathbf{y}, \forall \mathbf{x} \in \mathcal{D}$  where  $\mathbf{y}$  is the ground truth associated with the input data  $\mathbf{x}$ . Then, an interpretable or explainable model  $F_{\theta}$  fulfils one or more of the following properties (Gilpin et al., 2018; Du et al., 2019):

- The internals of the model  $F_{\theta}$  can be explained in a way that is understandable to humans.
- The output  $\hat{\mathbf{y}}$  can be explained in terms of the properties of the input  $\mathbf{x}$ , the input data distribution  $\mathcal{D}$ , and the model parameters  $\theta$ .
- The failure of a model on a given input data can be explained.
- For a certain distance metric  $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ , theoretical bounds on the expected error  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\Delta(\mathbf{y}, F_{\theta}(\mathbf{x}))]$  can be established based on the description of  $F_{\theta}$ , supported by the assumptions on the input data distribution  $\mathcal{D}$ .
- It can be identified whether the model  $F_{\theta}$  is susceptible to training bias, and to what extent.

### 6.1 Explaining $\Psi_{\mathcal{S}}$

In Fig. 7, we depict two consecutive snapshots of a subgraph, focused on node  $v$ . The dashed circle highlights the corresponding 1-hop neighbourhood  $\mathcal{U}_v$ . At any time  $t$ , we draw green and red arrows next to the nodes to depict whether its node feature value increased or decreased, respectively.

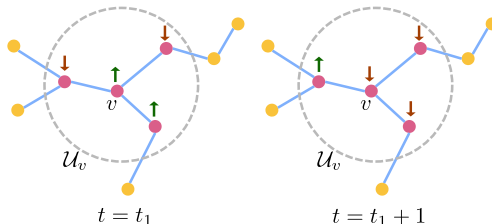


Figure 7: Consecutive subgraph snapshots.

306 The design of  $\Psi_{\mathcal{S}}$  was inspired by the correlation dynamics of the stock market (Caraiani, 2014), where the  
 307 inter-connectedness of various stocks exerts mutual influence on their respective prices. For instance, within  
 308 the semiconductor sector, stocks such as NVDA, AMD, and TSMC often exhibit synchronised movements, with  
 309 slight lead or lag. Similarly, the performance of gold mining stocks can offer insights into the future value of  
 310 physical gold and companies engaged in precious metal trade. This concept transcends individual industries  
 311 and encompasses competition across multiple sectors.

312 Let us record the states at two consecutive time-steps  $\mathbf{s}_{t_1} = [1 \ -1 \ 1 \ -1]^\top$ , and  $\mathbf{s}_{t_1+1} =$   
 313  $[-1 \ -1 \ -1 \ 1]^\top$ . At the state-level, we iterate through the time-steps, and collect all the states succeeding  
 314  $\mathbf{s} = [1 \ -1 \ 1 \ -1]^\top$ . If we then draw a random sample from this collection of succeeding states, we can  
 315 predict whether the node feature value is more likely to *increase or decrease*. However, we are interested in pre-  
 316 dicting the *amount* of change. Therefore, at every time step when the state  $\mathbf{s}_t$  matches  $\mathbf{s} = [1 \ -1 \ 1 \ -1]^\top$ ,  
 317 we collect the succeeding shock  $\varepsilon_{t+1}^{(U_v)}$  in a queue  $\mathcal{Q}_v(\mathbf{s})$ , i.e., at time  $\tau$ ,  $\mathcal{Q}_v(\mathbf{s}) = \{\varepsilon_{t+1}^{(U_v)} : \mathbf{s}_t = \mathbf{s}, \forall t < \tau\}$  with  
 318  $|\mathcal{Q}_v(\mathbf{s})| \leq M$ . The queue entries are then used to approximate a distribution from which a random sample is  
 319 drawn during forecast.

320 In Fig. 8, we plot the normalized histogram of the trace  $\text{tr}(\cdot)$  of the covariance matrix  $\Sigma(\mathbf{s})$  of all the states  
 321  $\mathbf{s} \in \mathcal{S}_v, v \in [n]$  for all the datasets used in multi-step forecasting. We notice that in both PEMS04 and PEMS08  
 322 the distribution of values is skewed to the left, with a concentration of data points at values close to zero.  
 323 This explains the better-than-SoTA performance of `mSPACE-S $\mu$`  on these datasets. In contrast, the histogram  
 324 of METRLA is completely away from zero, while for PEMS03, and PEMS07 there are peaks near zero, but a major  
 325 mass of the histogram is skewed away from zero. This explains the poor performance of `mSPACE-S $\mu$`  on these  
 326 datasets.

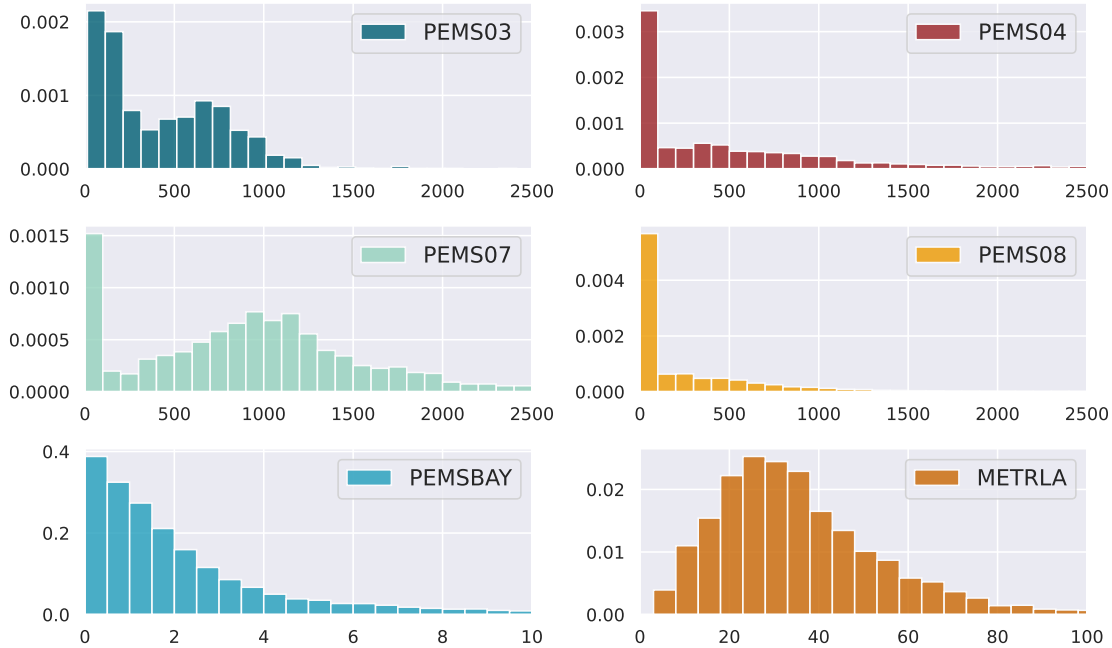


Figure 8: Normalized histogram of  $\{\text{tr}(\Sigma(\mathbf{s})) : \forall \mathbf{s} \in \mathcal{S}_v, \forall v \in [n]\}$  for different datasets.

## 327 6.2 Explaining $\Psi_{\mathcal{T}}$

328 Next, we discuss the rationale behind  $\Psi_{\mathcal{T}}$ , which is designed to identify periodic patterns. For instance, in  
 329 many traffic networks, trends exhibit weekly cycles, with distinct patterns on weekdays compared to weekends.  
 330 Moreover, on an annual basis, the influence of holidays on traffic can be discerned, as people engage in

331 shopping and other leisure activities. In Fig. 9, we have shown the traffic flow value of PEMS04 with weekly (a)  
 332 and daily (b) periodicity. For the weekly periodic view (a), the trend is more pronounced with less deviation  
 from the mean while for the daily view (b), a scattered trend is visible with high variance across states.

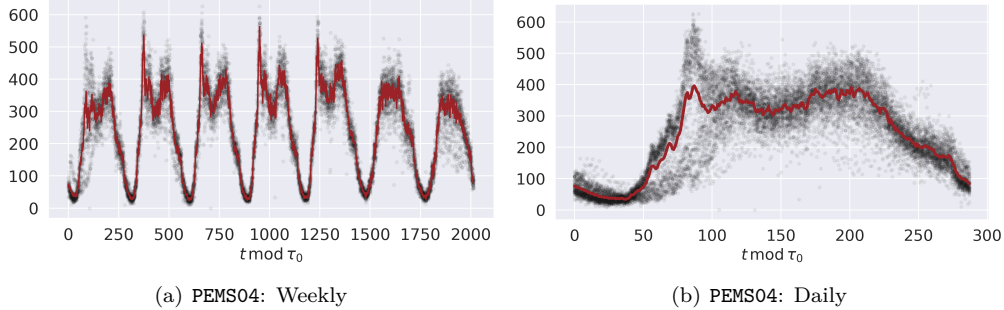


Figure 9: Periodic trends in the traffic dataset PEMS04; the black points represent the data-points, and the red line is the mean estimate for each state  $t \bmod \tau_0$ .

333

### 334 6.3 Error Bounds

335 We present the error bounds of `mSPACE` in the following theorem, a detailed proof of which can be found in  
 336 Appendix A.

337 **Theorem 6.1.** *The RMSE of `mSPACE` for a  $q$ -step node feature forecast is upper bounded as  $\text{RMSE}(q) \leq$   
 338  $\sqrt{\alpha q^2 + (3\alpha + \beta)q + (2\alpha + \beta)}$ , where  $\alpha, \beta \in \mathbb{R}^+$  are constants that depend on the data, as well as the variant  
 339 of the `mSPACE` algorithm.*

340 **Corollary 6.1.** *In the asymptotic case of large  $q$ , the RMSE grows linearly with  $q$ :  $\text{RMSE}(q) = \mathcal{O}(q)$ .*

### 341 6.4 Complexity Analysis

342 We denote the *computational complexity* operator as  $\mathfrak{C}(\cdot)$ , and the *space complexity* operator as  $\mathfrak{M}(\cdot)$ , where  
 343 the argument of each operator is an algorithm or a portion of an algorithm. The optional offline part of  
 344 `mSPACE` is denoted by **A**, while the online part is denoted by **B**. In Table 4, we exhibit the computational  
 345 and space complexities of the different `mSPACE` variants, where  $b \triangleq \max_{v \in [n]} |\mathcal{U}_v|$  is the maximum degree. For  
 more details please refer to Appendix B.

Table 4: Computational and space complexity of different `mSPACE` variants.

	$\Psi_S$	$\Psi_T$
$\Omega_{\mathcal{N}}$	$\mathfrak{C}(\mathbf{A}) = \mathcal{O}(ndb(rT + dbM \min\{rT, 2^{bd}\}))$ $\mathfrak{C}(\mathbf{B}) = \mathcal{O}\left((1-r)Tnd^2b^2\left(qdb + M \min\left\{\frac{(1+r)}{2}T, 2^{bd}\right\}\right)\right)$ $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(db(M + db) \min\{T, 2^{bd}\})$	$\mathfrak{C}(\mathbf{A}) = \mathcal{O}(nrT + d^2Mn\tau_0)$ $\mathfrak{C}(\mathbf{B}) = \mathcal{O}((1-r)Tnd^2(qd + M\tau_0))$ $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(d(M + d)\tau_0)$
$\Omega_{\mu}$	$\mathfrak{C}(\mathbf{A}) = \mathcal{O}(ndb(rT + M \min\{rT, 2^{bd}\}))$ $\mathfrak{C}(\mathbf{B}) = \mathcal{O}\left((1-r)Tndb(q + M) \min\left\{\frac{(1+r)}{2}T, 2^{bd}\right\}\right)$ $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(Mdb \min\{T, 2^{bd}\})$	$\mathfrak{C}(\mathbf{A}) = \mathcal{O}(nrT + dMn\tau_0)$ $\mathfrak{C}(\mathbf{B}) = \mathcal{O}((1-r)Tnd(q + M)\tau_0)$ $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(Md\tau_0)$

346

347 **Theorem 6.2.** *For asymptotically large number of nodes  $n$  and timesteps  $T$ , the computational complexity  
 348 of `mSPACE` is  $\mathcal{O}(nT)$ , and the space complexity is  $\mathcal{O}(1)$  across all variants.*

349 The proof is detailed in Appendix B.2.

## 7 Discussion

In this section we discuss the limitations of `mspace` and how they can be overcome. Firstly, `mspace` only considers binary edges, i.e.,  $\mathbf{A} \in \{0, 1\}^{n \times n}$  instead of a weighted adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . This does not imply that we have used datasets with binary edges, rather it means that we have used a binarized version of the adjacency matrix as input to `mspace` while the baselines exploited weighted edges. Secondly, we assume that the graph structure is fixed throughout, while for a truly dynamic graph, the graph structure should also be dynamic. Lastly, we have proposed two state functions: one that focuses on cross-correlation among the nodes, and the other that considers seasonality. Therefore, a state function which combines both can be studied in an extension of our work in the future.

**On incorporating edge weights** We now investigate how we can incorporate edge weights in `mspace`, and if it has any potential benefits. In addition to Assumption 2.2, consider the following:

**Assumption 7.1.** For nodes  $v, u, u' \in \mathcal{V}$ , if  $|\mathbf{A}_{v,u}| \geq |\mathbf{A}_{v,u'}|$  then  $|\rho(\mathbf{x}(v), \mathbf{x}(u))| \geq |\rho(\mathbf{x}(v), \mathbf{x}(u'))|$ , where  $\mathbf{A}_{v,u} \in \mathbb{R}$  denotes the edge weight of  $(v, u) \in \mathcal{E}$ .

Assumption 2.2 can be applied to both weighted and unweighted graphs, while Assumption 7.1 is applicable only to weighted graphs. It is also evident that Assumption 7.1 is stronger than Assumption 2.2. Therefore, we base `mspace-S*` on Assumption 2.2 and the correlation between the connected nodes are determined intrinsically through the conditional distributions, as the state  $\Psi_{\mathbf{s}}(\boldsymbol{\varepsilon}^{(\mathcal{U}_v)})$  encodes the structural information of a node  $v$  w.r.t its neighbours. However, we can enforce Assumption 7.1 through:

$$\mathbf{s}^* \sim \left\{ \mathbf{s} \in \mathcal{S}_v : \left\| \mathbf{A}_v^{(\mathcal{U}_v)} \odot (\mathbf{s} - \mathbf{s}_t) \right\| < \delta \right\}, \text{ where } \delta \in \mathbb{R}^+.$$

**On adapting to dynamic graph structures** Algorithms that exploit dynamic graph structures are based on the temporal extension of Assumption 7.1, formulated as:

**Assumption 7.2.** For nodes  $u, v \in \mathcal{V}$ , and time-steps  $t, t' \in [T]$ , if  $|\mathbf{A}_{u,v}(t)| \geq |\mathbf{A}_{u,v}(t')|$ , then  $|\rho(\mathbf{x}_t(u), \mathbf{x}_t(v))| \geq |\rho(\mathbf{x}_{t'}(u), \mathbf{x}_{t'}(v))|$ .

Finding the matched state as  $\mathbf{s}^* \sim \left\{ \mathbf{s} \in \mathcal{S}_v : \left\| \mathbf{A}_v^{(\mathcal{U}_v)}(t) \odot (\mathbf{s} - \mathbf{s}_t^{(\mathcal{U}_v)}) \right\| < \delta \right\}, \delta \in \mathbb{R}^+$  makes `mspace` compatible with dynamic graph structure. However, the number of nodes in the graph must remain fixed, i.e., `mspace` cannot deal with node addition or deletion.

**On creating a state function which combines  $\Psi_{\mathbf{s}}$  and  $\Psi_{\mathbf{T}}$**  We can define  $\Psi_{\mathbf{ST}} : \mathbb{R}^{|\mathcal{U}|d} \times \mathbb{N} \rightarrow \{-1, 1\}^{|\mathcal{U}|d} \times \{0, 1, \dots, \tau_0 - 1\}$  as  $\Psi_{\mathbf{ST}}(\boldsymbol{\varepsilon}^{(\mathcal{U})}, t) \triangleq [\text{sign}(\boldsymbol{\varepsilon}^{(\mathcal{U})})^\top \quad t \bmod \tau_0]^\top$ . In essence, the queues  $\mathcal{Q}_v(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}_v, \forall v \in [n]$  in `mspace-ST` would have lesser entries compared to `mspace-S` which might lead to poor estimates and consequently make the algorithm data-intensive. Furthermore, in the step where we find the closest state  $\mathbf{s}^*$ , the spatial and temporal parts can be assigned different weights:  $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \left\| [\mathbf{1}_{d|\mathcal{U}_v|} \quad \gamma]^\top \odot (\mathbf{s} - \mathbf{s}_t^{(\mathcal{U}_v)}) \right\|$ , where  $\gamma \in \mathbb{R}^+$ .

**On benchmarking using diverse datasets** Experiments on more diverse datasets would help establish the performance of the proposed algorithm. In this work, we have used 4 non-traffic datasets for single-step forecasting, and 6 traffic datasets for multi-step. The proposed algorithm `mspace` has a general formulation, and is not designed specifically for traffic datasets; `mspace` can be applied to any graph whose node features (of any dimension) evolve with time. We also proposed a synthetic temporal graph generation method in Appendix C to alleviate the data scarcity issue in temporal graph learning.

## 8 Conclusion

In conclusion, our proposed algorithm, `mspace`, performs at par with the SoTA TGNN models across various spatio-temporal datasets. As an online learning algorithm, `mspace` is adaptive to changes in data distribution and is suitable for deployment in scenarios where training samples are limited. The interpretability of `mspace` sets it apart from black-box deep learning models, allowing for a clearer understanding of the underlying

393 mechanisms driving predictions. This emphasis on interpretability represents a significant step forward in the  
394 field of temporal graph learning. In Sec. 7, we discussed the potential limitations of `mSPACE`, and suggested  
395 design changes through which they can be overcome.

396 In addition to the algorithm, we also introduce a synthetic temporal graph generator (see Appendix C) in  
397 which the features of the nodes evolve with the influence of their neighbours in a non-linear manner. These  
398 synthetic datasets can serve as a valuable resource for benchmarking algorithms.

## 399 References

- 400 David Barber. *Bayesian reasoning and machine learning*. Cambridge University Press, 2012.
- 401 Claudio DT Barros, Matheus RF Mendonça, Alex B Vieira, and Artur Ziviani. A survey on embedding  
402 dynamic graphs. *ACM Computing Surveys (CSUR)*, 55(1):1–37, 2021.
- 403 Ferenc Béres, Róbert Pálovics, Anna Oláh, and András A Benczúr. Temporal walk based centrality metric  
404 for graph streams. *Applied network science*, 3:1–26, 2018.
- 405 G. E. P. Box and David A. Pierce. Distribution of Residual Autocorrelations in Autoregressive-Integrated  
406 Moving Average Time Series Models. *Journal of the American Statistical Association*, 65(332):1509–1526,  
407 December 1970. ISSN 0162-1459. doi: 10.1080/01621459.1970.10481180.
- 408 Petre Caraiani. The predictive power of singular value decomposition entropy for stock market dynamics.  
409 *Physica A: Statistical Mechanics and its Applications*, 393:571–578, 2014.
- 410 Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical Evaluation of Gated  
411 Recurrent Neural Networks on Sequence Modeling, December 2014. URL [http://arxiv.org/abs/1412.](http://arxiv.org/abs/1412.3555)  
412 [3555](http://arxiv.org/abs/1412.3555). arXiv:1412.3555 [cs].
- 413 Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Learning dynamic context graphs for predicting social  
414 events. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery &*  
415 *Data Mining*, pp. 1007–1016, 2019.
- 416 Mengnan Du, Ninghao Liu, and Xia Hu. Techniques for interpretable machine learning. *Communications of*  
417 *the ACM*, 63(1):68–77, 2019.
- 418 Stefanos Eleftheriadis, Tom Nicholson, Marc Deisenroth, and James Hensman. Identification of gaussian  
419 process state space models. *Advances in neural information processing systems*, 30, 2017.
- 420 Zheng Fang, Qingqing Long, Guojie Song, and Kunqing Xie. Spatial-Temporal Graph ODE Networks for  
421 Traffic Flow Forecasting. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery &*  
422 *Data Mining*, pp. 364–373, August 2021. doi: 10.1145/3447548.3467430. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2106.12931)  
423 [2106.12931](http://arxiv.org/abs/2106.12931). arXiv:2106.12931 [cs].
- 424 Chen Gao, Xiang Wang, Xiangnan He, and Yong Li. Graph neural networks for recommender system. In  
425 *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pp. 1623–1625,  
426 2022.
- 427 Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining  
428 explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International*  
429 *Conference on data science and advanced analytics (DSAA)*, pp. 80–89. IEEE, 2018.
- 430 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In  
431 *International Conference on Learning Representations*, 2017. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=SJU4ayYgl)  
432 [SJU4ayYgl](https://openreview.net/forum?id=SJU4ayYgl).
- 433 Zhichen Lai, Dalin Zhang, Huan Li, Christian S Jensen, Hua Lu, and Yan Zhao. LightCTS: A lightweight  
434 framework for correlated time series forecasting. *Proceedings of the ACM on Management of Data*, 1(2):  
435 1–26, 2023.

- 436 Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network:  
437 Data-driven traffic forecasting. In *International Conference on Learning Representations*, 2018. URL  
438 <https://openreview.net/forum?id=SJiHXGWAZ>.
- 439 Zibo Liu, Parshin Shojaee, and Chandan K. Reddy. Graph-based multi-ODE neural networks for spatio-  
440 temporal traffic forecasting. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL  
441 <https://openreview.net/forum?id=0q5XKRVPq>.
- 442 Antonio Longa, Veronica Lachi, Gabriele Santin, Monica Bianchini, Bruno Lepri, Pietro Lio, Andrea Passerini,  
443 et al. Graph neural networks for temporal graphs: State of the art, open challenges, and opportunities.  
444 *Transactions on Machine Learning Research*, 2023.
- 445 Geoffrey J McLachlan, Sharon X Lee, and Suren I Rathnayake. Finite Mixture Models. 2019.
- 446 Alessio Micheli and Domenico Tortorella. Discrete-time dynamic graph echo state networks. *Neurocomputing*,  
447 496:85–95, 2022. Publisher: Elsevier.
- 448 Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler,  
449 Tao Schardl, and Charles Leiserson. Evolvegcn: Evolving graph convolutional networks for dynamic graphs.  
450 In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 5363–5370, 2020. Issue: 04.
- 451 Xuan Rao, Hao Wang, Liang Zhang, Jing Li, Shuo Shang, and Peng Han. FOGS: First-Order Gradient  
452 Supervision with Learning-based Graph for Traffic Flow Forecasting. In *Proceedings of the Thirty-First  
453 International Joint Conference on Artificial Intelligence*, pp. 3926–3932, Vienna, Austria, July 2022.  
454 International Joint Conferences on Artificial Intelligence Organization. ISBN 978-1-956792-00-3. doi:  
455 10.24963/ijcai.2022/545. URL <https://www.ijcai.org/proceedings/2022/545>.
- 456 Benedek Rozemberczki, Paul Scherer, Yixuan He, George Panagopoulos, Alexander Riedel, Maria Aste-  
457 fanoaei, Oliver Kiss, Ferenc Beres, Guzman Lopez, Nicolas Collignon, et al. Pytorch geometric temporal:  
458 Spatiotemporal signal processing with neural machine learning models. In *Proceedings of the 30th ACM  
459 international conference on information & knowledge management*, pp. 4564–4573, 2021a.
- 460 Benedek Rozemberczki, Paul Scherer, Oliver Kiss, Rik Sarkar, and Tamas Ferenci. Chickenpox cases in  
461 hungary: A benchmark dataset for spatiotemporal signal processing with graph neural networks. In  
462 *Workshop on Graph Learning Benchmarks@ TheWebConf 2021*, 2021b.
- 463 Robert H. Shumway and David S. Stoffer. *Time Series Analysis and Its Applications: With R Exam-  
464 ples*. Springer Texts in Statistics. Springer International Publishing, Cham, 2017. ISBN 978-3-319-  
465 52451-1 978-3-319-52452-8. doi: 10.1007/978-3-319-52452-8. URL [http://link.springer.com/10.1007/  
466 978-3-319-52452-8](http://link.springer.com/10.1007/978-3-319-52452-8).
- 467 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser,  
468 and Illia Polosukhin. Attention is All you Need. *Neural Information Processing Systems*, 2017.
- 469 Yifei Wang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin. Dissecting the diffusion process in linear  
470 graph convolutional networks. In *NeurIPS*, pp. 5758–5769, 2021. URL [https://proceedings.neurips.  
471 cc/paper/2021/hash/2d95666e2649fcfc6e3af75e09f5adb9-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/2d95666e2649fcfc6e3af75e09f5adb9-Abstract.html).
- 472 Greg Welch. An Introduction to the Kalman Filter. 1997.
- 473 Xinle Wu, Dalin Zhang, Chenjuan Guo, Chaoyang He, Bin Yang, and Christian S Jensen. AutoCTS:  
474 Automated correlated time series forecasting. *Proceedings of the VLDB Endowment*, 15(4):971–983, 2021a.
- 475 Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-  
476 temporal graph modeling. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*,  
477 pp. 1907–1913, 2019.
- 478 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A Comprehensive  
479 Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*,  
480 32(1):4–24, January 2021b. ISSN 2162-237X, 2162-2388. doi: 10.1109/TNNLS.2020.2978386. URL  
481 <https://ieeexplore.ieee.org/document/9046288/>.

- 482 Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: a deep learning  
483 framework for traffic forecasting. In *Proceedings of the 27th International Joint Conference on Artificial  
484 Intelligence*, pp. 3634–3640, 2018.
- 485 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In  
486 *Proceedings of the AAAI conference on Artificial Intelligence*, volume 37, pp. 11121–11128, 2023.
- 487 Ling Zhao, Yujiao Song, Chao Zhang, Yu Liu, Pu Wang, Tao Lin, Min Deng, and Haifeng Li. T-gcn: A  
488 temporal graph convolutional network for traffic prediction. *IEEE transactions on intelligent transportation  
489 systems*, 21(9):3848–3858, 2019. Publisher: IEEE.

## 490 A Error Bounds

491 **Upper Bound** We derive the upper bound on the RMSE for  $q$ -step iterative forecast below.

492 *Proof of Theorem 6.1.* For nodes in  $\mathcal{U}_v, v \in [n]$ , the shock at time  $t$  is sampled from a Gaussian distribution,  
493 the parameters of which depend on the previous shock  $\hat{\varepsilon}_{t-1}^{\langle \mathcal{U}_v \rangle}$  through the state function:

$$\hat{\varepsilon}_t^{\langle \mathcal{U}_v \rangle} \sim \mathcal{N}\left(\hat{\varepsilon}; \boldsymbol{\mu}\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t-1}^{\langle \mathcal{U}_v \rangle}\right)\right), \boldsymbol{\Sigma}\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t-1}^{\langle \mathcal{U}_v \rangle}\right)\right)\right) \quad (5)$$

494 We denote the shock estimated for node  $v$  at time  $t$  as:

$$\hat{\varepsilon}_t(v) = \hat{\varepsilon}_t^{\langle \mathcal{U}_v \rangle}(v) \sim \mathcal{N}\left(\hat{\varepsilon}; \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t-1}^{\langle \mathcal{U}_v \rangle}\right)\right), \boldsymbol{\Sigma}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t-1}^{\langle \mathcal{U}_v \rangle}\right)\right)\right) \quad (6)$$

495 The mean square error for  $q$ -step iterative node feature forecasting is defined as:

$$\begin{aligned} \text{MSE}(q) &\triangleq \frac{1}{ndq} \mathbb{E} \left[ \sum_{v \in [n]} \sum_{i \in [q]} \left\| \sum_{j \in [i]} \hat{\varepsilon}_{t+j}(v) - \varepsilon_{t+j}(v) \right\|^2 \right] \\ &= \frac{1}{ndq} \sum_{v \in [n]} \sum_{i \in [q]} \mathbb{E} \left[ \left\| \sum_{j \in [i]} \hat{\varepsilon}_{t+j}(v) - \varepsilon_{t+j}(v) \right\|^2 \right]. \end{aligned} \quad (7)$$

496 The shock difference between the true shock and predicted shock also follows a Gaussian distribution:

$$\hat{\varepsilon}_{t+j}(v) - \varepsilon_{t+j}(v) \sim \mathcal{N}\left(\boldsymbol{\varepsilon}; \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v), \boldsymbol{\Sigma}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right)\right). \quad (8)$$

497 Since, the sum of Gaussian r.v.s is also Gaussian, we have:

$$\sum_{j \in [i]} \hat{\varepsilon}_{t+j}(v) - \varepsilon_{t+j}(v) \sim \mathcal{N}\left(\boldsymbol{\varepsilon}; \sum_{j \in [i]} \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v), \sum_{j \in [i]} \boldsymbol{\Sigma}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right)\right). \quad (9)$$

498 Moreover, for a Gaussian r.v.  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbb{E} \left[ \|\mathbf{x}\|^2 \right] = \|\boldsymbol{\mu}\|^2 + \text{tr}(\boldsymbol{\Sigma})$ .

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{j \in [i]} \hat{\varepsilon}_{t+j}(v) - \varepsilon_{t+j}(v) \right\|^2 \right] &= \left\| \sum_{j \in [i]} \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v) \right\|^2 \\ &\quad + \sum_{j \in [i]} \text{tr}\left(\boldsymbol{\Sigma}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right)\right). \end{aligned} \quad (10)$$

499

$$\begin{aligned} \left\| \sum_{j \in [i]} \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v) \right\| &\leq \sum_{j \in [i]} \left\| \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v) \right\| \\ &\leq i \cdot \max_{j \in [i]} \left\| \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v) \right\| \\ &\leq i \cdot \max_{t, j \in \mathbb{N}} \left\| \boldsymbol{\mu}_v\left(\Psi_{\mathbf{S}}\left(\hat{\varepsilon}_{t+j-1}^{\langle \mathcal{U}_v \rangle}\right)\right) - \varepsilon_{t+j}(v) \right\| \\ &= i \cdot \sqrt{\alpha_{v,1}}. \end{aligned} \quad (11)$$



500

$$\sum_{j \in [i]} \text{tr} \left( \boldsymbol{\Sigma}_v \left( \Psi_{\mathbf{S}} \left( \hat{\boldsymbol{\epsilon}}_{t+j-1}^{\langle \mathcal{U}_v \rangle} \right) \right) \right) \leq i \cdot \max_{j \in [i]} \text{tr} \left( \boldsymbol{\Sigma}_v \left( \Psi_{\mathbf{S}} \left( \hat{\boldsymbol{\epsilon}}_{t+j-1}^{\langle \mathcal{U}_v \rangle} \right) \right) \right) \leq i \cdot \alpha_{v,2}. \quad (12)$$

$$\mathbb{E} \left[ \left\| \sum_{j \in [i]} \hat{\boldsymbol{\epsilon}}_{t+j}(v) - \boldsymbol{\epsilon}_{t+j}(v) \right\|^2 \right] \leq \alpha_{v,1} \cdot i^2 + \alpha_{v,2} \cdot i, \quad \alpha_{v,1}, \alpha_{v,2} \in \mathbb{R}^+. \quad (13)$$

$$\begin{aligned} \text{MSE}(q) &\leq \frac{1}{ndq} \sum_{v \in [n]} \sum_{i \in [q]} \alpha_{v,1} \cdot i^2 + \alpha_{v,2} \cdot i \\ &= \frac{\sum_{v \in [n]} \alpha_{v,1}}{6nd} (q+1)(q+2) + \frac{\sum_{v \in [n]} \alpha_{v,2}}{2nd} (q+1). \end{aligned} \quad (14)$$

501 Let  $\alpha \triangleq \frac{1}{6nd} \sum_{v \in [n]} \alpha_{v,1}$ , and  $\beta \triangleq \frac{1}{2nd} \sum_{v \in [n]} \alpha_{v,2}$ , then

$$\text{MSE}(q) \leq \alpha q^2 + (3\alpha + \beta)q + (2\alpha + \beta). \quad (15)$$

502 By Jensen's inequality,

$$\text{RMSE}(q) \leq \sqrt{\text{MSE}(q)} \leq \sqrt{\alpha q^2 + (3\alpha + \beta)q + (2\alpha + \beta)}. \quad (16)$$

503

□

504 The above proof is for `mspace-SN` and also applies to `mspace-TN`. For `mspace-Sμ` and `mspace-Tμ`,  $\beta = 0$ .

505 **Lower Bound** Similarly, we can find a lower bound on the MSE for  $q$ -step iterative forecast:

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{j \in [i]} \hat{\boldsymbol{\epsilon}}_{t+j}(v) - \boldsymbol{\epsilon}_{t+j}(v) \right\|^2 \right] &\geq \sum_{j \in [i]} \text{tr} \left( \boldsymbol{\Sigma}_v \left( \Psi_{\mathbf{S}} \left( \hat{\boldsymbol{\epsilon}}_{t+j-1}^{\langle \mathcal{U}_v \rangle} \right) \right) \right) \\ &\geq i \cdot \min_{j \in [i]} \text{tr} \left( \boldsymbol{\Sigma}_v \left( \Psi_{\mathbf{S}} \left( \hat{\boldsymbol{\epsilon}}_{t+j-1}^{\langle \mathcal{U}_v \rangle} \right) \right) \right) = i \cdot \alpha_{v,3}. \end{aligned} \quad (17)$$

506

$$\text{MSE}(q) \geq \frac{1}{ndq} \sum_{v \in [n]} \sum_{i \in [q]} i \cdot \alpha_{v,3} = \underbrace{\left( \frac{1}{nd} \sum_{v \in [n]} \alpha_{v,3} \right)}_{\triangleq \beta'} \cdot (q+1) = \beta' q + \beta'. \quad (18)$$

## 507 B Complexity Analysis

### 508 B.1 Computational Complexity

509 We denote the computational complexity operator as  $\mathfrak{C}(\cdot)$ , the argument of which is an algorithm or part of  
510 an algorithm. The optional offline part of the algorithm is denoted as A while the online part is denoted as B.

---

#### Algorithm 2 mspace-S $\mathcal{N}$

---

**Input**  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$ ,  $r \in [0, 1)$ ,  $q, M$

**Output**  $\hat{\boldsymbol{\varepsilon}}_t(v)$ ,  $\forall v \in \mathcal{V}, t \in [[r \cdot T], T]$

1:  $\boldsymbol{\varepsilon}_t(v) \leftarrow \mathbf{x}_t(v) - \mathbf{x}_{t-1}(v)$ ,  $\forall v \in \mathcal{V}, t \in [T]$   
*Offline training (A):*  
2: **for**  $t \in [[r \cdot T], T]$  **do**  
3:   **for**  $v \in \mathcal{V}$  **do**  
4:      $\mathbf{s}_t^{(\mathcal{U}_v)} \leftarrow \Psi(\boldsymbol{\varepsilon}_t^{(\mathcal{U}_v)})$   $\triangleright \sum_{v \in \mathcal{V}} d|\mathcal{U}_v|$   
5:      $\mathcal{S}_v \leftarrow \mathcal{S}_v \cup \{\mathbf{s}_t^{(\mathcal{U}_v)}\}$   $\triangleright n$   
6:      $\mathcal{Q}_v(\mathbf{s}_t^{(\mathcal{U}_v)}) \leftarrow \text{enqueue } \boldsymbol{\varepsilon}_{t+1}^{(\mathcal{U}_v)}$   $\triangleright n$   
7:   **end for**  
8: **end for**  
9:  $\boldsymbol{\mu}_v(\mathbf{s}) \leftarrow \text{mean}(\mathcal{Q}_v(\mathbf{s}))$ ,  $\forall \mathbf{s} \in \mathcal{S}_v, v \in \mathcal{V}$   $\triangleright \sum_{v \in \mathcal{V}} d|\mathcal{U}_v||\mathcal{S}_v|M$   
10:  $\boldsymbol{\Sigma}_v(\mathbf{s}) \leftarrow \text{covariance}(\mathcal{Q}_v(\mathbf{s}))$ ,  $\forall \mathbf{s} \in \mathcal{S}_v, v \in \mathcal{V}$   $\triangleright \sum_{v \in \mathcal{V}} (d|\mathcal{U}_v|)^2 |\mathcal{S}_v|M$   
*Online learning (B):*  
11: **for**  $t \in [[r \cdot T], T - q]$  **do**  
12:   **for**  $v \in \mathcal{V}$  **do**  
13:      $\mathbf{s}_t^{(\mathcal{U}_v)} \leftarrow \Psi(\boldsymbol{\varepsilon}_t^{(\mathcal{U}_v)})$   $\triangleright \sum_{v \in \mathcal{V}} d|\mathcal{U}_v|$   
14:      $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \|\mathbf{s} - \mathbf{s}_t^{(\mathcal{U}_v)}\|$   $\triangleright \sum_{v \in \mathcal{V}} d|\mathcal{U}_v||\mathcal{S}_v|$   
15:      $\hat{\boldsymbol{\varepsilon}}_{t+1}^{(\mathcal{U}_v)} \sim \mathcal{N}(\boldsymbol{\varepsilon}; \boldsymbol{\mu}_v(\mathbf{s}^*), \boldsymbol{\Sigma}_v(\mathbf{s}^*))$   $\triangleright \sum_{v \in \mathcal{V}} (|\mathcal{U}_v|d)^3$   
16:     **for**  $k \in [2, q]$  **do**  
17:        $\mathbf{s}^* \leftarrow \arg \min_{\mathbf{s} \in \mathcal{S}_v} \|\mathbf{s} - \Psi(\hat{\boldsymbol{\varepsilon}}_{t+k-1}^{(\mathcal{U}_v)})\|$   $\triangleright (q-1) \times \sum_{v \in \mathcal{V}} d|\mathcal{U}_v|(1 + |\mathcal{S}_v|)$   
18:        $\hat{\boldsymbol{\varepsilon}}_{t+k}^{(\mathcal{U}_v)} \sim \mathcal{N}(\boldsymbol{\varepsilon}; \boldsymbol{\mu}_v(\mathbf{s}^*), \boldsymbol{\Sigma}_v(\mathbf{s}^*))$   $\triangleright (q-1) \times \sum_{v \in \mathcal{V}} (|\mathcal{U}_v|d)^3$   
19:     **end for**  
20:      $\hat{\boldsymbol{\varepsilon}}_{t+k}(v) \leftarrow \hat{\boldsymbol{\varepsilon}}_{t+k}^{(\mathcal{U}_v)}(v)$ ,  $\forall k \in [q]$   
21:     Update  $\mathcal{S}_v, \mathcal{Q}_v$   $\triangleright 2n$   
22:     Update  $\boldsymbol{\mu}_v(\mathbf{s}), \boldsymbol{\Sigma}_v(\mathbf{s})$ ,  $\forall \mathbf{s} \in \mathcal{S}_v$   $\triangleright \sum_{v \in \mathcal{V}} (d|\mathcal{U}_v| + d^2|\mathcal{U}_v|^2)|\mathcal{S}_v|M$   
23:   **end for**  
24: **end for**

---

511 Computational complexity of offline training for mspace-S $\mathcal{N}$  can be written as:

$$\mathfrak{C}(\text{A}) = \mathcal{O} \left( \underbrace{[rT]d \sum_v |\mathcal{U}_v|}_{[4]} + \underbrace{[rT]2n}_{[5],[6]} + \underbrace{dM \sum_v |\mathcal{U}_v||\mathcal{S}_v|}_{[9](\text{mean})} + \underbrace{d^2M \sum_v |\mathcal{U}_v|^2|\mathcal{S}_v|}_{[10](\text{covariance})} \right). \quad (19)$$

512 Computational complexity of online learning for `mSPACE-SN` can be written as:

$$\begin{aligned} \mathfrak{C}(\mathbf{B}) = \mathcal{O} \left( \sum_{t=\lceil rT \rceil}^{T-q} \left\{ \underbrace{dq \sum_v |\mathcal{U}_v|}_{[13],[17]} + \underbrace{dq \sum_v |\mathcal{U}_v| |\mathcal{S}_v|}_{[14],[17]} + \underbrace{d^3 q \sum_v |\mathcal{U}_v|^3}_{[15],[18](\text{sampling})} + \underbrace{2n}_{[21]} \right. \right. \\ \left. \left. + \underbrace{dM \sum_v |\mathcal{U}_v| |\mathcal{S}_v|}_{[22](\text{mean})} + \underbrace{d^2 M \sum_v |\mathcal{U}_v|^2 |\mathcal{S}_v|}_{[22](\text{covariance})} \right\} \right). \end{aligned} \quad (20)$$

513 **Lemma B.1.** *The computational complexity of `mSPACE-SN` is:*

$$\begin{aligned} \mathfrak{C}(\mathbf{A}) &= \mathcal{O} \left( dbnrT + d^2 b^2 Mn \cdot \min\{rT, 2^{bd}\} \right), \\ \mathfrak{C}(\mathbf{B}) &= \mathcal{O} \left( (1-r)Tnd^2 b^2 \left( qdb + M \cdot \min \left\{ \frac{(1+r)}{2} T, 2^{bd} \right\} \right) \right), \end{aligned}$$

514 where  $b = \max_{v \in [n]} |\mathcal{U}_v|$ .

515 *Proof.* We denote the maximum degree of a node as  $b \triangleq \max_{v \in [n]} |\mathcal{U}_v| < n$  which does not necessarily scale  
516 with  $n$  unless specified by the graph definition. Furthermore, the total number of states observed for a node  
517 till time step  $t \in \mathbb{N}$  cannot exceed  $t$ , i.e.,  $|\mathcal{S}_v| \leq t$ . We also know the total number of states theoretically  
518 possible for node  $v$  is  $2^{|\mathcal{U}_v|d}$  for  $\Psi_{\mathcal{S}}(\cdot)$ . Therefore, the number of states observed till time  $t$  for node  $v$  is upper  
519 bounded as:  $|\mathcal{S}_v| \leq \min\{t, 2^{bd}\}$ . Based on this, we can simplify equation 19, and equation 20 as follows:

$$\begin{aligned} \mathfrak{C}(\mathbf{A}) &= \mathcal{O} \left( dbnrT + 2nrT + (dbM + d^2 b^2 M) \cdot n \min\{rT, 2^{bd}\} \right) \\ &= \mathcal{O} \left( dbnrT + d^2 b^2 Mn \cdot \min\{rT, 2^{bd}\} \right). \\ \mathfrak{C}(\mathbf{B}) &= \mathcal{O} \left( \sum_{t=\lceil rT \rceil}^{T-q} qdbn + qd^3 b^3 n + 2n + db(q+M)n \cdot \min\{t, 2^{bd}\} + d^2 b^2 Mn \cdot \min\{t, 2^{bd}\} \right) \\ &= \mathcal{O} \left( \sum_{t=\lceil rT \rceil}^{T-q} qd^3 b^3 n + (db(q+M) + d^2 b^2 M)n \cdot \min\{t, 2^{bd}\} \right) \\ &= \mathcal{O} \left( (1-r)T \cdot qd^3 b^3 n + d^2 b^2 Mn \cdot \min\{(1-r^2)T^2, 2^{bd}(1-r)T\} \right) \\ &= \mathcal{O} \left( (1-r)Tn \left( qd^3 b^3 + d^2 b^2 M \cdot \min \left\{ \frac{(1+r)}{2} T, 2^{bd} \right\} \right) \right). \end{aligned}$$

520

□

521 **Lemma B.2.** *The computational complexity of `mSPACE-Sμ` is:*

$$\begin{aligned} \mathfrak{C}(\mathbf{A}) &= \mathcal{O} \left( dbnrT + dbMn \cdot \min\{rT, 2^{bd}\} \right), \\ \mathfrak{C}(\mathbf{B}) &= \mathcal{O} \left( (1-r)Tndb(q+M) \cdot \min \left\{ \frac{(1+r)}{2} T, 2^{bd} \right\} \right). \end{aligned}$$

522

523 *Proof.* The sampling steps [15], and [18] in Algorithm 2 are replaced with  $\hat{\boldsymbol{\epsilon}}_t^{(\mathcal{U}_v)} \leftarrow \boldsymbol{\mu}(\mathbf{s}^*)$  which has a  
524 computational complexity of  $\mathcal{O}(d|\mathcal{U}_v|)$ . Moreover,  $\Omega_{\mu}(\cdot)$  does not require the covariance matrix, therefore we

525 do not need to compute it. We simplify the computational complexity expressions as:

$$\begin{aligned}
\mathfrak{C}(\mathbf{A}) &= \mathcal{O}\left(\lfloor rT \rfloor d \sum_v |\mathcal{U}_v| + \lfloor rT \rfloor 2n + dM \sum_v |\mathcal{U}_v| |\mathcal{S}_v|\right) \\
&= \mathcal{O}\left(dbnrT + dbMn \cdot \min\{rT, 2^{bd}\}\right). \\
\mathfrak{C}(\mathbf{B}) &= \mathcal{O}\left(\sum_{t=\lfloor rT \rfloor}^{T-q} \left\{ dq \sum_v |\mathcal{U}_v| + dq \sum_v |\mathcal{U}_v| |\mathcal{S}_v| + \underbrace{dq \sum_v |\mathcal{U}_v|}_{(\text{sampling})} + 2n + dM \sum_v |\mathcal{U}_v| |\mathcal{S}_v| \right\}\right) \\
&= \mathcal{O}\left(\sum_{t=\lfloor rT \rfloor}^{T-q} 2qdbn + 2n + db(q+M)n \cdot \min\{t, 2^{bd}\}\right) \\
&= \mathcal{O}\left((1-r)Tndb(q+M) \cdot \min\left\{\frac{(1+r)}{2}T, 2^{bd}\right\}\right).
\end{aligned}$$

526

□

527 **Lemma B.3.** *The computational complexity of  $m\text{space-TN}$  is:*

$$\begin{aligned}
\mathfrak{C}(\mathbf{A}) &= \mathcal{O}(nrT + d^2Mn\tau_0), \\
\mathfrak{C}(\mathbf{B}) &= \mathcal{O}\left((1-r)Tnd^2 \cdot (M\tau_0 + qd)\right).
\end{aligned}$$

528

529 *Proof.* For the state function  $\Psi_{\mathbf{T}}$ , the total number of states for any node is the period  $\tau_0 \in \mathbb{N}$ , i.e.,  $|\mathcal{S}_v| \leq \tau_0$ .  
530 Moreover, the state calculation  $\mathbf{s}_t \leftarrow \Psi(t)$  has computational complexity of  $\mathcal{O}(1)$ . Most importantly, for  $\Psi_{\mathbf{T}}$ ,  
531  $b = 1$  as it only focuses on the seasonal trends.

$$\begin{aligned}
\mathfrak{C}(\mathbf{A}) &= \mathcal{O}\left(\lfloor rT \rfloor \sum_v 1 + \lfloor rT \rfloor 2n + dM \sum_v |\mathcal{U}_v| |\mathcal{S}_v| + d^2M \sum_v |\mathcal{U}_v|^2 |\mathcal{S}_v|\right) \\
&= \mathcal{O}(3nrT + dMn\tau_0 + d^2Mn\tau_0) = \mathcal{O}(nrT + d^2Mn\tau_0). \\
\mathfrak{C}(\mathbf{B}) &= \mathcal{O}\left(\sum_{t=\lfloor rT \rfloor}^{T-q} \left\{ q \sum_v 1 + dq \sum_v |\mathcal{U}_v| |\mathcal{S}_v| + d^3q \sum_v |\mathcal{U}_v|^3 + 2n \right. \right. \\
&\quad \left. \left. + dM \sum_v |\mathcal{U}_v| |\mathcal{S}_v| + d^2M \sum_v |\mathcal{U}_v|^2 |\mathcal{S}_v| \right\}\right) \\
&= \mathcal{O}\left(\{q + dq\tau_0 + qd^3 + 2 + dM\tau_0 + d^2M\tau_0\} \cdot n(1-r)T\right) \\
&= \mathcal{O}\left((1-r)Tnd^2 \cdot (M\tau_0 + qd)\right).
\end{aligned}$$

532

□

533 **Lemma B.4.** *The computational complexity of  $m\text{space-T}\mu$  is:*

$$\begin{aligned}
\mathfrak{C}(\mathbf{A}) &= \mathcal{O}(nrT + dMn\tau_0), \\
\mathfrak{C}(\mathbf{B}) &= \mathcal{O}\left((1-r)Tn \cdot d(q+M)\tau_0\right).
\end{aligned}$$

534

535 *Proof.* Based on the explanation provided for `mSPACE-TN`, we simplify the computational complexity expres-  
536 sions for `mSPACE-Tμ` as:

$$\begin{aligned}\mathfrak{C}(\mathbf{A}) &= \mathcal{O}\left(\lfloor rT \rfloor \sum_v 1 + \lfloor rT \rfloor 2n + dM \sum_v |\mathcal{S}_v|\right) \\ &= \mathcal{O}(3nrT + dMn\tau_0) = \mathcal{O}(nrT + dMn\tau_0). \\ \mathfrak{C}(\mathbf{B}) &= \mathcal{O}\left(\sum_{t=\lfloor rT \rfloor}^{T-q} \left\{q \sum_v 1 + dq \sum_v |\mathcal{S}_v| + 2n + dM \sum_v |\mathcal{S}_v|\right\}\right) \\ &= \mathcal{O}\left(\{q + dq\tau_0 + 2 + dM\tau_0\} \cdot n(1-r)T\right) = \mathcal{O}\left((1-r)Tn \cdot d(q+M)\tau_0\right).\end{aligned}$$

537

□

## 538 B.2 Space Complexity

539 We denote the space complexity operator as  $\mathfrak{M}(\cdot)$ , the argument of which is an algorithm or part of an  
540 algorithm. The variables in offline training  $\mathbf{A}$  are re-used in online learning  $\mathbf{B}$ . Therefore, we can say that  
541  $\mathfrak{M}(\mathbf{B}) = \mathfrak{M}(\mathbf{A} \cup \mathbf{B})$ .

542 In an implementation of `mSPACE` where forecasting is sequentially performed for each node  $v \in [n]$ , memory  
543 space can be efficiently reused, except for storing the outputs. This approach optimises memory usage,  
544 resulting in a space complexity characterised by:

$$\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}\left(\max_{\substack{v \in [n], \\ t \in [T]}} \underbrace{d|\mathcal{U}_v||\mathcal{S}_v|}_{\mathcal{S}_v} + \underbrace{cMd|\mathcal{U}_v||\mathcal{S}_v|}_{\mathcal{Q}_v(\mathbf{s}) \forall \mathbf{s} \in \mathcal{S}_v} + \underbrace{cd|\mathcal{U}_v||\mathcal{S}_v|}_{\mu_v(\mathbf{s}) \forall \mathbf{s} \in \mathcal{S}_v} + \underbrace{c(d|\mathcal{U}_v|)^2|\mathcal{S}_v|}_{\Sigma_v(\mathbf{s}) \forall \mathbf{s} \in \mathcal{S}_v} + \underbrace{d|\mathcal{U}_v|}_{\mathbf{s}^*}\right). \quad (21)$$

545 **Lemma B.5.** *The space complexity of `mSPACE-SN` is  $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(db(M+db) \cdot \min\{T, 2^{bd}\})$ .*

546 *Proof.* Simplifying equation 21 results in:

$$\begin{aligned}\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) &= \mathcal{O}\left(\max_{\substack{v \in [n], \\ t \in [T]}} (db + cMdb + cdb + cd^2b^2)|\mathcal{S}_v| + db\right) \\ &= \mathcal{O}\left((cMdb + cd^2b^2) \cdot \max_{t \in [T]} \min\{t, 2^{bd}\}\right) = \mathcal{O}\left(db(M+db) \cdot \min\{T, 2^{bd}\}\right).\end{aligned}$$

547

□

548 **Lemma B.6.** *The space complexity of `mSPACE-Sμ` is  $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(Mdb \cdot \min\{T, 2^{bd}\})$ .*

549 *Proof.* Some space is saved in `mSPACE-Sμ`, as we do not need to store the covariance matrices.

$$\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}\left(\max_{\substack{v \in [n], \\ t \in [T]}} (db + cMdb + cdb)|\mathcal{S}_v| + db\right) = \mathcal{O}\left(Mdb \cdot \min\{T, 2^{bd}\}\right).$$

550

□

551 **Lemma B.7.** *The space complexity of `mSPACE-TN` is  $\mathfrak{M}(\mathbf{A} \cup \mathbf{B}) = \mathcal{O}(d(M+d)\tau_0)$ .*

552 *Proof.* As explained earlier, for the state function  $\Psi_T$ ,  $b = 1$ . Therefore, the queues only store the shock  
 553 vectors for a single node, and not the neighbours. The space complexity expression is simplified as:

$$\mathfrak{M}(A \cup B) = \mathcal{O} \left( \max_{\substack{v \in [n], \\ t \in [T]}} (d + cMd + cd + cd^2) |\mathcal{S}_v| + db \right) = \mathcal{O} \left( d(M + d)\tau_0 \right).$$

554

□

555 **Lemma B.8.** *The space complexity of `mspace-Tμ` is  $\mathcal{O}(Md\tau_0)$ .*

556 *Proof.*  $\mathfrak{M}(A \cup B) = \mathcal{O} \left( \max_{\substack{v \in [n], \\ t \in [T]}} (d + cMd + cd) |\mathcal{S}_v| + d \right) = \mathcal{O}(Md\tau_0)$ . □

557 **Asymptotic Analysis** Theorem 6.2 states that *for asymptotically large number of nodes  $n$  and timesteps*  
 558  *$T$ , the computational complexity of `mspace` is  $\mathcal{O}(nT)$ , and the space complexity is  $\mathcal{O}(1)$  across all variants.*

559 *Proof.* We analyse the lemmas B.1-B.8 introduced in this section for the asymptotic case of very large  $n$   
 560 and  $T$ . For very large  $T$ ,  $\min \left\{ \frac{(1+r)}{2}T, 2^{bd} \right\} \rightarrow 2^{bd}$ . Similarly,  $\min \{T, 2^{bd}\} \rightarrow 2^{bd}$ . Considering the terms  
 561  $r, d, M, q, \tau_0, b$  as constants, the computational complexity for both offline and online parts of all the `mspace`  
 562 variants becomes  $\mathcal{O}(nT)$  for asymptotically large  $n, T$ .

563 Furthermore, the space complexity terms lack  $n$  or  $T$  for very large  $T$ , which allows us to conclude that the  
 564 space complexity of all the variants of `mspace` is constant, i.e.,  $\mathcal{O}(1)$ . □

## 565 C Synthetic Datasets & Experiments

566 In traffic datasets, seasonality outweighs cross-nodal correlation, making it challenging to assess the efficacy  
 567 of a TGL algorithms on node feature forecasting task. To address this gap, we propose a synthetic dataset  
 568 generation technique in line with the design idea of `mspace` which is described in Algorithm 3.

---

### Algorithm 3 Synthetic Data Generation

---

**Input**  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ ,  $d$ ,  $\mu_{\min}$ ,  $\mu_{\max}$ ,  $\sigma_{\min}^2$ ,  $\sigma_{\max}^2$ ,  $\mu_0$ ,  $\sigma_0^2$ ,  $\tau$ ,  $\mu_\tau$ ,  $\sigma_\tau^2$ .

```

1:  $\varepsilon_0 \sim \text{Bernoulli}^{nd} \left( \frac{1}{2} \right)$ 
2:  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{x}; \mu_0 \mathbf{1}, \sigma_0^2 \mathbf{I})$ 
3: for  $t \in [T]$  do
4:    $\mathbf{s}_{t-1} \leftarrow \Psi_S(\varepsilon_{t-1})$ 
5:   if  $\mathbf{s}_{t-1} \notin \mathcal{S}$  then
6:      $\mathcal{S} \leftarrow \mathcal{S} \cup \{\mathbf{s}_{t-1}\}$ 
7:      $\boldsymbol{\mu}(\mathbf{s}_{t-1}) \sim \text{Uniform}^{nd}(\mu_{\min}, \mu_{\max})$ 
8:      $\tilde{\boldsymbol{\Sigma}} \sim \text{Uniform}^{nd \times nd}(\sigma_{\min}^2, \sigma_{\max}^2)$ 
9:      $\hat{\boldsymbol{\Sigma}} \leftarrow \frac{1}{2} (\tilde{\boldsymbol{\Sigma}} + \tilde{\boldsymbol{\Sigma}}^\top)$ 
10:     $\boldsymbol{\Sigma}(\mathbf{s}_{t-1}) \leftarrow \hat{\boldsymbol{\Sigma}} \odot (\mathbf{A} \otimes \mathbf{1}_{d \times d})$ 
11:   end if
12:    $\varepsilon_t \sim \mathcal{N}(\varepsilon; \boldsymbol{\mu}(\mathbf{s}_{t-1}), \boldsymbol{\Sigma}(\mathbf{s}_{t-1}))$ 
13:    $\mathbf{x}_t = \mathbf{x}_{t-1} + \varepsilon_t$ 
14: end for
15: if  $\tau > 0$  then
16:    $\mathbf{y}_t \sim \mathcal{N}(\mathbf{y}; \mu_\tau \mathbf{1}, \sigma_\tau^2 \mathbf{I}) \quad \forall t \in [\tau]$ 
17:    $\mathbf{x}_t \leftarrow \mathbf{x}_t + \mathbf{y}_{t \bmod \tau} \quad \forall t \in [T]$ 
18: end if
```

---

569 In steps 8-10, we construct a covariance matrix adhering to Assumption 2.2, and in step 12, we sample the  
 570 shock from a multivariate normal distribution. In steps 16-17, a random signal  $\mathbf{y}$  is tiled with period  $\tau$  and  
 571 added to the node features to introduce seasonality into the dataset.

572 The synthetic datasets can be utilized to analyze how various factors such as graph structure, periodicity,  
 573 connectivity, sample size, and other parameters affect error metrics.

574 We generate datasets through Algorithm 3 by supplying the parameters outlined in Table 5. For each dataset,  
 575 we create multiple random instances and report the mean and standard deviation of the metrics in the results.

Table 5: Parameters for different synthetic dataset packages.

Dataset	$\mathcal{G} \sim$	$d$	$T$	$\mu_{\min}$	$\mu_{\max}$	$\sigma_{\min}$	$\sigma_{\max}$	$\mu_0$	$\sigma_0$	$\tau$	$\mu_\tau$	$\sigma_\tau$
SYN01	$\mathcal{G}_{\text{ER}}(20, 0.2)$	1	$10^3$	-200	200	40	50	$2 \times 10^4$	5000	100	100	20
SYN02	$\mathcal{G}_{\text{ER}}(20, 0.2)$	1	$10^3$	-200	200	40	50	$2 \times 10^4$	5000	0		
SYN03	$\mathcal{G}_{\text{ER}}(40, 0.5)$	1	$10^3$	-400	400	30	40	$10^4$	2000	0		
SYN04	$\mathcal{G}_{\text{ER}}(40, 0.5)$	1	$10^4$	-400	400	30	40	$10^4$	2000	0		

576



Figure 10: Exemplary synthetic dataset samples shown for 5 nodes.

### 577 C.1 Periodicity

578 The generator parameters for SYN01 and SYN02 are same except for the periodic component added to SYN01  
 579 which has a period of  $\tau = 100$  timesteps consisting of shocks sampled from  $\mathcal{N}(100, 20)$ . An algorithm which  
 580 can exploit the periodic influence in the signal should perform better on SYN01 compared to SYN02. The  
 581 models which perform worse on periodic dataset are marked **red**.

Table 6: Impact of data periodicity on RMSE achieved by different models.

	SYN01		SYN02		% increase
	mean	std. dev.	mean	std. dev.	$\left(\frac{\text{SYN02}-\text{SYN01}}{\text{SYN01}}\right)$
mSPACE-S $\mu$	299.18	$\pm$ 6.55	294.99	$\pm$ 8.81	-0.63
mSPACE-S $\mathcal{N}$	400.99	$\pm$ 3.74	395.33	$\pm$ 3.24	-1.52
STGODE	420.86	$\pm$ 103.29	420.25	$\pm$ 52.17	-9.87
GRAM-ODE	921.94	$\pm$ 537.63	853.77	$\pm$ 340.45	-18.18
LightCTS	419.43	$\pm$ 176.5	334.59	$\pm$ 79.01	-30.6
Kalman- $x$	781.94	$\pm$ 32.35	776.75	$\pm$ 30.38	-0.88
Kalman- $\epsilon$	393.76	$\pm$ 4.72	390.45	$\pm$ 3.54	-1.13

## 582 C.2 Training Samples

583 The generator parameters for SYN03 and SYN04 are same except for the total number of samples being ten  
584 times more in SYN04. If a model perform better on SYN04 compared to SYN03, it would indicate that it is  
585 training intensive, requiring more samples to infer the trends. On the other hand, if the model performs  
586 worse on SYN04, it would indicate that there are scalability issues, or the training caused overfitting. An  
587 ideal model is expected to have similar performance on SYN03 and SYN04. The models with ideal behaviour  
588 are marked teal, and the models susceptible to overfitting are marked red. Moreover, model(s) that require  
589 more training samples are marked violet.

Table 7: Impact of number of training samples on RMSE achieved by different models.

	SYN03		SYN04		% increase
	mean	std. dev.	mean	std. dev.	$\left(\frac{\text{SYN04}-\text{SYN03}}{\text{SYN03}}\right)$
mSPACE-S $\mu$	793.41	$\pm$ 5.86	789.36	$\pm$ 3	-0.86
mSPACE-S $\mathcal{N}$	793.93	$\pm$ 5.73	792.61	$\pm$ 2.02	-0.63
STGODE	830.63	$\pm$ 127	931.33	$\pm$ 191.87	+17.29
GRAM-ODE	1382.48	$\pm$ 80.78	1423.93	$\pm$ 190.13	+10.31
LightCTS	769.34	$\pm$ 196.6	998.01	$\pm$ 319.72	+36.42
Kalman- $x$	785.7	$\pm$ 8.95	721.88	$\pm$ 1.73	-8.94
Kalman- $\epsilon$	782.6	$\pm$ 6.5	783.36	$\pm$ 1.45	-0.54

## 590 D Evaluation

### 591 D.1 Metrics

592 The root mean squared error (RMSE) of  $q$  consecutive predictions for all the nodes is:

$$\text{RMSE}(q) \triangleq \mathbb{E} \left[ \sqrt{\frac{1}{ndq} \sum_{v \in \mathcal{V}} \sum_{i \in [q]} \left\| \sum_{j \in [i]} \epsilon_{t+j}(v) - \hat{\epsilon}_{t+j}(v) \right\|_2^2} \right]. \quad (22)$$

593 The mean absolute error (MAE) of  $q$  consecutive predictions for all the nodes is:

$$\text{MAE}(q) \triangleq \frac{1}{ndq} \mathbb{E} \left[ \sum_{v \in \mathcal{V}} \sum_{i \in [q]} \left\| \sum_{j \in [i]} \epsilon_{t+j}(v) - \hat{\epsilon}_{t+j}(v) \right\|_1 \right]. \quad (23)$$

### 594 D.2 Datasets

595 In Table 8, we list the datasets commonly utilised in the literature for single and multi-step node feature  
596 forecasting.



597 **tennis (Béres et al., 2018)** represents a discrete-time dynamic graph showing the hourly changes in the  
 598 interaction network among Twitter users during the 2017 Roland-Garros (RG17) tennis match. The input  
 599 features capture the structural attributes of the vertices, with each vertex symbolizing a different user and  
 600 the edges indicating retweets or mentions within an hour <sup>4</sup>.

601 **wikimath (Rozemberczki et al., 2021a)** tracks daily visits to Wikipedia pages related to popular  
 602 mathematical topics over a two-year period. Static edges denote hyperlinks between the pages <sup>5</sup>.

603 **pedalme (Rozemberczki et al., 2021a)** reports weekly bicycle package deliveries by Pedal Me in London  
 604 throughout 2020 and 2021. The nodes are different locations, and the edge weight encodes the physical  
 605 proximity. The count of weekly bicycle deliveries in a location forms the node feature footnote <sup>6</sup>.

606 **cpox (Rozemberczki et al., 2021b)** tracks the weekly number of chickenpox cases for each county of  
 607 Hungary between 2005 and 2015. Different counties form the nodes, and are connected if any two counties  
 608 share a border <sup>6</sup>.

609 **PEMS03/04/07/08 (Rao et al., 2022)** The four datases are collected from four districts in California  
 610 using the California Transportation Agencies (CalTrans) Performance Measurement System (PeMS) and  
 611 aggregated into 5-minutes windows<sup>7</sup>. The spatial adjacency matrix for each dataset is constructed using the  
 612 length of the roads. PEMS03 is collected from September 2018 to November 2018. PEMS04 is collected from  
 613 San Francisco Bay area from July 2016 to August 2016. PEMS07 is from Los Angeles and Ventura counties  
 614 between May 2017 and August 2017. PEMS08 is collected from San Bernardino area between July 2016 to  
 615 August 2016.

616 *Variables:* The **flow** represents the number of vehicles that pass through the loop detector per time interval  
 617 (5 minutes). The **occupancy** variable represents the proportion of time during the time interval that the  
 618 detector was occupied by a vehicle. It is measured as a percentage. Lastly, the **speed** variable represents the  
 619 average speed of the vehicles passing through the loop detector during the time interval. It is measured in  
 620 miles per hour (mph).

621 **PEMSBAY (Li et al., 2018)** is a traffic dataset collected by CalTrans PeMS. It is represented by a  
 622 network of 325 traffic sensors in the Bay Area with 6 months of traffic readings ranging from January 2017 to  
 623 May 2017 in 5 minute intervals<sup>8</sup>.

624 **METRLA (Li et al., 2018)** is a traffic dataset based on Los Angeles Metropolitan traffic conditions.  
 625 The traffic readings are collected from 207 loop detectors on highways in Los Angeles County over 5 minute  
 626 intervals between March 2012 to June 2012<sup>9</sup>.

### 627 D.3 Baselines

628 **DCRNN (Li et al., 2018)** The Diffusion Convolutional Recurrent Neural Network (DCRNN) models the  
 629 node features as a diffusion process on a directed graph, capturing spatial dependencies through bidirectional  
 630 random walks. Additionally, it addresses nonlinear temporal dynamics by employing an encoder-decoder  
 631 architecture with scheduled sampling.

632 **TGCN (Zhao et al., 2019)** Temporal Graph Convolutional Network (TGCN) combines the graph convolu-  
 633 tional network (GCN) with a gated recurrent unit (GRU), where the former learns the spatial patterns, and  
 634 the latter learns the temporal.

<sup>4</sup><https://github.com/ferencberes/online-centrality>

<sup>5</sup>wikimath dataset from PyTorch Geometric Temporal

<sup>6</sup>[https://github.com/benedekrozemberczki/spatiotemporal\\_datasets](https://github.com/benedekrozemberczki/spatiotemporal_datasets)

<sup>7</sup><https://github.com/guoshnBJTU/ASTGNN/tree/main/data>

<sup>8</sup>PEMSBAY dataset from PyTorch Geometric Temporal

<sup>9</sup>METRLA dataset from PyTorch Geometric Temporal

Table 8: Real world datasets for single and multi-step forecasting.

Name	$n$	$\mathbf{x}$	time-step	$T$
tennis	1,000	# tweets	1 hour	120
wikimath	1,068	# visits	1 day	731
pedalme	15	# deliveries	1 week	35
cpox	20	# cases	1 week	520
PEMS03	358	flow	5 min	26,208
PEMS04	307	flow, occupancy, speed	5 min	16,992
PEMS07	883	flow	5 min	28,224
PEMS08	170	flow, occupancy, speed	5 min	17,856
PEMSBAY	325	speed	5 min	52,116
METRLA	207	speed	5 min	34,272

635 **EGCN (Pareja et al., 2020)** EvolveGCN (EGCN) adapts a GCN model without using node embeddings.  
 636 The evolution of the GCN parameters is learnt through an RNN. EGCN has two variants: ECGN-H which uses  
 637 a GRU, and ECGN-O which uses an LSTM.

638 **DynGESN (Micheli & Tortorella, 2022)** Dynamic Graph Echo State Networks (DynGESN) employ echo  
 639 state networks (ESNs) a special type of RNN in which the recurrent weights are conditionally initialized,  
 640 while a memory-less readout layer is trained. The ESN evolves through state transitions where the states  
 641 belong to a compact space. For more details please refer to the original text.

642 **GWNet (Wu et al., 2019)** GraphWave Net (GWNet) consists of an adaptive dependency matrix which  
 643 is learnt through node embeddings, which is capable of capturing the hidden spatial relations in the data.  
 644 GWNet can handle long sequences owing to its one-dimensional convolutional component whose receptive field  
 645 grows exponentially with the number of layers.

646 **STGODE (Fang et al., 2021)** Spatial-temporal Graph Ordinary Differential Equation (STGODE) employs  
 647 tensor-based ordinary differential equations (ODEs) to model the temporal evolution of the node features.

648 **GRAM-ODE (Liu et al., 2023)** Graph-based Multi-ODE (GRAM-ODE) improves upon STGODE by con-  
 649 necting multiple ODE-GNN modules to capture different views of the local and global spatiotemporal  
 650 dynamics.

651 **FOGS (Rao et al., 2022)** FOGS utilises first-order gradients to train a predictive model because the traffic  
 652 data distribution is irregular.

653 **LightCTS (Lai et al., 2023)** LightCTS stacks temporal and spatial operators in a computationally-efficient  
 654 manner, and uses lightweight modules L-TCN and GL-Former.

655 **ARIMA (Box & Pierce, 1970)** ARIMA is a multivariate time series forecasting technique that combines  
 656 autoregressive, integrated, and moving average components. It models the relationship between observations  
 657 and their lagged values, adjusts for non-stationarity in the data, and accounts for short-term fluctuations.

658 **Kalman (Welch, 1997)** Since `mspace` is a state-space algorithm, we also use the Kalman filter as a  
 659 baseline. We introduce two variants of the Kalman filter: `Kalman-x`, which considers the node features as  
 660 observations, and `Kalman-ε`, which operates on the shocks.