

# Uncertainty-aware Traffic Prediction under Missing Data

Hao Mei

*School of Computing and Augmented Intelligence  
Arizona State University  
Tempe, US  
hmei7@asu.edu*

Zhiming Liang

*Faculty of Electronic and Information Engineering  
Xi'an Jiaotong University  
Xi'an, China  
zhimingliang@stu.xjtu.edu.cn*

Bin Shi

*Faculty of Electronic and Information Engineering  
Xi'an Jiaotong University  
Xi'an, China  
shibin@xjtu.edu.cn*

Junxian Li

*Faculty of Electronic and Information Engineering  
Xi'an Jiaotong University  
Xi'an, China  
ljx201806@stu.xjtu.edu.cn*

Guanjie Zheng

*School of Electronic Information and Electrical Engineering  
Shanghai Jiaotong University  
Shanghai, China  
gjzheng@sjtu.edu.cn*

Hua Wei\*

*School of Computing and Augmented Intelligence  
Arizona State University  
Tempe, US  
hwei27@asu.edu*

**Abstract**—Traffic prediction is a crucial topic because of its broad scope of applications in the transportation domain. Recently, various studies have achieved promising results. However, most studies assume the prediction locations have complete or at least partial historical records and cannot be extended to non-historical recorded locations. In real-life scenarios, the deployment of sensors could be limited due to budget limitations and installation availability, which makes most current models not applicable. Though few pieces of literature tried to impute traffic states at the missing locations, these methods need the data simultaneously observed at the locations with sensors, making them not applicable to prediction tasks. Another drawback is the lack of measurement of uncertainty in prediction, making prior works unsuitable for risk-sensitive tasks or involving decision-making. To fill the gap, inspired by the previous inductive graph neural network, this work proposed an uncertainty-aware framework with the ability to 1) extend prediction to missing locations with no historical records and significantly extend spatial coverage of prediction locations while reducing deployment of sensors and 2) generate probabilistic prediction with uncertainty quantification to help the management of risk and decision making in the down-stream tasks. Through extensive experiments on real-life datasets, the result shows our method achieved promising results on prediction tasks, and the uncertainty quantification gives consistent results which highly correlated with the locations with and without historical data. We also show that our model could help support sensor deployment tasks in the transportation field to achieve higher accuracy with a limited sensor deployment budget.

**Index Terms**—1. Traffic prediction, 2. Data mining 3. Uncertainty quantification

## I. INTRODUCTION

Traffic prediction is an important problem in urban computing with numerous applications ranging from urban mobility systems to autonomous vehicle operations [1], [2]. In recent years, the availability of large-scale traffic datasets, coupled with advancements in information collection infrastructures, has led to increased attention and analysis of these datasets [3]–[5]. Notably, significant progress has been made in traffic prediction accuracy, primarily driven by data-driven approaches and the proliferation of deep learning models. For instance, Ma et al. [6] leverage Convolutional Neural Networks to model spatial correlations, while Liu et al. [7] employ Graph Neural Networks to capture spatiotemporal dependencies using the diffusion mechanism.

Despite the promising results achieved in recent studies on traffic prediction, there are two main limitations in most existing works. Firstly, many assume that all prediction locations are equipped with sensors and can access historical data. However, in real-life scenarios, due to budget constraints or limited accessibility, it is common for specific locations of interest to be absent from sensor coverage. Consequently, existing prediction models only apply to locations with historical data observable, limiting prediction capabilities at locations missing historical observation. There is a critical need for models that can forecast observable and missing locations. While some methods have focused on the kriging problem, which involves data imputation using surrounding location data [3], [8], [9], these approaches require simultaneous access to current data

\*Corresponding author.

and cannot be directly applied to predict future data.

Secondly, most current traffic prediction studies primarily report the accuracy of their models, neglecting the crucial aspect of prediction uncertainty. In many downstream tasks, especially those involving risk management or decision-making based on prediction results, the quantification of uncertainty is essential in determining whether a method can be implemented in real-life. Failure to assess potential risks can lead to significant costs in high-stakes transportation tasks. Some recent works have attempted to address this issue. For example, Wu et al. [10] quantified uncertainty in spatiotemporal prediction tasks by discussing different uncertainty quantification methods from frequentist and Bayesian perspectives. Similarly, Prob-GNNs [11] proposed a probabilistic graph neural network framework and investigated various probabilistic assumptions specific to ride-sharing demand tasks. However, these methods assume the availability of historical data at each location for evaluating prediction uncertainty. Importantly, locations missing historical observation exhibit distinct topology properties within their surrounding sub-graphs, which locations with historical data may not adequately represent. As a result, predictions at these locations are characterized by low confidence and need to be discerned among all points of interest. Thus, there is an urgent need for a new framework that enables predictions while uncovering the spatiotemporal pattern of uncertainty and extends to locations missing historical observation.

Motivated by the concept of inductive graph neural networks [9], initially employed to solve spatial kriging problems, we adapt this idea to the context of traffic prediction and propose an **Uncertainty-awarded Inductive Graph Neural Network (UIGNN)** to attack this problem. Our framework can extend to new locations missing historical data while integrating uncertainty quantification methods to assess predictions comprehensively. The details description is shown in Fig. 1. The main contributions of our proposed framework are as follows:

- Our work is the first to investigate traffic prediction problems concerning the ability to be spatially extended to locations missing historical data in the road network. Our proposed single-step model (directly predict) in traffic prediction problems outperforms the other two-step (impute then predict) methods on the missing data locations. At the same time, keep the same level of performance at observable data locations.
- To investigate the model performance at different locations, we integrate uncertainty quantification to evaluate the prediction quality at each road network location since the error-based metric is not accessible at missing data locations. The result shows the uncertainty reported by our model can successfully distinguish observable and missing locations and reflect the prediction accuracy. And we also justify our model in the downstream active sensing tasks by using uncertainty to assist deployment of sensors.
- We investigate our model on three real-world datasets with historical data observable at partial locations. The result shows our model performs better than two-step approaches and is effective in real-world sensor deployment tasks.

## II. RELATED WORK

### A. Traffic Prediction

Traffic prediction has been a challenging problem due to the complex dependencies in the time and space domain. Recently, with the help of data-driven approaches, many successes have been achieved and significantly reduced the gap for its application for downstream tasks in real life. However, most of these methods are either suffering from either strong stationary or linear dependencies assumptions [12], [13], making them not suitable for real-life tasks with highly complicated and stochastic characterization in nature. Deep learning models also provide promising results and more flexibility for spatiotemporal dependencies modeling. Convolutional neural networks(CNN) have been used to model the regular 2-D structured traffic network in the transportation domain and have promising results in predicting crowd flow [14]. Graph neural network is another counterpart to model more complex spatial structure [15]. The graph convolutional network (GCN) proposed by Bruna et al. [16] first bridged the gap between the deep neural network and the graph spectral theory [16] and then introduced in the traffic prediction problem by [17], [18]. Graph attention network (GAT) is an alternative approach that can learn the attentional weights on each graph node and is broadly used in spatiotemporal forecasting problems [19].

Though plenty of research exists in solving the traffic prediction problem, most rely on the assumption that the historical data of each location of interest in the prediction task are accessible. However, in real life, this assumption does not always hold. One way to solve this problem is to kriging the unrecorded locations and then forecast the future time point. Some statical models, like using mean or zero to fill the unrecorded locations or regression-based models [20], matrix factorization [21], are also used. However, the completed data used in training the prediction model will also bring in errors and weaken the prediction accuracy at locations with historical data. To avoid these problems, inspired by the inductive graph neural network [9], we proposed a new framework that can directly execute prediction tasks at both locations with historical data and locations whose historical data are not accessible because of reasons like sensors not deployed at these points of interest.

### B. Inductive Graph Learning

In recent years, graph neural network has attracted much attention due to their high expressiveness in capturing intricate relationships and dependencies of structured data and their flexibility to model irregular data [22]–[24]. After the proven of both transductive and inductive representation capability [25], researchers have developed many inductive GNN methods in applications such as recommendation systems from learning embedding of each node in a fixed graph to learning to embed of node features that can be generated to unseen nodes [26]. For example, Zhang et al. [27] train a GCN by masking a part of the observed user and item embedding and learning to reconstruct the masked ones. And Zeng [25] trains GNNs

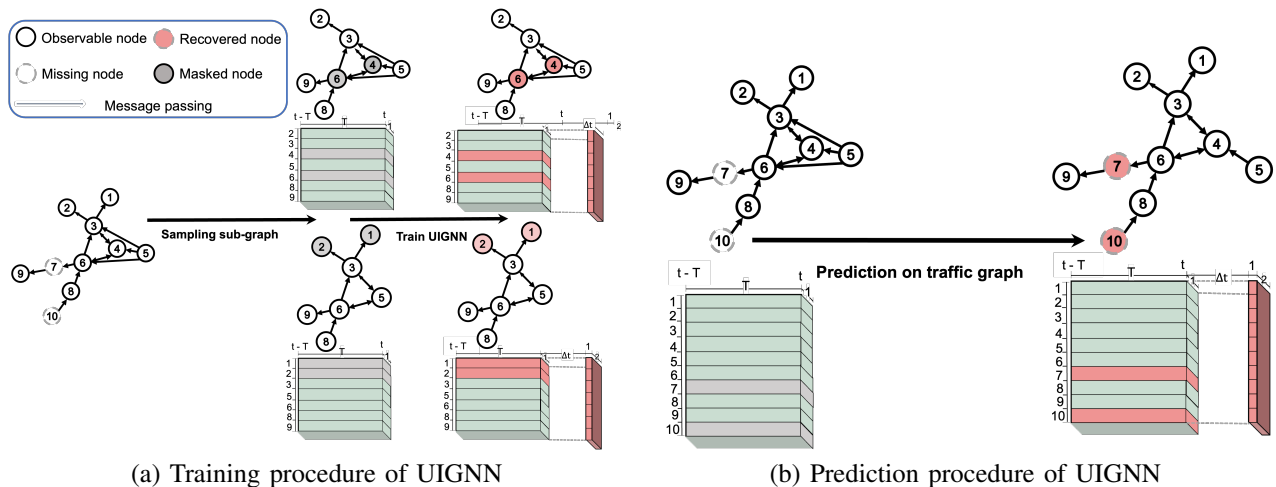


Fig. 1. Framework of UIGNN. (a) The missing locations indexes are 7,10. And during training, we random sample sub-graph and mask locations 4,6 (upper sample) and train GNN to impute mask locations from time  $t - T$  to  $t$  and predict all sampled locations at  $t + \Delta t$  and quantify the prediction uncertainty. (b) During prediction, UIGNN predicts both missing and observable locations at time  $t + \Delta t$  and quantifies the prediction uncertainty.

for large graphs by sampling and learning from the subgraph. In inductive learning, learning node embedding through the structural properties of a node's neighbor properties instead of learning each node's embedding in a fixed graph is suitable. It can be easily adapted to the prediction problem at locations with and without historical data, as we sample subgraphs at locations with historical data and train inductive GNN models. Then the trained GNN model could treat locations without historical data as unseen nodes and embed the unseen nodes using the learned topological structure of these node's neighborhoods and the distribution of the node features in the neighborhood [25]. Thus in this paper, we take the setting from the previous inductive graph neural network [9] and extend the task from spatiotemporal kriging into the prediction tasks with the ability to extend the perception to locations without historical data.

### C. Uncertainty Quantification

Research in uncertainty quantification related to deep learning has become a thriving field in recent years. There are two mainstream uncertainty quantification methods, mostly beyond consideration. The first category is the Bayesian method which models the posterior distribution of the network parameters optimized with the data. Methods including dropout [28], [29], ensembling [30], and other approaches [31], [32] are taken to quantify the variance and place the priors over network weights using variational inference [33]. A recent method, the evidential deep learning method, alternatively trains a deterministic model and places the uncertainty before the predictive distribution instead [34] and can simultaneously detect OOD and adversarial data. The second category is the frequentist uncertainty quantification methods which focus on the robustness, not the variations in the data [10]. Existing methods quantify uncertainty by taking the prediction intervals as the prediction objective [35] or with bootstrapping using the influence function [36] for example. In this work, we

investigate from both Bayesian and frequentist perspectives on our spatiotemporal forecasting problem. And we pay more attention to the evidential deep learning approach in the Bayesian framework as it can distinguish the OOD samples in the data, which can best help explain the high uncertainty predictions in prediction locations with no historical data.

### III. PRELIMINARY

In this section, we first introduce definitions and the original problem statement of traffic prediction problems following the previous works [7]. Then we will extend the original traffic prediction problem to missing historical data locations and give a formal problem statement. All the notation is summarized in Table I for brevity.

TABLE I  
NOTATIONS

Notations	Descriptions
$N$	Number of locations of interest in the network
$N_o, N_e$	Number of observable/missing locations
$\mathcal{V}$	Nodes representing the locations of interest in the traffic network
$\mathcal{A}$	Adjacent matrix representing the connectivity of the traffic network
$\mathcal{A}_o$	Adjacent matrix representing the connectivity of observable locations
$\mathcal{G}(\cdot, \cdot)$	Graph representation of the network
$\mathcal{V}_o, \mathcal{V}_e$	Nodes at observable/missing locations
$\mathbf{X}^{(t)}$	Node features at observable locations and 0 at missing locations at time step $t$
$\mathbf{X}_o^{(t)}$	Node features at observable locations at time step $t$
$\hat{\mathbf{X}}_o^{(t)}, \hat{\mathbf{X}}_e^{(t)}$	Predicted node features at observable/missing locations at time step $t$
$h(\cdot)$	The function to predict future node features at $\mathcal{V}_o$
$f(\cdot)$	The function to predict future node features at all locations of interest $\mathcal{V}_o$ and $\mathcal{V}_e$

Suppose we have predefined  $N$  locations of interest in the traffic network which are deemed critical and have a strong

desire to acquire timely information. We present this traffic network as a directed graph  $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ , where  $\mathcal{V}$  is a set of nodes and the cardinality of the nodes  $|\mathcal{V}|$  is  $N$ .  $\mathcal{A}$  is the adjacent matrix describing the connectivity of the network  $\mathcal{G}$  and  $\mathcal{A} \in \mathbb{R}^{N \times N}$ .

In real life, due to the accessibility of sensors at some given locations,  $N_o$  out of  $N$  locations, which are a subset  $\mathcal{V}_o \subseteq \mathcal{V}$  (also known as the observed locations) represented on  $\mathcal{G}$  have sensors deployed and historical data  $\mathbf{X}_o \in \mathbb{R}^{N_o \times P}$  are observable, where  $P$  is the number of features of each node.  $\mathbf{X}_o^{(t)}$  represents the traffic state observed at time  $t$  and location  $\mathcal{V}_o$ . And the sub-graph of these observed locations is described as  $\mathcal{G}_o = (\mathcal{V}_o, \mathcal{A}_o)$ , where  $\mathcal{A}_o$  is the adjacent matrix describing the connectivity between observed locations. For the rest of the locations of interest with no historical records at  $\mathcal{V}_e \subseteq \mathcal{V}$  (also known as the missing location), we denote historical data at time  $t$  as  $\mathbf{X}_e^{(t)} \in \mathbb{R}^{N_e \times P}$ , where  $N_e$  is the number of missing locations.

**Problem 1 (Traffic Prediction Problem).** *In the traffic prediction problem, at time step  $t$ , we have  $T$  length of historical traffic state  $\mathbf{X}_o^{(t-T+1)}, \mathbf{X}_o^{(t-T+2)}, \dots, \mathbf{X}_o^{(t)}$  at locations  $\mathcal{V}_o$  and the network structure  $\mathcal{G}$ , the goal is to learn a function  $h(\cdot)$  that takes the historical traffic states and graph structure and predict the traffic states at locations  $\mathcal{V}_o$  at the future time  $t + \Delta t$ :*

$$[\mathbf{X}_o^{(t-T+1)}, \mathbf{X}_o^{(t-T+2)}, \dots, \mathbf{X}_o^{(t)}, \mathcal{G}] \xrightarrow{h(\cdot)} \hat{\mathbf{X}}_o^{(t+\Delta t)} \quad (1)$$

In real-life scenarios, due to budget or physical restrictions, only part of the locations are observable and have accessibility to historical records. However, the information in other missing locations is also critical and of interest; for example, the predicted future speed could help plan routines for vehicles or to control traffic signals for better traffic regulation; we need to predict traffic state at both observable locations  $\mathcal{V}_o$  and extend the perception to missing locations  $\mathcal{V}_e$ . Based on the differences mentioned above, we extend the original problem to the traffic prediction with missing locations problem and formally define it as follows:

**Problem 2 (Traffic Prediction with Missing Locations Problem).** *In the traffic prediction with missing locations problem, following the original setting [7], we have the historical observations at locations  $\mathcal{V}_o$  and the network structure. In the new problem, the goal is to learn a function  $f(\cdot)$  that takes the historical traffic state at locations  $\mathcal{V}_o$  and traffic network graph structure  $\mathcal{G}$  and predicts the future traffic state at both observed locations  $\mathcal{V}_o$  and missing locations  $\mathcal{V}_e$  at the future time points  $t + \Delta t$ :*

$$[\mathbf{X}_o^{(t-T+1)}, \mathbf{X}_o^{(t-T+2)}, \dots, \mathbf{X}_o^{(t)}, \mathcal{G}] \xrightarrow{f(\cdot)} [\hat{\mathbf{X}}_o^{(t+\Delta t)}, \hat{\mathbf{X}}_e^{(t+\Delta t)}] \quad (2)$$

#### IV. METHOD

While traffic prediction is a well-studied field, making prediction extendable to locations with no historical data is still

challenging due to the lack of prior knowledge and the complex spatial dependencies on the traffic network. To overcome these problems, we adopt the inductive graph neural network, which is used initially in kriging tasks, and extend it to traffic prediction. At the same time, we integrate the uncertainty quantification component into our prediction model to evaluate the confidence of the prediction in the absence of ground truth data. The proposed uncertainty-aware inductive neural network (UIGNN) framework can predict traffic state at observable locations with historical data and extend prediction perception to missing locations with no historical data while quantifying the uncertainty of predictions. In the rest of this section, we will briefly introduce the UIGNN framework.

##### A. Traffic Prediction with Missing Locations

In the setting defined by Problem 2, the model  $f(\cdot)$  should be applicable to predict observed locations  $\mathcal{V}_o$  and missing locations  $\mathcal{V}_e$ . Most current methods in conventional traffic prediction problems cannot be directly used in this scenario because of the lack of training data at the missing locations. To stress this problem, inspired by graph sampling-based inductive learning methods [37] and its recent application in spatiotemporal kriging problem [9], we take an inductive learning method that models the traffic network with graph neural network and take a sampling-based method to learn the representation of nodes in the sampled graph and generalize those representations to unseen nodes and graphs. The inductive graph neural network (inductive GNN) can learn the topological structure and distribution of the node features in the neighborhood and is thus applicable to graphs both with and without node features w.r.t observable locations and missing locations, respectively. In this work, we follow the past IGNNK [9] and extend the task from kriging node features  $\hat{\mathbf{X}}_e^{(t)} = g(\mathbf{X}_o^{(t-T+1)}, \mathbf{X}_o^{(t-T+2)}, \dots, \mathbf{X}_o^{(t)}, \mathcal{G})$ , where  $g(\cdot)$  is the learned function to kriging the missing locations, into the traffic prediction with missing location problem. Below are the details of how we take the inductive learning method and train the model to learn the node representation only on observable locations and generalize the learned node representation to predict future traffic states on observable and missing locations the model has not seen.

1) **Inductive Graph Neural Network:** To overcome the lack of historical data and learn a node representation generalized to  $\mathcal{V}_e$ , we follow the implementation of IGNNK to train an inductive graph neural network by sampling sub-graph node features at  $\mathcal{V}_o$ . Following the inductive learning approach, there is an assumption that the observable locations represent the overall population distribution, and the message-passing mechanism could be generated to unseen locations  $\mathcal{V}_e$ .

Considering our task is to predict traffic network state with only partial locations observable, during training, our inductive model should also be able to learn node representation with part of node features available such that the learned model could predict all locations  $\mathcal{V}$  with only historical data at observed locations  $\mathcal{V}_o$ . Thus we take a sampling strategy to train our model on  $\mathcal{G}_s \subseteq \mathcal{G}_o$  instead of  $\mathcal{G}_o$ .

2) **Random Sample of Sub-graph:** To prepare the sample for training GNN, we first generate random integers  $n_s \leq n_o$ , which is the number of sampled nodes  $\mathcal{V}_s \subseteq \mathcal{V}_o$ . And next, we split sampled nodes  $\mathcal{V}_s$  into two groups  $\mathcal{V}_r$  and  $\mathcal{V}_m$  with carnality of  $|\mathcal{V}_r| = n_r, |\mathcal{V}_m| = n_m, n_r + n_m = n_s$ . And we reserve the observed historical data at  $\mathcal{V}_r$  and mask the observed historical data at  $\mathcal{V}_m$  to make node features partially accessible. Thus we have the set of sampled nodes, reserved nodes, and masked nodes in the network following the relationship:  $\mathcal{V}_s = \mathcal{V}_r \cup \mathcal{V}_m$ . To make annotation simple, we use  $\mathcal{V} \subset \mathcal{G}$  with slight abuse to represent the node  $\mathcal{V}$  itself and also its index in the graph  $\mathcal{G}$ . For example,  $i \in \mathcal{V}_r \subset \mathcal{G}_s$  means the set of index belongs to observed nodes in the sampled graph  $\mathcal{G}_s$ .

Then we generate the mask matrix  $\mathcal{M}_s \in \mathbb{R}^{n_s \times P}$  with

$$\mathcal{M}_s[i, :] = \begin{cases} 1, & \text{if } i \in \mathcal{V}_r \subset \mathcal{G}_s, \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

and construct adjacent matrix  $\mathcal{A}_s$  of sub-graph  $\mathcal{G}_s$  by selecting entries at  $\mathcal{A}_{i,j}$  of  $\mathcal{A}$ , where  $i, j \in \mathcal{V}_s \subset \mathcal{G}$ . And the sampled graph could be represent as  $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{A}_s)$

We finally randomly sample time step  $t$  and pick  $\mathbf{X}_s^{(t-T+1:t)}$  at  $\mathcal{V}_s$  from  $t-T+1$  to  $t$  time interval, where  $T$  is the length of the historical data taken by the model as input node features and  $\mathbf{X}_s^{(t+\Delta t)}$  as the target.

We iterate for  $I$  iterations, and in each iteration, we sample sub-graph  $\mathcal{G}_s$  for  $S$  times and train the GNN model  $f(\cdot)$  with the samples described above. The detail is shown in Algorithm 1. Then we apply the trained  $f(\cdot)$  in the prediction task on the graph  $\mathcal{G}$ .

3) **Diffusion Graph Convolutional Network:** To capture the stochastic nature of spatiotemporal dependencies in the directional network, we follow the past work and adopt Diffusion Graph Convolutional Network [38] (DGCNs) as the basic building block for the GNN architecture. The DGCNs model the message passing on the GNN as a diffusion process that explicitly captures the stochastic nature of the transition dynamics and is theoretically and experimentally justified applied to spatiotemporal forecasting problems. The basic building block could be represented as:

$$H_{l+1} = \sum_{k=1}^K T_k(\bar{\mathcal{A}}_f) H_l \Theta_{b,l}^k + T_k(\bar{\mathcal{A}}_b) H_l \Theta_{f,l}^k \quad (4)$$

where  $\bar{\mathcal{A}}_b$  is  $\mathcal{A}$  normalized across each row and  $\bar{\mathcal{A}}_b = \bar{\mathcal{A}}_f.\text{transpose}()$ . The two matrices represent the forward and backward transition matrices, respectively. We have  $T_k(\bar{\mathcal{A}}) = 2\bar{\mathcal{A}}T_{k-1}(\bar{\mathcal{A}}) - T_{k-2}(\bar{\mathcal{A}})$  and  $T_0(\bar{\mathcal{A}}) = I, T_1(\bar{\mathcal{A}}) = \bar{\mathcal{A}}$  for the first and second order, where  $k$  is the order of the Chebyshev polynomial used in the DGCN block. And the  $\Theta_{f,l}^k$  and  $\Theta_{b,l}^k$  are learnable parameters of GNN model. The input layer in our traffic forecasting problem is

$$H_0 = \mathbf{X}_s \otimes \mathcal{M}_s \quad (5)$$

---

**Algorithm 1:** Sub-graph sampling and random mask for GNN training

---

**Input:** Historical data  $\mathbf{X}$  of the traffic network, graph structure of traffic network  $\mathcal{G}(\mathcal{V}, \mathcal{A})$ , initialized GNN model  $f(\cdot)$   
**Output:** trained GNN model  $f(\cdot)$

```

1 for  $i = 1, 2, \dots, I$  do
2   create sampl mask list  $\mathcal{M}^l = []$ , adjacent matrix list
    $\mathcal{A}^l = []$ , feature list  $\mathbf{X}^l = []$ , target list  $\mathbf{Y}^l = []$ 
3   for  $s = 1, 2, \dots, S$  do
4     Generate number  $n_r, n_m$ 
5     Sample  $n_r$  of nodes  $\mathcal{V}_r \subset \mathcal{V}_o$ ,  $n_m$  of nodes
      $\mathcal{V}_m \subset \mathcal{V}_o \setminus \mathcal{V}_r$  and  $\mathcal{V}_s = \mathcal{V}_r \cup \mathcal{V}_m$ 
6     Construct adjacent matrix  $\mathcal{A}_s$  by selecting
     entries at  $\mathcal{A}_{i,j}$ , where  $i, j \in \mathcal{V}_s \subset \mathcal{G}$ , and
     construct sub-graph  $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{A}_s)$ .
      $\mathcal{A}^l.\text{append}(\mathcal{A}_s)$ 
7     Generate mask matrix  $\mathcal{M}_s$  with 1 at row
      $i \in \mathcal{V}_r \in \mathcal{G}_s$  and 0 at row  $j \in \mathcal{V}_m \subset \mathcal{G}_s$ ,
      $\mathcal{M}^l.\text{append}(\mathcal{M}_s)$ 
8     Random sample time step  $t$  and obtain node
     features  $\mathbf{X}_s^{(t-T+1:t)}$  at  $\mathcal{V}_s$  for  $T$  historical
     length and target  $\mathbf{X}_s^{(t+\Delta t)}$ ,
      $\mathbf{X}^l.\text{append}(\mathbf{X}_s^{(t-T+1:t)})$ ,  $\mathbf{Y}^l.\text{append}(\mathbf{X}_s^{(t+\Delta t)})$ 
9   end
10  Train GNN model  $f(\cdot)$  with sampled data  $\mathcal{M}^l, \mathcal{A}^l,$ 
    $\mathbf{X}^l, \mathbf{Y}^l$ .
11 end
```

---

, where  $\otimes$  is the Hadamard product. Since  $H_0[j, :]$ , where  $j \in \mathcal{V}_s \subset \mathcal{G}_s$  is all 0 after masked by  $\mathcal{M}_s$ , we follow [9] and define the second layer as:

$$H_2 = \sigma\left(\sum_{k=1}^K T_k(\bar{\mathcal{A}}_f) H_1 \Theta_{b,1}^k + T_k(\bar{\mathcal{A}}_b) H_1 \Theta_{f,1}^k\right) + H_1 \quad (6)$$

where  $\sigma(\cdot)$  is the nonlinear activation function for the DGCN layers. We add the  $H_1$  to the  $H_2$  since it contains the information at the  $\mathcal{V}_m$  locations with no sensors. The prediction target is the output of the last layer from DGCN, changing into the right shape through a fully connected layer. During training, the forward and backward transition adjacent matrix is  $\bar{\mathcal{A}}_s$  calculated from  $\mathcal{A}_s$  for learning the graph representation and parameter  $\{\Theta_{f,l}^k, \Theta_{b,l}^k\}_{\{1,2,\dots,K\}}^{\{1,2,\dots,L\}}$ , where  $K$  is the order of Chebyshev polynomial and  $L$  is the number of layers of GNN model.

Since our task is to predict future traffic state on observable and missing locations on  $\mathcal{G}$ , the training procedure on the sampled network  $\mathcal{G}_s$  should be able to predict reserved locations observable after masking and recover locations masked from observation such that the message passing mechanism could be generalized to all nodes. As a result, the learned inductive GNN model trained on  $\mathcal{G}_s$  could be generalized to network  $\mathcal{G}$ . To achieve this, we design the prediction loss function as follows:

$$J_{pre} = \sum_{s \in S} \text{Loss}(X_s^{t+\Delta t}, \text{MLP}(f(\mathbf{X}_s^{(t-T+1:t)} \otimes \mathcal{M}_s, \mathcal{A}_s))) \quad (7)$$

where MLP means the fully connected layer for final outputs.

To keep more structure and feature information, we also minimize the difference between recovered data and initial model inputs by defining a recovery loss as follows:

$$J_{rec} = \sum_{s \in S} \text{Loss}(f(\mathbf{X}_s^{(t-T+1:t)} \otimes \mathcal{M}_s, \mathcal{A}_s), H_0) \quad (8)$$

And we can define our total loss as follows:

$$J_{total} = J_{pre} + \alpha J_{rec} \quad (9)$$

where  $\alpha$  is a hyperparameter for balancing the two loss functions.

### B. Uncertainty Aware Traffic Forecasting

In this section, we will introduce the limitations of most current models based on point estimation [10] in our Problem 2 and propose a new method that can predict the future traffic state and quantify the uncertainty at each location in the traffic network and show how the down-stream tasks can benefit from it.

1) *Model Uncertainty on Sensor Network:* The past works in traffic prediction problems usually use point estimation rather than probabilistic prediction with built-in uncertainty. However, accurate point estimations often require evaluation metrics like root mean square error (RMSE), which is not retrievable at missing locations with no historical data. As a result, model performance can only be evaluated at observable locations, and we will have no sense of the model performance at missing locations. And the resulting predictions will not be accredited in the downstream tasks that need assessing potential risks and involve decision-making.

To overcome the limitations of point estimation and assessable at missing locations, we need to quantify the uncertainty at different locations of traffic networks, especially at the extended missing ones. Specifically, we would like to output the uncertainty value along with the prediction value at each location:

$$[\mathbf{X}_o^{(t-T+1)}, \mathbf{X}_o^{(t-T+2)}, \dots, \mathbf{X}_o^{(t)}, \mathcal{G}] \xrightarrow{f(\cdot)} [\hat{\mathbf{X}}_o^{(t+\Delta t)}, \hat{\mathbf{X}}_e^{(t+\Delta t)}], \quad [\mathbf{U}_o^{(t+\Delta t)}, \mathbf{U}_e^{(t+\Delta t)}] \quad (10)$$

where the  $\mathbf{U}_o^{(t+\Delta t)}$ ,  $\mathbf{U}_e^{(t+\Delta t)}$  are uncertainty at observable and missing locations respectively.

Any uncertainty quantification method applies to our proposed framework as long as it can quantify the uncertainty of the deep neural network in regression tasks such as evidential deep regression [34] (EDL), Dropout [28] (Dropout), Quantile Regression [39] (QR), etc. In this paper, we use EDL to quantify model uncertainty and investigate the uncertainty

at both observable and missing locations, since it is most compatible with our Problem 2 setting.

Intuitively, the inductive graph neural network enables our model to predict future traffic states in both observable and missing locations; However, due to the absence of node features at missing locations, the information may be lost and thus bring in additional data uncertainty compared to observable locations with node features. Besides, due to the topological structure and distribution of node features in the neighborhood in the traffic network graph, some of the nodes' representation may not be well-represented during the training process and, as a result, have higher model uncertainty compared to other locations. And finally, the out-of-distribution nodes at missing locations will have high prediction uncertainty.

To quantify the uncertainty at each location, we incorporate the uncertainty quantification layer with the output layer and output the negative log-likelihood along with the predicted traffic states.

### C. Computational Complexity Analysis

To demonstrate the efficiency of our proposed model, we discuss the advantage of taking diffusion graph neural network (DGCN) as our model backbone here. In general, computing a convolution on the graph is expensive. However, since  $\mathcal{G}$  is sparse due to the nature of traffic network connection and use of Gaussian thresholded kernel, Equation 4 could be computed efficiently using  $\mathcal{O}(K)$  recursive sparse-dense matrix multiplication [7]. The overall total time complexity is  $\mathcal{O}(K|\mathcal{E}|) \ll \mathcal{O}(|\mathcal{V}|^2)$ , which makes our model able to be deployed on real-time hardware with a limited computational resource.

## V. EXPERIMENTAL EVALUATION

To evaluate our model's performance in the Problem 2 and justify that our new problem setting holds significant relevance and could serve as guidance for different downstream tasks, we focus on the following research questions:

**RQ1:** Comparing to other approaches applicable in our problem setting, how does our model UIGNN perform?

**RQ2:** Does our model veritably reflect the prediction uncertainty at different sensor locations?

**RQ3:** Does the model trained in our new problem setting applicable and guide the downstream task?

### A. Dataset

We conduct our experiments on three real-world, large-scale traffic datasets. **METR-LA:** This dataset contains traffic information from loop detectors located in the Los Angeles highway road net. We follow the settings of IGNNK [9] and select 207 sensors with a 6-month record in total to construct a spatiotemporal graph. **PeMS-BAY:** It's a dataset similar to METR-LA, collected from highway networks with 325 sensors. Different from METR-LA, the weighted adjacency matrix is not constructed by latitude and longitude for the sensors listed directly, but by matching absolute postmile markers to find connectivity information. **SeData:** This dataset

is collected from 323 locations with loop detectors on highway road net, following the same format as METR-LA. Its adjacency matrix is a simple binary matrix representing whether two nodes are connected or not.

In these datasets, we choose 70% of the data for training and the rest for validation and testing. During the sampling period, we randomly select a subset of sensors as unknown ones removed from the full set (the number changes when using different datasets). This removed set is defined as missing locations ( $\mathcal{V}_e$ ) without historical data available and cannot be used as the target in training. And the remained set is denoted as the observable locations ( $\mathcal{V}_o$ ) with historical data available during training and inference.

The adjacent matrix is computed with the thresholded Gaussian kernel [40], paired wisely from the distances along the road network between the two locations to construct the traffic sensor network graph. It is related to the direction of the road network. Each entry  $A_{i,j} = \exp(-\frac{dist(v_i, v_j)^2}{\sigma^2})$  and trucked to be 0 if  $dist(v_i, v_j) \leq \kappa$ , where  $\sigma$  is the standard deviation of distances and  $\kappa$  is the threshold.

## B. Experiment Setup

1) **Baselines**: Since our problem predicts future traffic state with historical data at partial locations, this is of major difference to past traffic prediction problems, and most past models cannot be directly adopted in this scenario. Thus, we take a two-step approach first, predict the historical traffic state at locations without sensors (missing locations) and then forecast the future traffic state based on the completed historical data. To make the comparison fair, we use the diffusion graph convolution neural network (DGCN) as the building block, the same as our UIGNN model. We pick 1) **Kriging + DGCN**: kriging is a spatial interpolation model [3] widely used in the geostatistical field to recover data at missing locations and then use DGCN to forecast future traffic state at time  $t + \Delta t$ . We use the ordinary Kriging method, and the correlation is solely dependent on the network road distances instead of the spatial distance since it is on a road network structure. 2) **KNN + DGCN**: K spatially nearest neighbors of the missing locations are aggregated and averaged to estimate their traffic states. And the prediction step is the same as Kriging + DGCN. 3) **ST-MVL + DGCN**: ST-MVL [41] is a spatiotemporal multiview-based learning algorithm to recover missing data of geo-sensory time series data. In our setting, since missing locations, the historical data is unavailable; we only use spatial view to complete data. The prediction step is the same as the previous methods. 4) **IGNNK + DGCN**: We first train an IGNNK model to complete data at missing locations following the original work procedure [9]. The prediction step is the same as previous methods. Notice matrix factorization-based methods do not apply to our problem since they need future traffic states at observable locations to recover the missing locations at the test dataset. Thus it is not fair to compare.

2) **Hyperparameters**: The dimension of the hidden layers of our model is set to 100. We implement all our batch sizes in our experiments set to 4. We use the Adam optimizer, set

the learning rate to  $1e-4$ , and train the UIGNN model for 750 epochs. We choose time slice  $h=24$  for METR-LA, PeMS-bay, and SeData.

3) **Metrics**: To compare our model and other baseline model performance, we use mean RMSE at 30 mins for all three traffic speed datasets and split the locations into two groups: observed and missing. To directly show our prediction and uncertainty quantification results, we report the root mean square error (RMSE) and negative log-likelihood (NLL) loss of all sensors in our experiments and split them into two groups: observable and missing locations. We report 15-, 30- and 60-min results for all three datasets to investigate their prediction performance at different time scales.

## C. Implementation of Downstream Task

To answer **RQ3**, we evaluate our model for the active sensing downstream task in a traffic sensor network [42], [43]. This task involved deploying sensors at a subset of locations and gradually adding sensors based on the model’s output at each step. Due to budget constraints and communication costs, this active sensing scenario is crucial in real-life situations, and the current traffic prediction problem has not yet been stressed.

The experiment is conducted on the METR-LA dataset and predicts the traffic state for the upcoming 30 minutes. During the initial phase, we randomly selected 50 sample locations from 207 locations in the traffic sensor network. These selected locations  $\mathcal{V}_o$  among all locations  $\mathcal{V}$  are designated as "deployed with sensors," indicating that historical data was observable. We proceeded to train our model using the available data from these 50 locations. In the active sensing phase, new sensors are available with a fixed budget (10 sensors are deployed at each step). At each stage, we sample ten new locations from  $\mathcal{V} \setminus \mathcal{V}_o$  based on the highest uncertainty indicated by our model. We then made the historical data available for these newly selected locations. Subsequently, we retrained our model using both the previously observable locations and the newly acquired data.

## VI. EXPERIMENTAL RESULTS

### A. Overall Performance of UIGNN (RQ1)

To evaluate traffic prediction performance, we separately report the results of our method and four other two-step approaches on three traffic speed datasets in terms of root mean square error (RMSE), mean absolute error (MAE). We additionally measure the coefficient of determination ( $R^2$ ) to describe the proportion of the variation in prediction errors over historical mean values.

The result in Table II shows that our UIGNN achieves the best or second best performance in all three datasets compared to all two-step-based (first impute then predict) approaches on both observable and unobservable locations. On unobservable locations, our result shows UIGNN can accurately forecast future traffic states at locations without any historical records in the sensor network. While on observable locations, it also achieves promising results without hurting its accuracy on the traditional traffic prediction problem. We also find the performance of two-step approaches largely relies on the

TABLE II

OVERALL PERFORMANCE OF UGINN AND OTHER BASELINE METHODS ON 30 MINS PREDICTION TASK. FOR MAE AND RMSE, THE LOWER, THE BETTER. FOR  $R^2$ , THE HIGHER, THE BETTER. THE **BEST** AND SECOND BEST PERFORMANCES ARE HIGHLIGHTED. '-' MEANS THE PERFORMANCE CANNOT CONVERGE.

Dataset	Location	Metric	Method				
			UGINN	Kriging+DGCN	KNN+DGCN	ST-MVL+DGCN	IGNNK+DGCN
METR-LA	Observable	MAE	<b>6.5677</b>	7.5525	<u>7.2840</u>	7.5547	7.2916
		RMSE	<b>3.9850</b>	4.1465	<u>4.0180</u>	4.3918	4.0263
		$R^2$	<b>0.9141</b>	0.8917	<u>0.8992</u>	0.8916	0.8991
	Missing	MAE	<b>10.5004</b>	16.7824	14.4969	11.8830	<u>10.8898</u>
		RMSE	<b>7.3052</b>	12.0907	10.5916	8.0088	<u>7.7083</u>
		$R^2$	<b>0.7997</b>	0.3950	0.6214	0.7456	<u>0.7865</u>
PeMS-BAY	Observable	MAE	4.8879	—	4.4959	4.5611	<b>4.4065</b>
		RMSE	<b>1.9735</b>	—	2.1616	2.2434	<u>2.1282</u>
		$R^2$	0.8062	—	0.8006	0.7948	<b>0.8085</b>
	Missing	MAE	<b>6.9762</b>	—	7.9449	17.6647	<u>6.9922</u>
		RMSE	<b>4.1054</b>	—	5.0034	9.8226	4.1738
		$R^2$	0.4518	—	0.3138	-4.5843	<b>0.4688</b>
SeData	Observable	MAE	<b>6.4863</b>	6.6840	6.7028	6.9860	<u>6.5252</u>
		RMSE	<u>3.8114</u>	3.9150	3.8350	4.2137	<b>3.7636</b>
		$R^2$	<b>0.7299</b>	0.7132	0.7116	0.6867	<u>0.7281</u>
	Missing	MAE	8.1671	10.6495	12.3653	13.0136	<b>7.8984</b>
		RMSE	<b>4.9139</b>	7.0653	9.9310	7.6135	<u>4.9779</u>
		$R^2$	<b>0.6103</b>	0.2693	0.2567	0.0911	<u>0.6002</u>

performance of the imputation model in the first step. While on missing locations, since the imputed input is not accurate and thus brings in errors in node representation learning hence weakens the performance at these locations. Since the IGNNK + DGCN method is a learning-based model and achieves higher accuracy in the data imputation step. As a result, it outperforms other baseline models. And kriging+DGCN method cannot converge since the imputation error on missing locations is too large.

### B. Uncertainty on Sensor Network (RQ2)

One of the distinct advantages of the proposed method is the capability to quantify uncertainties for its prediction. To investigate the capability of uncertainty quantification of our method at different sensor locations, we conduct experiments on three traffic datasets and report NLL and RMSE regarding model uncertainty and accuracy at each location in Table III. The results show that observable ones have higher accuracy and less model uncertainty among all three datasets than missing locations. This is consistent with the detailed result on the METR-LA dataset shown in Fig. 2 (a) and (b). Most of the missing locations (stars in the figure) have high RMSE and NLL compared to observable ones (dots in the figure). In Fig. 2 (c), we use normalized neighbor weights which reflect the distance between the location and its surrounding observable neighbors. The larger the weight, the closer this location is to its nearby observable locations; the result shows the epistemic uncertainty at observable locations is not affected by its neighboring, while at missing locations, the epistemic uncertainty negatively correlates with its neighboring weights. This is due to the node feature at available observable locations, and our GNN model can learn better node representations. And at the missing locations, the result is weakened since these

locations have no node features, and the node representations are aggregated from the neighboring locations.

From Fig. 2 (d), we found the uncertainty successfully reflects model performance at both observed and missing locations, as error steadily decreases with epistemic uncertainty decreasing.

TABLE III  
PERFORMANCE OF UGINN ON 15 MINS, 30 MINS, AND 1-HOUR PREDICTION TASKS ON THREE REAL-WORLD DATASETS. RMSE AND NLL REFLECT MODEL ACCURACY AND UNCERTAINTY, RESPECTIVELY.

	Time	RMSE		NLL	
		Observable	Missing	Observable	Missing
METR-LA	15 min	3.2874	7.1424	5.1535	8.6832
	30 min	3.9450	7.3052	5.5716	9.0121
	1 hour	5.0083	7.8019	6.3048	8.6404
PeMS-BAY	15 min	1.6456	3.9945	3.5171	6.5895
	30 min	1.9735	4.1254	3.7642	7.0601
	1 hour	2.7178	4.2870	5.1537	7.2614
SeData	15 min	3.1805	4.7916	6.2262	7.9262
	30 min	3.8114	4.9139	6.8129	8.1601
	1 hour	4.8856	5.8468	7.6020	8.5669

### C. Case Study: Active Sensing Task (RQ3)

Fig. 3 shows the performance of 1-hour traffic prediction on the METR-LA dataset under the uncertainty sampling method directed by the UGINN model (**blue lines**) and random sampling method (**orange lanes**):

- In random and uncertainty sampling methods, the RMSE decreases with the number of observable locations increasing,



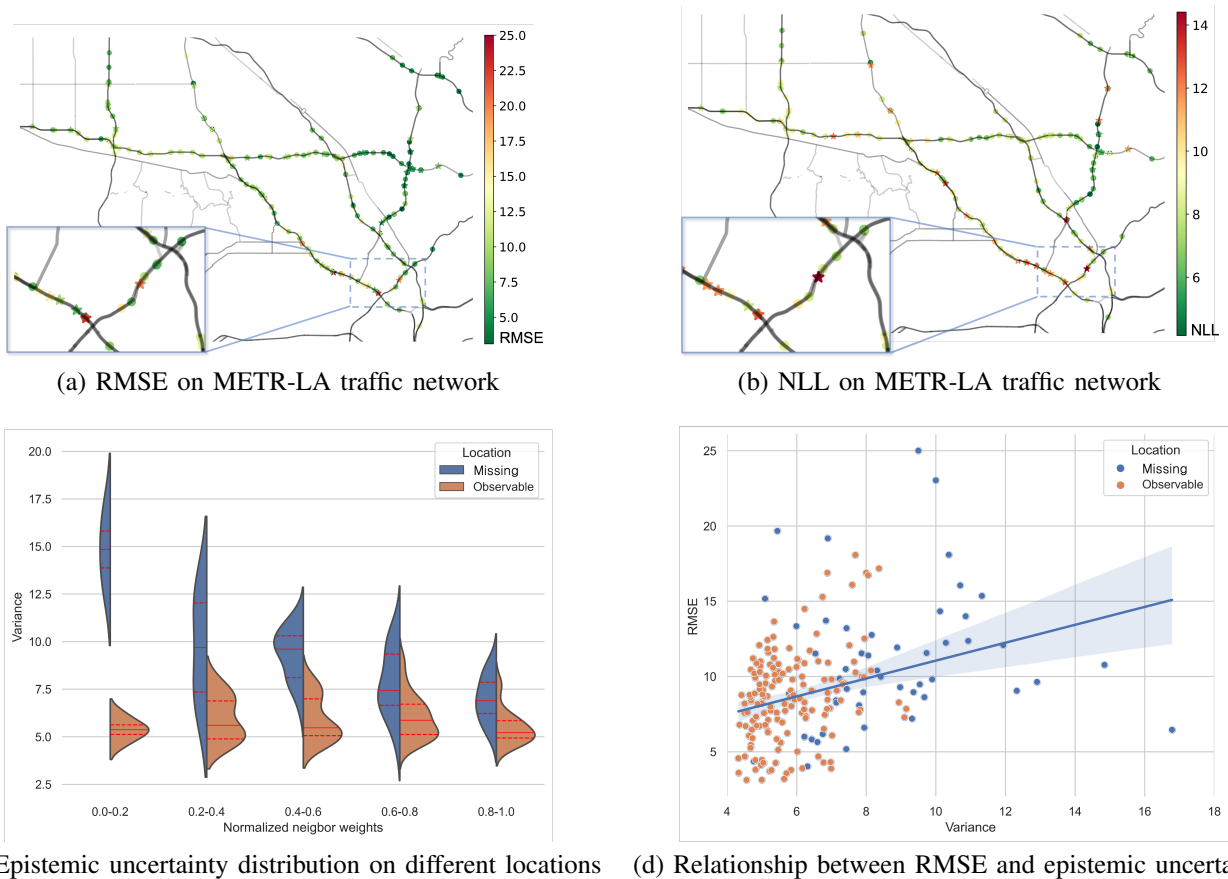


Fig. 2. Model uncertainty of 30 mins task prediction task on different locations of the METR-LA road network. (a) and (b) shows RMSE and NLL on different locations of the road network, and the stars and dots mean missing and observable locations. Compared to observable locations, missing locations have higher error and model uncertainty. (c) shows the Epistemic uncertainty of predictions at different locations. Missing locations with more and closer observable neighbors have lower model uncertainty, and observable locations do not have this effect. (d) shows that the error increases with the epistemic uncertainty increasing in both observable and missing locations.

and they finally achieve the same level of performance. This shows the performance of both methods benefits from more training data as the traffic sensing network becomes complete. And the same level of performance of different methods is also explainable as in the final few steps, the observable locations of the two methods are highly overlapping.

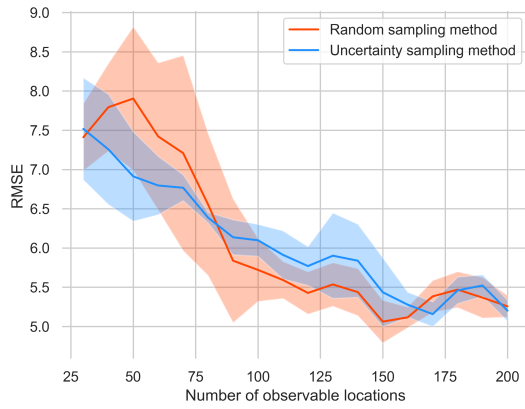
- At unobserved locations, the RMSE of the uncertainty sampling method decreases faster than the random sampling method at the beginning. This is because, at high-uncertainty locations, the topology of the locations in the sensing network and traffic states of its surrounding locations are more likely to be out-of-distribution (OOD) samples which the learned model has difficulty generalizing to [34], [37]. Picking at these locations will make our model learn from the unseen data, thus decreasing the RMSE even if it is exposed to limited locations. And random sample method cannot benefit from uncertainty quantification at unobserved locations and is therefore exposed to these OOD data later than the uncertainty sampling method. As a result, the RMSE of this method decreases slower.

- At the observable locations, the RMSE of uncertainty and random sampling method have no major differences, and both

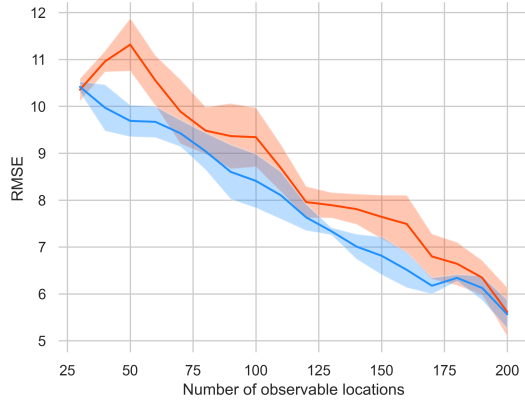
of these methods achieve the best performance at around 100 observable locations. This is because of features at each location, and few of its neighbors can approximate the future traffic state accurately; thus, when nearly half of the locations are observable, the forecasting model achieves the best performance. Since uncertainty helps the sensing network pick up the most likely unseen data at high-uncertainty locations, it does not provide much improvement at locations with data already seen.

## VII. CONCLUSION

We investigate the traffic prediction problem in a realistic setting where data is missing at part of the locations. To extend prediction perception to these locations missing historical records, we propose an uncertainty-aware inductive graph neural network that can: 1) predict future traffic state at both observable locations and missing locations; 2) quantify the uncertainty of the prediction which can help evaluate the model performance at both observable and missing locations. We conduct extensive environments using three real-world datasets and testify our method performs better than intuitive solutions based on two-steps (impute then forecast) approaches. In



(a) RMSE of  $\mathcal{V}_o$  at different  $N_e$  on METR-LA



(b) RMSE of  $\mathcal{V}_e$  at different  $N_e$  METR-LA

Fig. 3. Performance of Uncertainty sampling and random sampling method on active sensing task. Lower RMSE means better prediction performance. (a) It shows different sampling methods have similar performance on observable locations. (b) The uncertainty sampling method performs better compared to the random sampling method. And two methods finally converge and achieve the same level of accuracy.

addition, we show in the case study on the METR-LA dataset that the uncertainty can help evaluate prediction accuracy and assist downstream tasks like sensor deployment on traffic networks.

## REFERENCES

- [1] Y. Zheng, L. Capra, O. Wolfson, and H. Yang, "Urban computing: concepts, methodologies, and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 5, no. 3, pp. 1–55, 2014.
- [2] I. Lana, J. Del Ser, M. Velez, and E. I. Vlahogianni, "Road traffic forecasting: Recent advances and new challenges," *IEEE Intelligent Transportation Systems Magazine*, vol. 10, no. 2, pp. 93–109, 2018.
- [3] N. Cressie and C. K. Wikle, *Statistics for spatio-temporal data*. John Wiley & Sons, 2015.
- [4] H. V. Jagadish, J. Gehrke, A. Labrinidis, Y. Papanikolaou, J. M. Patel, R. Ramakrishnan, and C. Shahabi, "Big data and its technical challenges," *Communications of the ACM*, vol. 57, no. 7, pp. 86–94, 2014.
- [5] M. X. Hoang, Y. Zheng, and A. K. Singh, "Fccf: forecasting citywide crowd flows based on big data," in *Proceedings of the 24th ACM SIGSPATIAL international conference on advances in geographic information systems*, 2016, pp. 1–10.
- [6] X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, and Y. Wang, "Learning traffic as images: A deep convolutional neural network for large-scale transportation network speed prediction," *Sensors*, vol. 17, no. 4, p. 818, 2017.

- [7] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion convolutional recurrent neural network: Data-driven traffic forecasting," *arXiv preprint arXiv:1707.01926*, 2017.
- [8] M. T. Bahadori, Q. R. Yu, and Y. Liu, "Fast multivariate spatio-temporal analysis via low rank tensor learning," *Advances in neural information processing systems*, vol. 27, 2014.
- [9] Y. Wu, D. Zhuang, A. Labbe, and L. Sun, "Inductive graph neural networks for spatiotemporal kriging," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 5, 2021, pp. 4478–4485.
- [10] D. Wu, L. Gao, M. Chinazzi, X. Xiong, A. Vespignani, Y.-A. Ma, and R. Yu, "Quantifying uncertainty in deep spatiotemporal forecasting," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 1841–1851.
- [11] Q. Wang, S. Wang, D. Zhuang, H. Koutsopoulos, and J. Zhao, "Uncertainty quantification of spatiotemporal travel demand with probabilistic graph neural networks," *arXiv preprint arXiv:2303.04040*, 2023.
- [12] H.-F. Yu, N. Rao, and I. S. Dhillon, "Temporal regularized matrix factorization for high-dimensional time series prediction," *Advances in neural information processing systems*, vol. 29, 2016.
- [13] Y. Zhao, L. Ye, P. Pinson, Y. Tang, and P. Lu, "Correlation-constrained and sparsity-controlled vector autoregressive model for spatio-temporal wind power forecasting," *IEEE Transactions on Power Systems*, vol. 33, no. 5, pp. 5029–5040, 2018.
- [14] J. Zhang, Y. Zheng, and D. Qi, "Deep spatio-temporal residual networks for citywide crowd flows prediction," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017.
- [15] J. Ye, J. Zhao, K. Ye, and C. Xu, "How to build a graph-based deep learning architecture in traffic domain: A survey," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 5, pp. 3904–3924, 2020.
- [16] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- [17] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson, "Structured sequence modeling with graph convolutional recurrent networks," in *Neural Information Processing: 25th International Conference, ICONIP 2018, Siem Reap, Cambodia, December 13-16, 2018, Proceedings, Part I 25*. Springer, 2018, pp. 362–373.
- [18] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," in *International conference on machine learning*. PMLR, 2017, pp. 1243–1252.
- [19] Y. Liu, Z. Liu, C. Lyu, and J. Ye, "Attention-based deep ensemble net for large-scale online taxi-hailing demand prediction," *IEEE transactions on intelligent transportation systems*, vol. 21, no. 11, pp. 4798–4807, 2019.
- [20] C. F. Ansley and R. Kohn, "On the estimation of arima models with missing values," in *Time Series Analysis of Irregularly Observed Data: Proceedings of a Symposium held at Texas A & M University, College Station, Texas February 10–13, 1983*. Springer, 1984, pp. 9–37.
- [21] E. Acar, D. M. Dunlavy, T. G. Kolda, and M. Mørup, "Scalable tensor factorizations for incomplete data," *Chemometrics and Intelligent Laboratory Systems*, vol. 106, no. 1, pp. 41–56, 2011.
- [22] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [23] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," *Advances in neural information processing systems*, vol. 28, 2015.
- [24] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [25] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, and V. Prasanna, "Graphsaint: Graph sampling based inductive learning method," *arXiv preprint arXiv:1907.04931*, 2019.
- [26] M. Zhang and Y. Chen, "Inductive matrix completion based on graph neural networks," *arXiv preprint arXiv:1904.12058*, 2019.
- [27] J. Zhang, X. Shi, J. Xie, H. Ma, I. King, and D.-Y. Yeung, "Gaan: Gated attention networks for learning on large and spatiotemporal graphs," *arXiv preprint arXiv:1803.07294*, 2018.
- [28] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*. PMLR, 2016, pp. 1050–1059.
- [29] Y. Gal, J. Hron, and A. Kendall, "Concrete dropout," *Advances in neural information processing systems*, vol. 30, 2017.

- [30] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," *Advances in neural information processing systems*, vol. 30, 2017.
- [31] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural network," in *International conference on machine learning*. PMLR, 2015, pp. 1613–1622.
- [32] J. M. Hernández-Lobato and R. Adams, "Probabilistic backpropagation for scalable learning of bayesian neural networks," in *International conference on machine learning*. PMLR, 2015, pp. 1861–1869.
- [33] D. P. Kingma, T. Salimans, and M. Welling, "Variational dropout and the local reparameterization trick," *Advances in neural information processing systems*, vol. 28, 2015.
- [34] A. Amini, W. Schwarting, A. Soleimany, and D. Rus, "Deep evidential regression," *Advances in Neural Information Processing Systems*, vol. 33, pp. 14 927–14 937, 2020.
- [35] T. Pearce, A. Brintrup, M. Zaki, and A. Neely, "High-quality prediction intervals for deep learning: A distribution-free, ensembled approach," in *International conference on machine learning*. PMLR, 2018, pp. 4075–4084.
- [36] A. Alaa and M. Van Der Schaar, "Frequentist uncertainty in recurrent neural networks via blockwise influence functions," in *International Conference on Machine Learning*. PMLR, 2020, pp. 175–190.
- [37] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.
- [38] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," *Advances in neural information processing systems*, vol. 29, 2016.
- [39] R. Koenker and K. F. Hallock, "Quantile regression," *Journal of economic perspectives*, vol. 15, no. 4, pp. 143–156, 2001.
- [40] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE signal processing magazine*, vol. 30, no. 3, pp. 83–98, 2013.
- [41] X. Yi, Y. Zheng, J. Zhang, and T. Li, "St-mvl: filling missing values in geo-sensory time series data," in *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, 2016.
- [42] A. Krause, A. Singh, and C. Guestrin, "Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies," *Journal of Machine Learning Research*, vol. 9, no. 2, 2008.
- [43] C. Guestrin, A. Krause, and A. P. Singh, "Near-optimal sensor placements in gaussian processes," in *Proceedings of the 22nd international conference on Machine learning*, 2005, pp. 265–272.