# Adversarial Diffusion Attacks on Graph-based Traffic Prediction Models

Lyuyi Zhu,  Kairui Feng,  Ziyuan Pu, *Member, IEEE* Wei Ma, *Member, IEEE*

arXiv:2104.09369v1 [cs.LG] 19 Apr 2021

*Abstract*—Real-time traffic prediction models play a pivotal role in smart mobility systems and have been widely used in route guidance, emerging mobility services, and advanced traffic management systems. With the availability of massive traffic data, neural network-based deep learning methods, especially the graph convolutional networks (GCN) have demonstrated outstanding performance in mining spatio-temporal information and achieving high prediction accuracy. Recent studies reveal the vulnerability of GCN under adversarial attacks, while there is a lack of studies to understand the vulnerability issues of the GCN-based traffic prediction models. Given this, this paper proposes a new task – diffusion attack, to study the robustness of GCN-based traffic prediction models. The diffusion attack aims to select and attack a small set of nodes to degrade the performance of the entire prediction model. To conduct the diffusion attack, we propose a novel attack algorithm, which consists of two major components: 1) approximating the gradient of the black-box prediction model with Simultaneous Perturbation Stochastic Approximation (SPSA); 2) adapting the knapsack greedy algorithm to select the attack nodes. The proposed algorithm is examined with three GCN-based traffic prediction models: ST-GCN, T-GCN, and A3T-GCN on two cities. The proposed algorithm demonstrates high efficiency in the adversarial attack tasks under various scenarios, and it can still generate adversarial samples under the drop regularization such as DROPOUT, DROPNODE, and DROPEDGE. The research outcomes could help to improve the robustness of the GCN-based traffic prediction models and better protect the smart mobility systems.

*Index Terms*—Traffic Prediction, Deep Learning, Graph Convolutional Network, Adversarial Attack, Intelligent Transportation Systems

## I. INTRODUCTION

PEOPLE'S activities and movements in smart cities rely on accurate, robust, and real-time traffic information. With massive data collected in the intelligent transportation systems (ITS), various methods, such as time series models, state-space models, and deep learning, have been developed to carry out the short-term prediction for traffic operation and management [1]. Among these methods, deep learning methods, especially the graph convolutional networks (GCN), achieve state-of-the-art accuracy and are widely employed in industry-level smart mobility applications. For example, Deepmind has partnered with Google Maps to improve the accuracy of real-time Estimated Time of Arrival (ETA) prediction using GCN [2].

The predicted traffic information plays a critical role in our daily traveling, and travelers take for granted that the predicted results are accurate and trustworthy. However, the robustness and vulnerability issues of these deep learning models have not been investigated for traffic prediction models. Recent studies have shown that neural networks are vulnerable to deliberately designed samples, which are known as adversarial samples. In general, the adversarial samples could be generated by adding imperceptible perturbations to the original data sample. Though the adversarial sample is very similar to its original counterpart, it can significantly change the performance of the deep learning models. Szegedy et al. (2013) firstly discovered this phenomenon on deep neural networks (DNN), and they found that adversarial samples are low-probability but densely distributed [3]. Goodfellow et al. [4] also showed neural networks are vulnerable to the adversarial samples in the sense that it is sufficient to generate adversarial samples when DNNs demonstrate linear behaviors in high-dimensional spaces.

Due to the existence of the adversarial samples, potential attackers could take advantage of the deep learning models and degrade the model performance. Though related theories and applications have been studied in various areas such as computer vision [5], social networks [6], traffic signs [7] and recommendation systems [8], few of the studies have investigated the vulnerability and robustness of the traffic prediction systems. It has been shown that industry-level traffic information systems can be "attacked" easily. Recently, a German artist walked slowly with a handcart, which was loaded with 99 smartphones. On each smartphone, the mobile application Google Maps was turned on. The 99 cell phones virtually created 99 vehicles on the roads, and all the "vehicles" were slowly moving along the road. As Google Maps estimated and predicted the traffic states based on the data sent back from those cell phones, it wrongly identified an empty street (green) to be a congested road (red), as shown in Fig. 1. Though it is unclear which model Google Maps is using, this experiment indeed reveals the possibility of adversarial attacks on real-world and industry-level traffic information systems.

Adversarial attacks on traffic prediction models can affect every aspect of the smart mobility systems. We summarize the following four scenarios in which adversarial attacks can significantly degrade the performance of the systems.

- **Smartphone-based mobility applications.** Mobile phone-based mapping services such as Google Maps and AutoNavi make traffic state estimation and prediction based on the GPS trajectories sent from their users [10].

L. Zhu is with the College of Civil Engineering and Architecture, Zhejiang University, Hangzhou, China (E-mail: 3170103586@zju.edu.cn).

K. Feng is with the Department of Civil and Environmental Engineering, Princeton University, NJ, U.S.A (E-mail: kairuif@princeton.edu)

Z. Pu is with the School of Engineering, Monash University, Jalan Lagoon Selatan, 47500 Bandar Sunway, Malaysia (Email: ziyuan.pu@monash.edu).

W. Ma is with the Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University, Hong Kong SAR, China (E-mail: wei.w.ma@polyu.edu.hk).
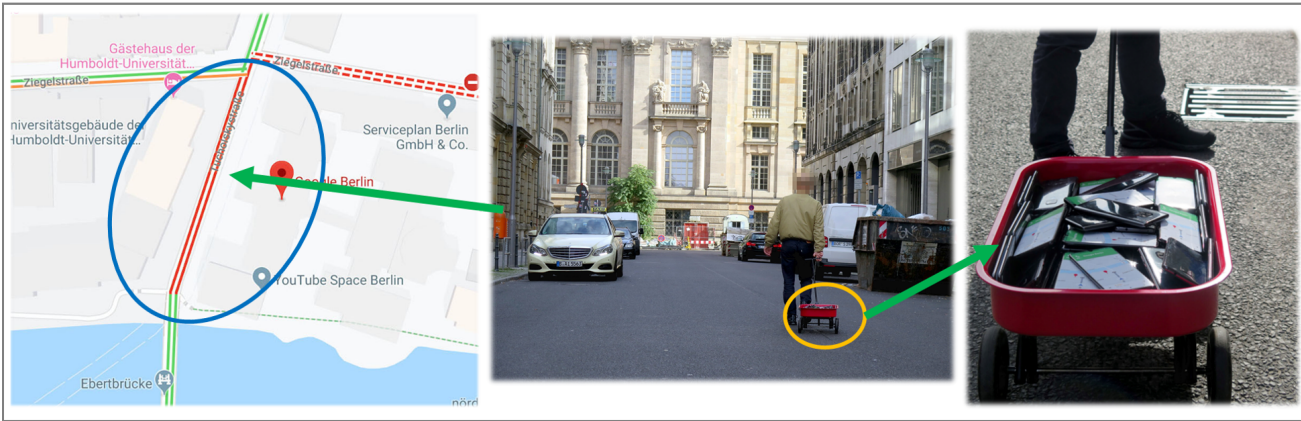
Manuscript received: April 20, 2021.

Fig. 1: Google Maps Hacks using 99 cell phones [9].

However, users' mobile phones can be hacked and the information can be deliberately altered to attack the systems [11]. In general, these individual mobility applications are vulnerable to adversarial attacks because the difficulties of hacking users' mobile phones are way lower than hacking the application server.

- **Connected vehicle (CV) systems.** In the CV systems, traffic information is collected by the roadside units (RSUs) and sent to the traffic control center [12]. CVs use information from either RSU or control center to plan routes and avoid congestion. However, it is possible to hack some of the RSUs to send the adversarial samples to manipulate the predicted traffic information from the control center. Attackers could make use of adversarial attacks to benefit a group of vehicles while causing unexpected delays to other vehicles.

- **Emerging mobility services.** Transportation network companies (TNCs) depend on accurate traffic speed prediction for vehicle dispatching, routing, and relocation on the central platform. It is possible that a group of vehicles collide with each other and falsely report their GPS trajectories and speed to confuse the central platform. By carefully designing the adversarial samples with purpose, this group of vehicles could take advantage of the central platform by receiving more orders and running on less congested roads.

- **Advanced Traffic Management Systems (ATMS).** Most of the network-wide transportation management systems [13] rely on user equilibrium (or stochastic user equilibrium) models to depict and predict travelers' behaviors and both models assume that travelers can acquire accurate (or nearly accurate) traffic information. Under adversarial attacks, this assumption no longer holds. The inaccurately predicted traffic information can reduce the effectiveness of the ATMS and degrade the efficiency of the entire network.

Adversarial attack for traffic prediction systems is a unique task that is different from existing literature. In this paper, we focus on the GCN-based traffic prediction models. Previous literature aims at attacking one node by modifying the features of all the nodes in the GCN-based neural networks, while

this is impractical on traffic prediction models. On real-world traffic networks, modifying the features on all nodes is challenging and costly and our objective is to degrade the network-wide system performance instead of a single node. A practical task is to degrade the overall system performance by perturbing the features on a small subset of nodes, which can be viewed as the opposite of the conventional adversarial attack problems. To this end, we propose a novel concept – diffusion attack, and its definition is presented in Definition 1.

**Definition 1** (Diffusion Attack). *Considering a GCN-based deep neural network, the input features are associated with each node of the corresponding graph. The diffusion attack aims to modify the input features on a limited set of nodes while keeping the graph topology unchanged, and the goal is to degrade the overall performance of the neural network on all the nodes.*

An illustration of the diffusion attack is presented in Fig. 2. On the left-hand side, the prediction model performs normally and can generate accurate predictions. After the diffusion attack, two nodes in red are attacked, and their nearby neighbors are strongly perturbed, followed by their 2-hop neighbors being slightly perturbed. One can see that the attack effects diffuse from the node being attacked to its neighbors, and later we will develop mathematical proofs and numerical experiments to verify this phenomenon. The purpose of attackers is to select the optimal attack nodes and to generate adversarial samples to maximize attack effects.
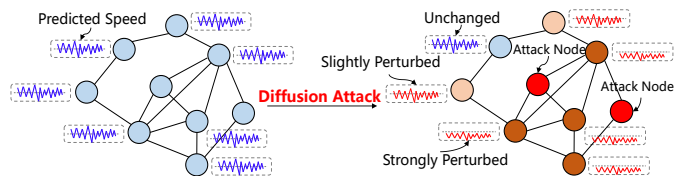


Fig. 2: An illustration of the diffusion attack.

To summarize, vulnerability issues of the traffic prediction models are critical to smart mobility systems, while the related research is still lacking. Given this, we explore the robustness of the GCN-based traffic prediction models. Based on the

unique characteristics of the adversarial attacks on traffic prediction models, we propose the concept of diffusion attack, which aims to attack a small set of nodes to degrade the performance of the entire network. On top of that, we develop a diffusion attack algorithm, which consists of two major components: 1) approximating the gradient of the black-box prediction model with Simultaneous Perturbation Stochastic Approximation (SPSA); 2) adapting the knapsack greedy algorithm to select nodes to attack. The proposed algorithm is examined with three GCN-based traffic prediction models: ST-GCN, T-GCN, and A3T-GCN on two cities: Los Angeles and Hong Kong. The proposed algorithm demonstrates high efficiency in the diffusion attack tasks under various scenarios, and the algorithm can still generate adversarial samples under different drop regularization. We further discuss how to improve the robustness of the GCN-based traffic prediction models, and the research outcomes could help to better protect the smart mobility systems in both the cyber and physical world. The contributions of this study are summarized as follows:

- Different from existing adversarial attack tasks on graphs, we propose a novel task of diffusion attack, which aims to select and attack a small set of nodes to degrade the performance of the entire traffic prediction models. This task is suitable for the traffic prediction context, while it is overlooked in the existing literature.
- To generate the adversarial samples on traffic prediction models, we propose the Simultaneous Perturbation Stochastic Approximation (SPSA) algorithm to efficiently approximate gradients of the black-box prediction models.
- To select the optimal attack nodes, we formulate the diffusion problem as a knapsack problem and then adapt the greedy algorithm to determine the priority of the attack nodes.
- We conduct extensive experiments to attack the widely-used traffic prediction models on Los Angeles and Hong Kong. Different drop regularization strategies (*e.g.*, DROPOUT, DROPEDGE, DROPNODE) for defending the adversarial attacks are also tested to ensure that the proposed algorithm can still generate effective adversarial samples in various scenarios.

The remainder of this paper is organized as follows. Section II reviews the related studies on traffic prediction, GCN models, and adversarial attacks on graphs. Section III rigorously formulates the diffusion attack problem and presents the developed attack algorithms. In Sections IV, three traffic prediction models and two real-world datasets are used to examine the proposed attack algorithms. Finally, conclusions and future research are summarized in Section V.

## II. RELATED WORKS

In this section, we first overview the traffic prediction tasks and then summarize recent studies on GCN-based traffic prediction models. Robustness issues of the GCN-based prediction models under adversarial attacks are also discussed.

### A. Traffic Prediction

The traffic prediction problem has been extensively studied for decades, and various statistical models have been developed to solve the problem, including History Average (HA) [14], Autoregressive Integrated Moving Average (ARIMA) [15]–[19], Support Vector Regression (SVR) [20], clustering [21], and Kalman filtering [22], [23]. In recent years, the data scale becomes large and the spatio-temporal correlation of the data becomes complicated, and hence traditional statistical methods reveal their limits in face of the massive and complex data. Instead, neural network models demonstrate potentials in traffic prediction with multi-source data on large-scale networks. Various neural network models have been used for traffic prediction, including Convolutional Neural Network (CNN) [24], Recurrent Neural Network (RNN) [25]–[27], attention [28], [29] and Graph Convolutional Network (GCN) [30]–[33].

Traffic prediction tasks can also be categorized into multiple purposes, such as traffic state prediction, demand prediction, and trajectory prediction [34]. Traffic state prediction includes the prediction of traffic flow [28], speed [30], and travel time [35]. Traffic demand prediction aims to make prediction of the number of users and traffic demand, such as taxi request [36], subway inflow/outflow [37], bike-sharing demand [38], [39], and origin-destination demand [40]. It is also possible to predict the trajectories of travelers and vehicles, and this task is used for dynamic positioning and resource allocation [41], [42]. Overall, most of the traffic prediction tasks can be carried out by neural network models, and hence it is crucial to study their vulnerability issues.

### B. GCN and its Applications on Traffic Prediction

Traffic data is closely associated with the topological structure of the road networks, and hence it is typical graph-based data. The graph-based data is represented in the non-Euclidean space, and conventional machine learning methods (*e.g.,* multi-layer perceptron) overlooks the graph-based inter-relationship among data [43]. In this paper, we summarize that traffic data consists of the following two types of information:

- **Spatial Information.** Traffic-related data can be presented on a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, 2, \cdots, N\}$ represents a set of nodes with $N = |\mathcal{V}|$, and $\mathcal{E}$ denotes a set of edges. In traffic prediction, one way to construct the graph is to make each node $v_i$ represent a road segment, and each edge represents the connectivity relationship between the road segments [30]. We further define the adjacency matrix $\boldsymbol{A} \in \{0, 1\}^{N \times N}$ on the graph $\mathcal{G}$, where $A_{ij} = 1$ when node $i$ connects node $j$, and $0$ otherwise.
- **Temporal Information.** Traffic data, such as traffic speed, density, and flow, on each node can be viewed as a time series. The traffic data on the graph is represented by a feature matrix $\boldsymbol{X} \in \mathbb{R}^{N \times S}$, where $S$ denotes the number of time intervals in the study period.

Graph Convolutional Network (GCN) demonstrates great capacities in learning graph-based information for various applications. Short-term traffic prediction is one of the most important and practical applications. The GCN can be used to

extract the spatial (non-Euclidean) relationship among nodes [30], and it can couple with recurrent neural networks (RNN) to learn the spatio-temporal patterns of the traffic data. For example, Long Short-Term Memory (LSTM) is widely used with GCN for traffic prediction [44], and the Gated Recurrent Unit (GRU) is also adopted to model the time-series on each node [30]. Some emerging models for time-series modeling, such as gated CNN and attention, can also be incorporated into the GCN-based deep learning framework [28], [31]–[33]. Readers are referred to [34] for a comprehensive review of the GCN-based traffic prediction methods.

As the backbone of a GCN-based model, the graph convolutional layer is defined as follows:

$$\boldsymbol{H}^{(l+1)} = \sigma(\hat{\boldsymbol{A}}\boldsymbol{H}^{(l)}\boldsymbol{W}^{(l)}),$$

where $l$ is the layer index, $\boldsymbol{H}^{(0)} = \boldsymbol{X}$, and $\hat{\boldsymbol{A}} = \widetilde{\boldsymbol{D}}^{-\frac{1}{2}}\widetilde{\boldsymbol{A}}\widetilde{\boldsymbol{D}}^{-\frac{1}{2}}$ is the Laplacian matrix. $\widetilde{\boldsymbol{D}}$ is the corresponding degree matrix $\widetilde{\boldsymbol{D}}_{ii} = \sum_j \boldsymbol{A}_{ij}$, and $\widetilde{\boldsymbol{A}} = \boldsymbol{A} + \boldsymbol{I}_N$, where $\boldsymbol{I}_N$ is a $N \times N$ identity matrix. $\sigma(\cdot)$ represents the activation function and $\boldsymbol{W}^{(l-1)}$ are parameters of the $l$th layer.

A $L$-layers GCN model can be expressed as following:

$$\boldsymbol{Y} = f(\boldsymbol{X}, \boldsymbol{A}) = g(\hat{\boldsymbol{A}}\cdots\sigma(\hat{\boldsymbol{A}}\boldsymbol{X}\boldsymbol{W}^{(0)})\cdots\boldsymbol{W}^{(L-1)}), \quad (1)$$

where $g(\cdot)$ is a generalized function, $\boldsymbol{Y} \in \mathbb{R}^{N\times T}$ and the parameters $\boldsymbol{W}^{(l-1)}$ for each layer could be learned by minimizing the loss between the estimated $\boldsymbol{Y}$ and true $\boldsymbol{Y}_{true}$, represented by $\mathcal{L}(\boldsymbol{Y}, \boldsymbol{Y}_{true})$. This paper will adopt the GCN as the backbone model and study its robustness and vulnerability issues under adversarial attacks.

### C. Adversarial Attacks On Graphs

Like other neural networks, GCN is vulnerable to adversarial attacks. Various attack concepts on graph have proposed, such as targeted and non-targeted attacks, structure and feature attacks, poisoning and evasion attacks. Targeted attacks aim to attack a target node [45], while the non-targeted attacks aim to compromise global performance of a model [46]–[49]. Structure attacks modify the structure of graph, such as adding or deleting nodes/edges [47], [50], [51], and feature attacks perturb the labels/features of nodes without changing the connectivity of graph [46], [52], [53]. Poisoning attacks modify the training data [46], and evasion attacks insert an adversarial samples when using the models [47], [48], [52], [54]. To be specific, attacks on traffic prediction systems are non-targeted, feature, and evasion attacks, which have not been well studied in the existing literature.

The attack algorithms can be further categorized into white-box attacks and black-box attacks. In white-box attacks, the neural network structure and parameters, training methods, and training samples are exposed to attackers. Attackers could utilize the neural network model to generate adversarial samples. Many white-box methods achieve great performance, such as fast gradient sign method (FGSM) [4], Jacobian-based saliency map approach (JSMA) [55], Carlini and Wagner Attacks (CW) [56], and Deepfool [57]. In contrast, in black-box attacks, attackers know little about internal information of the target neural network model, especially the model structure and

parameters, and only the model input and output are exposed to the attackers. For example, if attackers plan to attack the traffic prediction system, he/she is unlikely to know the internal structure or information of the prediction model.

Though it is more common in the real world, the black-box attack is much more challenging than the white-box attack. Black-box attacks can be achieved through response surface models [58] and meta-heuristic. For instance, one-pixel attack [5] uses Differential Evolution (DE) algorithm and decision-based attack [59] uses Covariance Matrix Adaptation Evolutionary Strategies (CMA-ES) to generate and improve the adversarial samples by iteration. It is also possible to approximate the gradient of the target models, and the representative models include Zeroth Order Optimization (ZOO) [60], Autoencoder-based Zeroth Order Optimization Method (AutoZOOM) [61], and Natural Evolutionary Strategies (NES) [62]. Besides, the semi-black-box attacks are in between, and it is assumed that the information of the prediction models is partially observed [63]. The above attack algorithms mainly focus on the classification task, while attack methods for traffic prediction models (*i.e.* regression task) are still lacking.

### III. PROPOSED WORKS

In this section, we first present the general formulation of the adversarial attack on graphs. Then, the new concept of diffusion attack is proposed for traffic prediction models. Lastly, we propose the new formulation and algorithm to construct the diffusion attack and discuss its implementations.

### A. Preliminaries

Traffic prediction on graphs can be regarded as a graph-based regression problem [34]. Using $f : \mathbb{R}^{N\times S} \to \mathbb{R}^{N\times T}$ in Equation 1 as a regression model (*e.g.,* a traffic prediction model) on graph $\mathcal{G}$, and $f$ contains a $L$-layer GCN, where $S$ represents the look-back time interval, and $T$ is the number of time intervals to predict. The feature matrix $\boldsymbol{X}$ contains the historical traffic states, and $\boldsymbol{Y}$ is the traffic states we want to predict. For traffic prediction problems, $\boldsymbol{Y}$ mainly depends on $\boldsymbol{X}$ as most of $\boldsymbol{A}$ are fixed [34]. We further denote $\boldsymbol{X} = (\boldsymbol{x}_1;\cdots;\boldsymbol{x}_N)$, $\boldsymbol{Y} = f(\boldsymbol{X}) = (\boldsymbol{y}_1;\cdots;\boldsymbol{y}_N)$, and $\boldsymbol{x}_i$ and $\boldsymbol{y}_i$ are the $i$th row of $\boldsymbol{X}$ and $\boldsymbol{Y}$, respectively. For node $i$, $\boldsymbol{x}_i$ is the $i$th row of the feature matrix $\boldsymbol{X}$, and $\boldsymbol{x}_i$ represents a time series of speed on node $i$. The prediction model for node $i$ can be written as $\boldsymbol{x}_i = (x_{i1}, x_{i2}, \cdots, x_{iS}) \mapsto \boldsymbol{y}_i = (y_{i1}, y_{i2}, \cdots, y_{iT})$, where $x_{i\cdot}$ is the historical traffic states, and $y_{i\cdot}$ is the future traffic states on node $i$.

### B. Diffusion Attacks on Traffic State Forecasting Models

This paper aims to attack the traffic prediction system by adding perturbations on the feature matrix $\boldsymbol{X}$, to maximally change the prediction results over a selected set of nodes. As discussed above, the graph structure is fixed and cannot be changed easily for the problem of traffic prediction, so we assume that $\boldsymbol{A}$ is fixed throughout the paper.

We construct an adversarial sample by adding the perturbation $\boldsymbol{U}$, which is the same size as $\boldsymbol{X}$, to the original input

feature matrix $\boldsymbol{X}$, as represented by $\boldsymbol{X}' = \boldsymbol{X} + \boldsymbol{U}$. Consequently, the corresponding output changes from $\boldsymbol{Y} = f(\boldsymbol{X})$ to $\boldsymbol{Y}' = f(\boldsymbol{X} + \boldsymbol{U})$.

We suppose that attackers select a set of nodes $\mathcal{P} \subseteq \mathcal{V}$ to attack. For each node $i \in \mathcal{P}$, $\boldsymbol{u}_i \begin{cases} \neq \boldsymbol{0} & i \in \mathcal{P} \\ = \boldsymbol{0} & i \notin \mathcal{P} \end{cases}$, where $\boldsymbol{u}_i$ is the $i$th row of $\boldsymbol{U}$. Then adversarial sample can be expressed as follows:

$$\boldsymbol{X}' = \boldsymbol{X} + \boldsymbol{U} = (\boldsymbol{x}_1 + \boldsymbol{u}_1; \cdots ; \boldsymbol{x}_i + \boldsymbol{u}_i; \cdots ; \boldsymbol{x}_N + \boldsymbol{u}_N).$$

By perturbing the node $i$, we change the original prediction result on node $i$ (denoted as $\boldsymbol{y}_i$) to $\boldsymbol{y}_i'$. The attack influence function $\phi_i(\boldsymbol{U})$ on node $i$ is defined as follows:

$$\phi_i(\boldsymbol{U}) = \mathcal{L}\left(\boldsymbol{y}_i'(\boldsymbol{X}'), \boldsymbol{y}_i(\boldsymbol{X})\right), \tag{2}$$

where we use $\boldsymbol{y}_i(\boldsymbol{X})$ to indicate that $\boldsymbol{y}_i$ is a function of $\boldsymbol{X}$. $\mathcal{L}(\cdot, \cdot)$ represents loss function between $\boldsymbol{y}_i'$ and $\boldsymbol{y}_i$. Here the attack influence evaluates the difference between original prediction (instead of true speed) and perturbed speed. The reason is that we assume the prediction is accurate enough, otherwise there is no need to attack. On the other hand, we could never know the true traffic condition in the future, so it is impossible to perturb the prediction against true values.

To mathematically characterize the diffusion phenomenon when attacking the GCN, we demonstrate that the following Proposition 1 holds.

**Proposition 1.** *Using the L-layer GCN model presented in Equation 1 for traffic prediction, the effect of perturbation $\boldsymbol{U}$ on each node $i$, which is denoted as $\phi_i(\boldsymbol{U})$, depends on the perturbations of its L-hop neighbors.*

*Proof.* Using $|\cdot|$ to represent the element-wise absolute value operator, and assuming that $g(\cdot)$ is Lipschitz continuous with constant $M$, $[\cdot]_i$ is the $i$th row of a matrix, $\sigma(\cdot)$ is the Relu function, and $\mathcal{L}$ represents the Mean Squared Error (MSE), we have

$$
\begin{aligned}
\phi_i(\boldsymbol{U}) &= \left\| \boldsymbol{y}_i'(\boldsymbol{X}') - \boldsymbol{y}_i(\boldsymbol{X}) \right\|_2^2 \\
&= \left\| \left[ g\left(\hat{\boldsymbol{A}} \boldsymbol{H}'^{(L-1)} \boldsymbol{W}^{(L-1)}\right) - g\left(\hat{\boldsymbol{A}} \boldsymbol{H}^{(L)} \boldsymbol{W}^{(L-1)}\right) \right]_i \right\|_2^2 \\
&\leq M \left\| \left[ \hat{\boldsymbol{A}}(\boldsymbol{H}'^{(L-1)} - \boldsymbol{H}^{(L-1)}) \boldsymbol{W}^{(L-1)} \right]_i \right\|_2^2 \\
&\leq M \left\| \left[ \hat{\boldsymbol{A}}\left(|\boldsymbol{H}'^{(L-1)} - \boldsymbol{H}^{(L-1)}|\right) |\boldsymbol{W}^{(L-1)}| \right]_i \right\|_2^2 \\
&\leq M \left\| \left[ \hat{\boldsymbol{A}}\left(|\boldsymbol{H}'^{(L-1)} - \boldsymbol{H}^{(L-1)}|\right) \right]_i \right\|_2^2 \left\| |\boldsymbol{W}^{(L-1)}| \right\|_2^2 \\
&\leq MW \left\| \left[ \hat{\boldsymbol{A}}\left(|\boldsymbol{H}'^{(L-1)} - \boldsymbol{H}^{(L-1)}|\right) \right]_i \right\|_2^2 \\
&\leq MW^2 \left\| \left[ \hat{\boldsymbol{A}}^2 \left(|\boldsymbol{H}'^{(L-2)} - \boldsymbol{H}^{(L-2)}|\right) \right]_i \right\|_2^2 \\
&\leq MW^L \left\| \left[ \hat{\boldsymbol{A}}^L \left(|\boldsymbol{H}'^{(0)} - \boldsymbol{H}^{(0)}|\right) \right]_i \right\|_2^2 \\
&= MW^L \left\| \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}| \right\|_2^2,
\end{aligned}
$$

where $\left\| |\boldsymbol{W}^{(l)}| \right\|_2^2 \leq W$ and $\hat{\boldsymbol{A}}_i^L$ represents the $i$th row of $\hat{\boldsymbol{A}}^L$. Here we can see that the upper bound of attack influence $\phi_i(\boldsymbol{U})$ associates closely with $\left\| \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}| \right\|_2^2$. Denoting $|\boldsymbol{U}|_{\cdot k}$ as the $k$th column of $|\boldsymbol{U}|$, we can expand $\left\| \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}| \right\|_2^2$ as follows:

$$
\begin{aligned}
\left\| \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}| \right\|_2^2 &= \left\| \left[ \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}|_{\cdot 1}, \cdots, \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}|_{\cdot k}, \cdots, \hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}|_{\cdot S} \right] \right\|_2^2 \\
&= \sum_{k=1}^{S} (\hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}|_{\cdot k})^2
\end{aligned}
$$

For each $k$ we can further expand it as the sum of perturbation on different nodes:

$$\hat{\boldsymbol{A}}_{i\cdot}^L |\boldsymbol{U}|_{\cdot k} = \sum_{h=1}^{N} \hat{\boldsymbol{A}}_{ih}^L |\boldsymbol{U}|_{hk},$$

where $|\boldsymbol{U}|_{hk}$ is the absolute value of perturbation added on $k$th historical traffic state of node $h$. $\hat{\boldsymbol{A}}_{ih}^L$ represents normalized connectivity weight between node $i$ and node $h$ (similar to $\boldsymbol{A}_{ih}^L$, which represents number of $L$-hop paths between $i$ and $h$, according to the graph theory). If $\hat{\boldsymbol{A}}_{ih}^L$ is large, node $h$ will have more impact on $\phi_i(\boldsymbol{U})$. Specially, if node $h$ is out of $L$-hop neighbors of node $i$, then $\hat{\boldsymbol{A}}_{ih}^L \equiv 0$, which means that attack effect will not diffuse to node $i$ when attacking $h$. For most of GCN-based traffic prediction models $L \leq 3$, so the effect of $\boldsymbol{U}$ only diffuses to its local neighbors.

$\square$

As discussed in the previous section, this paper focuses on the novel concept of diffusion attack, which aims at changing the network-wide prediction results by perturbing a small subset of node features. Mathematically, the diffusion attack problem can be expressed as to find the optimal $\mathcal{P}$ and the corresponding perturbation $\boldsymbol{U}$ such that the following influence function $\Phi(\boldsymbol{U})$ is maximized:

$$
\begin{aligned}
\max_{\boldsymbol{U}, \boldsymbol{z}} \quad & \Phi(\boldsymbol{U}) = \sum_{i \in \mathcal{V}} w_i \phi_i(\boldsymbol{U}) \\
s.t. \quad & -\varepsilon^- z_i \boldsymbol{x}_i \leq \boldsymbol{u}_i \leq \varepsilon^+ z_i \boldsymbol{x}_i \quad \forall i \in \mathcal{V} \\
& \sum_{i \in \mathcal{V}} b_i z_i \leq B \\
& z_i \in \{0, 1\} \quad\quad\quad\quad \forall i \in \mathcal{V}
\end{aligned}
\tag{3}
$$

where $z_i$ indicates whether node $i$ is attacked. To be precise, $z_i = \begin{cases} 1 & i \in \mathcal{P} \\ 0 & i \notin \mathcal{P} \end{cases}$. $b_i$ represents the cost of attacking node $i$, and $B$ is the total budget. The objective function $\Phi(\boldsymbol{U})$ represents attack influence function for the entire network, and $w_i$ is a pre-determined importance weight of node $i$. $\varepsilon^-, \varepsilon^+ > 0$ are used to control the scale of the perturbations.

Formulation 3 is a mixed integer programming (MIP), and it contains two components: 1) determining $\boldsymbol{U}$ given a fixed $\mathcal{P}$; 2) determining $\mathcal{P}$. The two components will be discussed in section III-C and III-D, respectively.

### C. Black-box attacks using SPSA

Given a fixed $\mathcal{P}$, Equation 3 reduces to a continuous optimization problem with affine constraints, as shown in the following equation:

$$
\begin{aligned}
\max_{\boldsymbol{U}} \quad & \Phi(\boldsymbol{U}) = \sum_{i \in \mathcal{V}} w_i \phi_i(\boldsymbol{U}) \\
s.t. \quad & -\varepsilon^- \boldsymbol{x}_i \leq \boldsymbol{u}_i \leq \varepsilon^+ \boldsymbol{x}_i \quad \forall i \in \mathcal{P}
\end{aligned}
\tag{4}
$$

Formulation 4 suggests that an ideal perturbation $\boldsymbol{U}^*$ should be small, and meanwhile it maximizes the overall influence function $\Phi(\boldsymbol{U}^*)$. In real-world applications, the internal information about the traffic prediction model is opaque, hence it is proper to consider the traffic prediction model as a black-box. We assume both input feature $\boldsymbol{X}$ and prediction results $\boldsymbol{Y}$ are known to attackers as both matrices represent the true and predicted traffic states in real-world, and hence Formulation 4 can be viewed as a black-box optimization problem. To solve it, we adopt the Simultaneous Perturbation

Stochastic Approximation (SPSA) method, which is featured with its efficiency and scalability [64]. In recent years, SPSA has been used for adversarial attacks on classification problems [65], while it has not been used for regression problems, in particular, the traffic prediction problem.

The SPSA method uses finite differences between two randomly perturbed inputs to approximate the gradient of the objective function. Mathematically, the gradient of $\Phi$ can be calculated as follows:

$$\widehat{\nabla\Phi}(\boldsymbol{U}_n) = \frac{\Phi(\boldsymbol{U}_n + c_n\boldsymbol{\Delta}_n) - \Phi(\boldsymbol{U}_n - c_n\boldsymbol{\Delta}_n)}{2c_n\boldsymbol{\Delta}_n}, \quad (5)$$

where $n$ is the index of the iteration, $\boldsymbol{\Delta}_n$ is a random perturbation vector whose elements are sampled from Rademacher distribution (Bernoulli $\pm 1$ distribution with probability $p = 0.5$, and we denote $\mathrm{RAD}_{1\times S} \in \{-1, 1\}^{1\times S}$ as a sample vector that follows Rademacher distribution). We further denote sequences $\{c_n\}$ and $\{a_n\}$ as follows:

$$a_n = \frac{a}{(\eta + n)^\alpha} \qquad c_n = \frac{c}{n^\gamma}, \quad (6)$$

where $a, c, \alpha, \gamma$ are hyper-parameters for SPSA [64]. Both sequences decrease when the iteration $n$ increases. Then the gradient ascent approach is utilized to maximize $\Phi(\boldsymbol{U})$, as shown in the following equation:

$$\boldsymbol{U}_{n+1} = \boldsymbol{U}_n + a_n\widehat{\nabla\Phi}(\boldsymbol{U}_n) \quad \forall n. \quad (7)$$

To summarized, the adversarial attack with fixed node set $\mathcal{P}$ is presented in Algorithm 1.

### D. Node Selection using Knapsack Greedy

This section focuses on determining the attack set $\mathcal{P}$ with a limited budget $B$. The cost $b_i$ is different across different nodes due to the level of difficulty in attacking the node. For example, urban roads may contain a more recent and secured information collection system ($b_i$ is high), while rural roads can be attacked easily ($b_i$ is low). In contrast, attacks on urban roads usually generate a higher impact on the traffic prediction methods because the urban traffic volumes are high. It can be seen that there is a trade-off between the cost and benefit when selecting the attack set $\mathcal{P}$.

We review the formulation of the diffusion attack in Equation 3, and it is actually similar to the 0-1 knapsack problem, except for that the utility of each node $i$ ($\phi_i(\boldsymbol{U})$) is unknown [66]. The attack set $\mathcal{P}$ can be viewed as a knapsack with maximum capacity $B$, and the nodes are items with their weight $w_i$. Node $i$ is added to $\mathcal{P}$ if $z_i = 1$. Due to the nature of integer programming, there is no provably efficient method to solve formulation 3, which is a NP-hard problem. Real-world networks contain hundreds or thousands of nodes, and hence it is impractical to enumerate all possible integer solutions. In this paper, we develop a family of Knapsack Greedy (KG) algorithms to solve for formulation 3, and those algorithms are inspired by the original greedy algorithm for the knapsack problem.

A trivial but insightful observation is that any perturbation $\boldsymbol{U}$ could reduce the performance of the prediction model. Proposition 2 shows that the convexity exists if we only perturb

---

**Algorithm 1** Determine the adversarial sample $\boldsymbol{X}'$ and optimal perturbation $\boldsymbol{U}$ given a fixed $\mathcal{P}$

---

**Input:** Traffic prediction model $f(\boldsymbol{X})$, input feature matrix $\boldsymbol{X}$, attack set $\mathcal{P}$, maximum iteration $\texttt{MaxIter}$.
**Output:** Adversarial sample $\boldsymbol{X}'$, and optimal perturbation $\boldsymbol{U}$.
  Initialize $\boldsymbol{U}_1 \in \mathbb{R}^{N\times S}$
  **for** $n = 1, 2, \cdots, \texttt{MaxIter}$ **do**
    Update $a_n$ and $c_n$ based on Equation 6.
    For $i \in \mathcal{V}$, random sample $\boldsymbol{\delta}_i = \begin{cases} \boldsymbol{0}_{1\times S} & \text{if } i \notin \mathcal{P} \\ \mathrm{RAD}_{1\times S} & \text{if } i \in \mathcal{P} \end{cases}$
    $\boldsymbol{\Delta} = (\boldsymbol{\delta}_1; \cdots; \boldsymbol{\delta}_i; \cdots; \boldsymbol{\delta}_N)$.
    $\boldsymbol{U}^+ \leftarrow \boldsymbol{U}_n + c_n\boldsymbol{\Delta}; \boldsymbol{U}^- \leftarrow \boldsymbol{U}_n - c_n\boldsymbol{\Delta}$.
    Compute $\widehat{\nabla\Phi}(\boldsymbol{U}_n)$ based on Equation 5.
    Compute $\boldsymbol{U}_{n+1}$ based on Equation 7.
    Set $(\boldsymbol{u}_1; \cdots; \boldsymbol{u}_i; \cdots; \boldsymbol{u}_N) \leftarrow \boldsymbol{U}_{n+1}$.
    **for** $i \in \mathcal{V}$ **do**
      **if** $i \in \mathcal{P}$ **then**
        $\boldsymbol{u}_i \leftarrow \min(\epsilon^+\boldsymbol{x}_i, \boldsymbol{u}_i)$
        $\boldsymbol{u}_i \leftarrow \max(-\epsilon^-\boldsymbol{x}_i, \boldsymbol{u}_i)$
      **else**
        $\boldsymbol{u}_i \leftarrow \boldsymbol{0}$
      **end if**
    **end for**
    $\boldsymbol{U}_{n+1} \leftarrow (\boldsymbol{u}_1; \cdots; \boldsymbol{u}_i; \cdots; \boldsymbol{u}_N)$
  **end for**
  $\boldsymbol{U} \leftarrow \boldsymbol{U}_{\texttt{MaxIter}+1}$
  $\boldsymbol{X}' \leftarrow \boldsymbol{X} + \boldsymbol{U}$
  Return $\boldsymbol{X}', \boldsymbol{U}$

---

one node and the perturbation $u$ is small enough, then the object function in Formulation 3 is locally convex.

**Proposition 2.** *The objective function $\Phi$ in Formulation 3 is locally convex under small perturbation $\boldsymbol{U}$.*

*Proof.* Given function $\Phi$ is smooth and attains global optimal when the perturbation $\boldsymbol{U} = \boldsymbol{0}$, there exists one region around $\boldsymbol{U} = \boldsymbol{0}$, in which $\Phi(\boldsymbol{U})$ is convex. $\quad\square$

Existing literature has shown that the convex separable nonlinear knapsack problems could be approximated with the greedy search framework described by Algorithm 2 [67]. Recalling Proposition 1, the perturbation on GCN-based models would only arise local effect, which means if the attack nodes are selected L-hop away from each other, the objective of Formulation 3 would be separable. We also observed that the objective is locally convex given attack on only one node in Proposition 2. Combing the enlightenment we get from Proposition 1 and 2, the greedy search framework would work efficiently for solving formulation 3.

The proposed solution procedure consists of two steps: 1) compute $\hat{\phi}_i$ to approximate $\phi_i$ for each $i$; 2) adopt the greedy algorithm for the standard knapsack problem with $\hat{\phi}_i$ as the utility. In Step 1, we proposed that the utility $\hat{\phi}_i$ can be obtained by SPSA. To be precise, we run Algorithm 1 with $\mathcal{P} = \mathcal{V}$, and the algorithm outcome is $\boldsymbol{U}_{\mathcal{V}}$. Then, $\phi_i$ is approximated by $\hat{\phi}_i = \phi_i(\boldsymbol{U}_{\mathcal{V}})$; in Step 2, we initialize the

attack set $\mathcal{P}$ as an empty set, then each node is added to the attack set sequentially with the highest utility over budget. The entire procedure is referred as KG-SPSA, and details of the algorithm are presented in Algorithm 2.

---

**Algorithm 2** KG-SPSA for the diffusion attack on traffic prediction models.

---

**Input:** Traffic prediction model $f(\boldsymbol{X})$, input feature matrix $\boldsymbol{X}$, total budget $B$, and cost of each node $\{b_i\}_{i \in \mathcal{V}}$.
**Output:** Adversarial sample $\boldsymbol{X}'$, optimal perturbation $\boldsymbol{U}$, and attack set $\mathcal{P}$.
Initialize $\mathcal{P} = \emptyset$.
Evaluate $\hat{\phi}_i = \phi_i(\boldsymbol{U}_{\mathcal{V}}), i \in \mathcal{V}$ with Algorithm 1.
**while** $\sum_{i \in \mathcal{P}} b_i \leq B$ **do**
  Set `max_utility` $= -\infty$, `max_idx` $= -\infty$.
  **for** $i \in \mathcal{V} \setminus \mathcal{P}$ **do**
    **if** $\frac{\hat{\phi}_i}{b_i} >$ `max_utility` **then**
      Set `max_utility` $= \frac{\hat{\phi}_i}{b_i}$.
      Set `max_idx` $= i$.
    **end if**
  **end for**
  **if** $\sum_{i \in \mathcal{P}} b_i + b_{\texttt{max\_idx}} \leq B$ **then**
    $\mathcal{P} = \mathcal{P} \cup \{\texttt{max\_idx}\}$
  **end if**
**end while**
Run Algorithm 1 with fixed $\mathcal{P}$, obtain $\boldsymbol{X}'$, $\boldsymbol{U}$.
Return $\boldsymbol{X}', \boldsymbol{U}, \mathcal{P}$.

---

It is possible to adopt other methods to estimate $\hat{\phi}_i$ in Step 1, including clustering methods and graph centrality measures. These algorithms will be viewed as baseline algorithms and compared with KG-SPSA in numerical experiments.

## IV. EXPERIMENTS

In this section, we evaluate the performance of the proposed diffusion attack algorithm under different scenarios using real-world data.

### A. Experiment Setup

We examine the proposed attack algorithm on three traffic prediction models and two datasets, and details are described as follows:

**Traffic Data.** We consider two real-world traffic datasets:

- LA: The LA dataset contains traffic speed obtained from 207 loop detectors in Los Angeles [30], and The data ranges from March 1st to March 7th, 2012. The average degree of the adjacency matrix is 14. The speed data are collected every 5 minutes, and the average speed is 58km/h. The study area is showed in the upper part of Fig. 4.
- HK: The HK speed data is collected from an open data platform initiated by the Hong Kong government, and overall 179 roads are considered in the Hong Kong island and Kowloon area. The data ranges from May 1st to May 31st, 2020. The average degree of the adjacency matrix

is 2. The speed data are collected every 5 minutes, and the average speed is 45km/h. The study area is showed in the lower part of Fig. 4.

**Traffic prediction models.** We evaluate the developed diffusion attack framework on three traffic prediction models: T-GCN [30], ST-GCN [31], and A3T-GCN [32], which are all based on GCN structures. For each model, we set $S = 12$ and $T = 1$. To conduct the comprehensive evaluation, we train four variants of each model, which are the original model, the original model with DROPOUT, DROPNODE, and DROPE-DGE regularization, respectively [68], [69]. For different drop regularization strategies, we set the drop probability to 30%. In Appendix A, the Accuracy and Root Mean Squared Error (RMSE) of each model are showed in TABLE III and TABLE IV, and both measures are defined in [30]. Overall, accuracy of the different prediction models are around 90% in testing data.

**Attack settings.** The parameter settings for the diffusion attack models and algorithms, as well as the evaluation criterion are discussed as follows:

- *Model specifications.* In Equation 3, we set the constraint $-z_i \boldsymbol{x}_i \leq \boldsymbol{u}_i \leq 0.5 z_i \boldsymbol{x}_i$, and $w_i = 1$, which means each node is equally important. The cost is defined as $b_i = \text{Degree}(i) = S_D(i)$, where $S_D(i) = \sum_j \boldsymbol{A}_{ij}$. $B = \{20, 50, 100, 150, 200\}$. The optimal target reduces to $\Phi(\boldsymbol{U}) = \sum_{i \in \mathcal{V}} \phi_i(\boldsymbol{U})$, where $\phi_i(\boldsymbol{U}) = y_i'(\boldsymbol{X}') - y_i(\boldsymbol{X})$. The attack algorithm aims to reduce the predicted speed and we intend to generate virtual "congestion" on the network.
- *Evaluation of Algorithms.* We use Average Attack Influence (AAI) and Average Attack Influence Ratio (AAIR) to evaluate the effect of the diffusion attacks. We define

$$
\begin{aligned}
\text{AAI} &= \frac{1}{N} \sum_{i \in \mathcal{V}} |\phi_i(\boldsymbol{U})| \\
\text{AAIR} &= \frac{1}{N} \sum_{i \in \mathcal{V}} \frac{\phi_i(\boldsymbol{U})}{|y_i(\boldsymbol{X})|}
\end{aligned} ,
$$

where AAI represents the average degradation and AAIR represents the average degradation ratio of the prediction on the entire network, respectively.
- *SPSA Setting.* For Algorithm 1, $a = 0.328, c = 0.1, \alpha = 0.202, \gamma = 0.101$, and $\eta = \frac{n}{10}$. For diffusion attack, `MaxIter` $= 30000$; for computing the $\hat{\phi}_i$, `MaxIter` $= 100$.

**Baseline algorithms.** Because the diffusion attack is a newly proposed task, there are very few existing methods that can be used as baseline algorithms. In addition to the proposed KG-SPSA approach, we modify and develop 8 algorithms for comparison. The major difference of each baseline algorithm lies in how to select the attack set $\mathcal{P}$ and whether to use the KG-$\star$ greedy algorithm. DEGREE selects nodes with highest degree $S_D(i)$. K-MEDOIDS selects nodes by clustering the nodes with geo-location features until reaching the total budget $B$ [70], PAGERANK selects nodes with highest pagerank scores $S_{PR}(i) = \frac{1-\alpha}{N} + \alpha \sum_{j \in \mathcal{N}_i} \frac{S_{PR}(j)}{|\mathcal{N}_j|}$, where $\alpha = 0.85$, and BETWEENNESS chooses nodes with high betweenness scores $S_{Bw}(i) = \sum_{j \neq i, k \neq i} \frac{\text{Path}_{jk}(i)}{\text{Path}_{jk}}$, where $\text{Path}_{jk}$ is the number of shortest path between node $j$ and $k$, and $\text{Path}_{jk}(i)$ is the number of shortest path that passes node $i$. The RANDOM

algorithm just selects the node randomly until meets the budget limit. SPSA selects the highest $\hat{\phi}_i = \phi_i(U_{\mathcal{V}})$ by running Algorithm 1 with $\mathcal{P} = \mathcal{V}$, and the greedy algorithm is not used in SPSA. When we use the greedy algorithm, it is also possible to use centrality measures such as pagerank and betweenness to represent $\hat{\phi}_i$. Different from PAGERANK and BETWEENNESS, which determine $\mathcal{P}$ by the highest scores, KG-PAGERANK and KG-BETWEENNESS approximate $\hat{\phi}_i$ by the pagerank and betweenness scores, followed by running Algorithm 2 with different $\hat{\phi}_i$.

### B. Results of diffusion attack

In this section, we present the experimental results. We first verify the diffusion effect when attacking a single node, then the performance of the proposed attack algorithm is evaluated and compared with baseline methods in different scenarios. Lastly, we examine the robustness of the proposed algorithm under different drop regularization strategies.

*1) Diffusion Effects of attacks on a single node:* To demonstrate the diffusion phenomenon, we construct an attack on a single node for the three traffic prediction models, and the attack effect of different hops of neighbors is presented in Fig. 3. As can be seen, the attack mainly influences the ego node and its 1~2-hop neighbors, and the influence will diminish as the number of hops increases. As proven in Proposition 1, for a $L$-layer GCN model, the diffusion only occurs within $L$-hop neighbors, which is verified in this experiment. The experimental results also indicate that the influence of a single node attack is localized, and a successful diffusion attack algorithm requires a scattered node selection strategy.
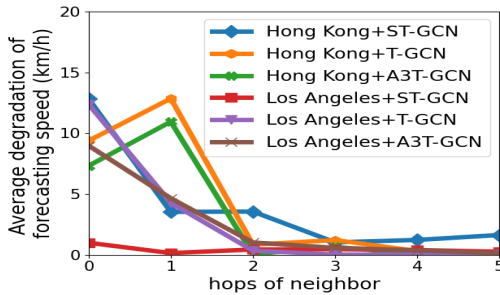


Fig. 3: Diffusion effects when attacking on a single node.

*2) Comparisons of different attack algorithms:* The proposed algorithm KG-SPSA is compared with different baseline algorithms with $B = 50$, and the corresponding AAI is presented in TABLE I (unit: km/hour), and the AAIR table is presented in Appendix B. The attack algorithms are categorized into two types: semi-black-box algorithms that know the graph structure, and black-box algorithms that only require inputs and outputs of the prediction models. From TABLE I one can see, the proposed algorithm KG-SPSA outperforms all the baseline algorithms on LA, and KG-⋆ generally outperforms the original counterparts. Comparing with other methods, KG-SPSA is a black-box method, which

does not rely on knowledge of the graph structure (*i.e.* adjacency matrix $A$). In real-world, it is challenging to obtain the information of $A$ in the prediction system as the graph can be generated by different configurations such as sensor layout, network topology, and causal relationship, and this information is hidden to users [34]. Fig. 4 presents the selected nodes by KG-SPSA for different prediction methods with $B = 50$, and the color (from green to red) represents AAIR of each node under the attack. To maximize the attack effect, the selected nodes distribute across the entire network, which is consistent with our previous conjecture.

In addition, it is observed that the robustness of the three prediction models is different. A3T-GCN demonstrates great vulnerability under attack algorithms, while both T-GCN and ST-GCN are more robust. This could be due to the strengthened connection among nodes by the attention layers in A3T-GCN, meanwhile, the strengthened connection can also increase the vulnerability of the prediction models. Comparing the dataset LA, predictions on HK are more vulnerable to adversarial diffusion attacks, which might be explained by the drastic changes of Hong Kong's traffic conditions within a day [71].

*3) Sensitivity analysis on the budget $B$:* To study the effect of budget $B$ on the diffusion attack, we run KG-SPSA with $B \in [20, 50, 100, 150, 200]$ for both LA and HK, and the corresponding AAI is presented in Fig. 5. The attack influence will increase when the total budget $B$ increases for both datasets while the trend is becoming marginal. It is also observed that prediction models on HK are more vulnerable, while A3T-GCN is the least robust predictions models for both datasets.

*4) Performance of the proposed algorithm on drop regularization strategies:* To better understand the performance of the proposed attack algorithms, we examine the attack effects on the prediction models with drop regularization. Existing studies have widely demonstrated that drop regularization strategies could improve the robustness of the GCN-based model [68], [69], hence it is crucial to show that the performance of the proposed methods remains effective under different drop regularization strategies. To this end, we conduct diffusion attacks with KG-SPSA on the prediction models trained with DROPOUT, DROPNODE, and DROPEGDE. DROPOUT randomly drops rows in feature matrix $X$, DROPNODE randomly drops a subset of the nodes on the graph, and DROPEGDE will randomly drop the edges of the graph for each epoch during the model training. Details of the models are presented in section IV-A. We run KG-SPSA and KG-PAGERANK for prediction models with different drop regularization on the two datasets, and the algorithm performance is presented in TABLE II. The reason we choose KG-SPSA and KG-PAGERANK is because both algorithms outperform other semi-black-box and black-box algorithms. For dataset LA, KG-SPSA outperforms KG-PAGERANK on all the prediction models, and KG-SPSA is slightly better on HK. In most cases, DROPOUT could degrade the performance of the attack algorithms, while the other two drop regularization strategies do not protect the prediction models. Overall, the proposed diffusion attack algorithms could still generate adversarial samples under various

TABLE I: Comparison of different diffusion attack algorithms in terms of AAI. ($B = 50$)

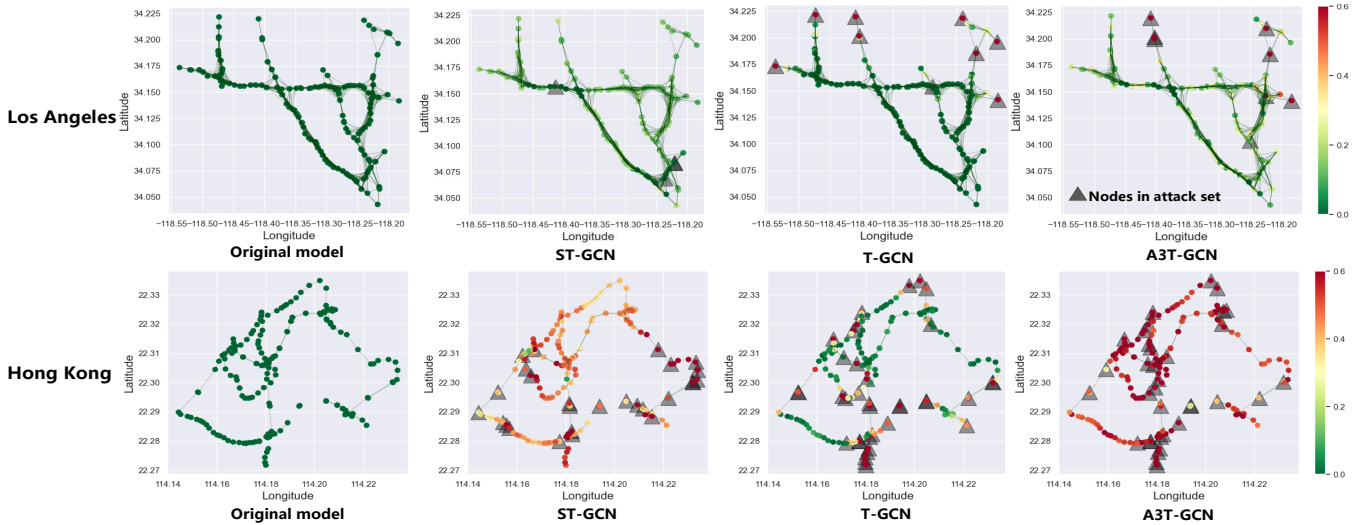| Types | Algorithm | LA | | | HK | | |
|---|---|---|---|---|---|---|---|
| | | ST-GCN | T-GCN | A3T-GCN | ST-GCN | T-GCN | A3T-GCN |
| Semi-black-box | DEGREE | 0.74 | 0.65 | 0.64 | 5.30 | 3.64 | 4.22 |
| | K-MEDOIDS | 0.69 | 0.87 | 1.69 | 12.41 | 5.72 | 7.74 |
| | PAGERANK | 2.41 | 2.11 | 2.61 | 9.42 | 4.84 | 5.54 |
| | BETWEENNESS | 2.35 | 1.56 | 2.06 | 16.28 | 7.84 | 19.88 |
| | KG-BETWEENNESS | 2.75 | 1.92 | 2.66 | 17.37 | 8.42 | 24.44 |
| | KG-PAGERANK | 4.74 | 3.69 | 5.06 | 22.63 | **12.99** | **34.47** |
| Black-box | RANDOM | 1.06 | 1.31 | 1.86 | 14.61 | 8.74 | 20.42 |
| | SPSA | 3.18 | 1.46 | 7.66 | 18.60 | 7.43 | 28.97 |
| | KG-SPSA | **5.46** | **4.36** | **12.74** | **23.34** | *12.21* | *32.86* |



Fig. 4: The distribution of selected nodes and AAIR of each node under the attack algorithm KG-SPSA. (the selected nodes are marked as triangle, and the color represents AAIR)
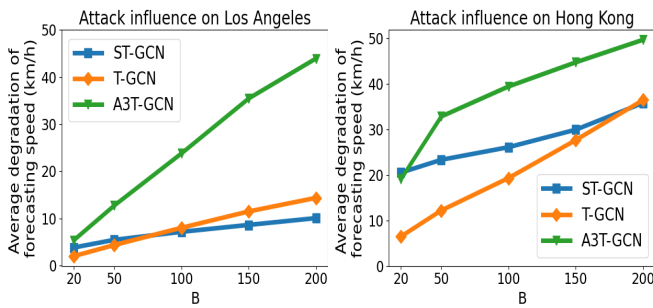


Fig. 5: Attack effect with different budget $B$.

TABLE II: Comparison of KG-SPSA and KG-PAGERANK on the three defense strategies in terms of AAI. ($B = 50$)

| Datasets | Model | Baseline | DROPOUT | DROPNODE | DROPEDGE |
|---|---|---|---|---|---|
| | | KG-SPSA | | | |
| LA | ST-GCN | **5.46** | **5.65** | **6.86** | **11.60** |
| | T-GCN | **4.36** | 3.43 | 3.49 | 2.79 |
| | A3T-GCN | **12.74** | 3.07 | **18.59** | 6.18 |
| HK | ST-GCN | **23.34** | 11.89 | 14.53 | **25.84** |
| | T-GCN | 12.21 | **12.34** | 11.63 | 11.43 |
| | A3T-GCN | 32.86 | **41.77** | 11.44 | **91.77** |
| | | KG-PAGERANK | | | |
| LA | ST-GCN | 4.74 | 2.71 | 4.65 | 6.67 |
| | T-GCN | 3.69 | 3.06 | 3.21 | 2.47 |
| | A3T-GCN | 5.06 | 2.68 | 6.07 | 3.42 |
| HK | ST-GCN | 22.63 | **12.84** | **16.02** | 24.84 |
| | T-GCN | **12.99** | **12.34** | **12.31** | **12.04** |
| | A3T-GCN | **34.47** | 28.34 | **24.10** | 74.93 |

drop regularization strategies.

## C. Discussions

In this section, we discuss the implications and suggestions for improving the robustness of the traffic prediction models. In the previous section, we carry out numerical experiments to demonstrate the performance of the proposed attack algorithms on different datasets, prediction models, and regularization strategies. Based on the experimental results, we provide the following suggestions to improve the model robustness during different phrases:

- **Model selection.** When choosing GCN-based models for speed prediction, RNN-based models are generally more robust than attention-based models. Depending on the data and city scale, it is suggested to choose models with simpler layers, as the complex layers in ST-GCN and A3T-GCN can sometimes degrade significantly under attacks. There is a trade-off between accuracy and robustness, so it is critical to balance the accuracy and robustness for practical usage.

- **Model regularization.** Based on the experimental results, it is suggested to adopt DROPOUT during the training, as the model accuracy remains high while the robustness can be improved after the DROPOUT training. It is also suggested to test different drop regularization strategies before the actual deployment.
- **Model privacy.** The graph structure should not be disclosed to the public, as it can significantly improve the efficiency of the attack algorithms. It is also suggested to frequently update the prediction model, as the attack models rely on multiple trials and errors on the prediction models. If the prediction model updates frequently, then the robustness of the entire prediction system can be significantly improved.
- **Active defending strategies.** Before actual deployment, it is necessary to comprehensively test the vulnerability of the prediction models and to identify the critical nodes with significant attack influence. For those important nodes, we can enhance the protection by regular patrol in the physical world and consistency checking in the cyber system. For example, if an RSU on a road segment is identified to be critical, then this device should be protected physically [72]. If the attack on this device indeed occurs, the traffic center should spot the anomaly in real-time and block the information sent from this device.

## V. CONCLUSION

In this paper, we explore the robustness and vulnerability issues of graph-based neural network models for traffic prediction. Different from existing adversarial attack tasks, adversarial attacks for traffic prediction require to degrade the model performance for the entire network, rather than a specific sample of nodes. Given this, we propose a novel concept of diffusion attack, which aims to reduce the prediction accuracy of the whole traffic network by perturbing a small number of nodes. To solve for the diffusion attack task, we develop an algorithm KG-SPSA, which consists of two major components: 1) using SPSA to generate the optimal perturbations to maximize the attack effects; 2) adapting the greedy algorithm in the knapsack problem to select the most critical nodes. The proposed algorithm is examined with three widely used GCN-based traffic prediction models (ST-GCN, T-GCN, and A3T-GCN) on the Los Angeles and Hong Kong datasets. The experimental results indicate that the proposed algorithm outperforms the baseline algorithms under various scenarios, which demonstrates the effectiveness and efficiency of the proposed algorithm. In addition, the proposed attack algorithms can still generate effective adversarial samples for traffic prediction models trained with drop regularization. This study could help the public agencies and private sectors better understand the robustness and vulnerability of GCN-based traffic prediction models under adversarial attacks, and strategies to improve the model robustness in different phrases are also discussed.

As for the future research directions, the proposed attack algorithms could be applied to not only road traffic prediction, but also other traffic modes such as urban railway transit systems, ride-sourcing services, and parking systems [73], [74]. It is also important to study the effect of adversarial attacks on flow prediction, origin-destination demand prediction, and other tasks relying on the GCN-based models. For the users of the traffic prediction models, it is critical to develop models for defending adversarial attacks and protecting traffic prediction results. Another way to protect the prediction model is through real-time anomaly detection and filtering of the incoming data stream, which could be a new research direction for improving the robustness of the traffic prediction models under adversarial attacks.

## SUPPLEMENTARY MATERIALS

The proposed diffusion attack algorithm and evaluation framework are implemented in Python and open-sourced on GitHub (https://github.com/LYZ98/Adversarial-Diffusion-Attacks-on-Graph-based-Traffic-Prediction-Models).

## ACKNOWLEDGMENT

## APPENDIX A
### MORE DETAILS ABOUT THE THREE TRAFFIC PREDICTION MODELS AND ATTACK RESULTS

To train the three traffic prediction models, we set the learning rate to be $0.001$, batch size to be $32$, and the number of epoch to be $300$. The two datasets are divided into two parts, in which $80\%$ and $20\%$ are training set and testing set, respectively. Mean Squared Error (MSE) is used as the loss function [30], and Adam is adopted as the optimizer. The testing accuracy in terms of accuracy and Root Mean Squared Error (RMSE) of the trained prediction models are presented in TABLE III and TABLE IV, respectively. Overall, all the prediction models could achieve high prediction accuracy on both datasets.

### TABLE III: Accuracy of trained models

| Dataset | Model | Baseline | DROPOUT | DROPNODE | DROPEDGE |
|---------|-------|----------|---------|----------|----------|
| LA | ST-GCN | 92.68% | 92.70% | 92.77% | 88.18% |
| | T-GCN | 90.44% | 89.24% | 90.01% | 90.05% |
| | A3T-GCN | 89.04% | 72.71% | 89.15% | 89.84% |
| HK | ST-GCN | 92.88% | 93.20% | 93.17% | 92.80% |
| | T-GCN | 88.50% | 87.70% | 86.78% | 88.08% |
| | A3T-GCN | 89.47% | 87.30% | 80.57% | 88.12% |

## APPENDIX B
### COMPARISONS OF DIFFERENT ATTACK ALGORITHMS IN TERMS OF AAIR. ($B = 50$)

Similar to TABLE I, we evaluate the performance of different attack algorithms in terms of AAIR in TABLE V. In general, similar arguments could be obtained based on the

TABLE IV: RMSE of trained models

| Dataset | Model | Baseline | DROPOUT | DROPNODE | DROPEDGE |
|---|---|---|---|---|---|
| LA | ST-GCN | 4.30 | 4.28 | 4.25 | 6.95 |
| | T-GCN | 5.61 | 6.32 | 5.86 | 5.84 |
| | A3T-GCN | 6.43 | 16.02 | 6.36 | 5.96 |
| HK | ST-GCN | 3.55 | 3.39 | 3.41 | 3.59 |
| | T-GCN | 5.74 | 6.13 | 6.59 | 5.95 |
| | A3T-GCN | 5.25 | 6.33 | 9.68 | 5.92 |

AAIR, and the proposed KG-SPSA outperforms other baseline models in LA, and performs similarly as the semi-black-box algorithms in HK.

Results in TABLE VI follow the same patterns as in TABLE II, and one can observe that the proposed attack algorithms can still generate effective adversarial samples in terms of AAIR.

## REFERENCES

[1] E. I. Vlahogianni, J. C. Golias, and M. G. Karlaftis, "Short-term traffic forecasting: Overview of objectives and methods," *Transport reviews*, vol. 24, no. 5, pp. 533–557, 2004.

[2] O. Lange and L. Perez, "Traffic prediction with advanced graph neural networks." [Online]. Available: https://deepmind.com/blog/article/traffic-prediction-with-advanced-graph-neural-networks

[3] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.

[4] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[5] J. Su, D. V. Vargas, and K. Sakurai, "One pixel attack for fooling deep neural networks," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 828–841, 2019.

[6] K. Zhou, T. P. Michalak, M. Waniek, T. Rahwan, and Y. Vorobeychik, "Attacking similarity-based link prediction in social networks," in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, ser. AAMAS '19. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2019, p. 305–313.

[7] Y. Li, X. Xu, J. Xiao, S. Li, and H. T. Shen, "Adaptive square attack: Fooling autonomous cars with adversarial traffic signs," *IEEE Internet of Things Journal*, pp. 1–1, 2020.

[8] M. Fang, G. Yang, N. Z. Gong, and J. Liu, "Poisoning attacks to graph-based recommender systems," *Proceedings of the 34th Annual Computer Security Applications Conference*, Dec 2018.

[9] S. Weckert, "Google maps hacks." [Online]. Available: http://www.simonweckert.com/googlemapshacks.html

[10] A. Bayen, J. Butler, A. Patire *et al.*, "Mobile millennium," Tech. Rep. UCB-ITS-CWP-2011-6, CCIT Research Report, UC Berkeley, Tech. Rep., 2011.

[11] M. T. Ahvanooey, Q. Li, M. Rabbani, and A. R. Rajput, "A survey on smartphones security: software vulnerabilities, malware, and attacks," *arXiv preprint arXiv:2001.09406*, 2020.

[12] N. Lu, N. Cheng, N. Zhang, X. Shen, and J. W. Mark, "Connected vehicles: Solutions and challenges," *IEEE Internet of Things Journal*, vol. 1, no. 4, pp. 289–299, 2014.

[13] B. K. J. Al-Shammari, N. Al-Aboody, and H. S. Al-Raweshidy, "Iot traffic management and integration in the qos supported network," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 352–370, 2018.

[14] Y. J. Edes, P. G. Michalopoulos, and R. A. Plum, "Improved estimation of traffic flow for real-time control," *Transportation Research Record*, vol. 7, no. 9, p. 28, 1980.

[15] M. S. Ahmed and A. R. Cook, "Analysis of freeway traffic time-series data by using box-jenkins techniques," *Transportation Research Record*, no. 722, pp. 1–9, 1979.

[16] M. M. Hamed, H. R. Al-Masaeid, and Z. M. B. Said, "Short-term prediction of traffic volume in urban arterials," *Journal of Transportation Engineering*, vol. 121, no. 3, pp. 249–254, 1995.

[17] M. Van Der Voort, M. Dougherty, and S. Watson, "Combining kohonen maps with arima time series models to forecast traffic flow," *Transportation Research Part C: Emerging Technologies*, vol. 4, no. 5, pp. 307–318, 1996.

[18] S. Lee and D. B. Fambro, "Application of subset autoregressive integrated moving average model for short-term freeway traffic volume forecasting," *Transportation Research Record*, vol. 1678, no. 1, pp. 179–188, 1999.

[19] B. M. Williams and L. A. Hoel, "Modeling and forecasting vehicular traffic flow as a seasonal arima process: Theoretical basis and empirical results," *Journal of transportation engineering*, vol. 129, no. 6, pp. 664–672, 2003.

[20] C.-H. Wu, J.-M. Ho, and D.-T. Lee, "Travel-time prediction with support vector regression," *IEEE transactions on intelligent transportation systems*, vol. 5, no. 4, pp. 276–281, 2004.

[21] F. G. Habtemichael and M. Cetin, "Short-term traffic flow rate forecasting based on identifying similar traffic patterns," *Transportation research Part C: emerging technologies*, vol. 66, pp. 61–78, 2016.

[22] I. Okutani and Y. J. Stephanedes, "Dynamic prediction of traffic volume through kalman filtering theory," *Transportation Research Part B: Methodological*, vol. 18, no. 1, pp. 1–11, 1984.

[23] C. P. Van Hinsbergen, T. Schreiter, F. S. Zuurbier, J. Van Lint, and H. J. Van Zuylen, "Localized extended kalman filter for scalable real-time traffic state estimation," *IEEE transactions on intelligent transportation systems*, vol. 13, no. 1, pp. 385–394, 2011.

[24] 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.

[25] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[26] A. Azzouni and G. Pujolle, "A long short-term memory recurrent neural network framework for network traffic matrix prediction," *arXiv preprint arXiv:1705.05690*, 2017.

[27] N. Ramakrishnan and T. Soni, "Network traffic prediction using recurrent neural networks," in *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2018, pp. 187–193.

[28] S. Guo, Y. Lin, N. Feng, C. Song, and H. Wan, "Attention based spatial-temporal graph convolutional networks for traffic flow forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 922–929.

[29] F. Zhou, Q. Yang, K. Zhang, G. Trajcevski, T. Zhong, and A. Khokhar, "Reinforced spatiotemporal attentive graph neural networks for traffic forecasting," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6414–6428, 2020.

[30] L. Zhao, Y. Song, C. Zhang, Y. Liu, P. Wang, T. Lin, M. Deng, and H. Li, "T-gcn: A temporal graph convolutional network for traffic prediction," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 9, pp. 3848–3858, 2020.

[31] B. Yu, H. Yin, and Z. Zhu, "Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 3634–3640.

[32] J. Zhu, Y. Song, L. Zhao, and H. Li, "A3t-gcn: attention temporal graph convolutional network for traffic forecasting," *arXiv preprint arXiv:2006.11583*, 2020.

[33] B. Yu, Y. Lee, and K. Sohn, "Forecasting road traffic speeds by considering area-wide spatio-temporal dependencies based on a graph convolutional neural network (gcn)," *Transportation Research Part C: Emerging Technologies*, vol. 114, pp. 189–204, 2020.

[34] 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*, p. 1–21, 2020.

[35] D. Wang, J. Zhang, W. Cao, J. Li, and Y. Zheng, "When will you arrive? estimating travel time based on deep neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[36] L. Bai, L. Yao, S. Kanhere, X. Wang, Q. Sheng *et al.*, "Stg2seq: Spatial-temporal graph to sequence model for multi-step passenger demand forecasting," *arXiv preprint arXiv:1905.10069*, 2019.

[37] L. Liu, J. Chen, H. Wu, J. Zhen, G. Li, and L. Lin, "Physical-virtual collaboration modeling for intra-and inter-station metro ridership prediction," *IEEE Transactions on Intelligent Transportation Systems*, 2020.

[38] L. Lin, Z. He, and S. Peeta, "Predicting station-level hourly demand in a large-scale bike-sharing network: A graph convolutional neural network approach," *Transportation Research Part C: Emerging Technologies*, vol. 97, pp. 258–276, Dec 2018.

[39] J. Yang, B. Guo, Z. Wang, and Y. Ma, "Hierarchical prediction based on network-representation-learning-enhanced clustering for bike-sharing

TABLE V: Comparison of different diffusion attack algorithms in terms of AAIR. ($B = 50$)

| Types | Algorithm | LA | | | HK | | |
|---|---|---|---|---|---|---|---|
| | | ST-GCN | T-GCN | A3T-GCN | ST-GCN | T-GCN | A3T-GCN |
| Semi-blackbox | DEGREE | 1.54% | 1.88% | 2.07% | 12.35% | 8.69% | 10.28% |
| | K-MEDOIDS | 1.37% | 1.77% | 2.33% | 25.19% | 11.84% | 17.09% |
| | PAGERANK | 3.79% | 3.77% | 3.80% | 19.90% | 10.77% | 12.57% |
| | BETWEENNESS | 3.83% | 3.70% | 3.34% | 30.95% | 14.46% | 43.92% |
| | KG-BETWEENNESS | 4.73% | 4.21% | 2.55% | 32.69% | 15.07% | 52.50% |
| | KG-PAGERANK | 7.88% | 6.99% | 8.38% | 42.27% | 23.91% | **72.37**% |
| Blackbox | RANDOM | 1.65% | 2.51% | 3.19% | 28.57% | 16.63% | 45.33% |
| | SPSA | 5.80% | 3.07% | 12.36% | 35.61% | 15.00% | 60.63% |
| | KG-SPSA | **8.32**% | **7.76**% | **22.77**% | **43.25**% | **24.26**% | 70.21% |

TABLE VI: Comparison of KG-SPSA and KG-PAGERANK on the three defense strategies im terms of AAIR. ($B = 50$)

| Dataset | Model | Baseline | DROPOUT | DROPNODE | DROPEDGE |
|---|---|---|---|---|---|
| | | KG-SPSA | | | |
| LA | ST-GCN | **8.32**% | **8.73**% | **11.36**% | **13.88**% |
| | T-GCN | **7.76**% | **6.21**% | **6.15**% | **4.98**% |
| | A3T-GCN | **22.77**% | **8.33**% | **38.11**% | **14.35**% |
| HK | ST-GCN | **43.25**% | 23.66% | 27.51% | **46.36**% |
| | T-GCN | **24.26**% | **23.47**% | **22.83**% | 22.75% |
| | A3T-GCN | 70.21% | **79.80**% | 24.34% | **230.85**% |
| | | KG-PAGERANK | | | |
| LA | ST-GCN | 7.88% | 4.47% | 7.84% | 9.38% |
| | T-GCN | 6.99% | 5.60% | 5.71% | 4.44% |
| | A3T-GCN | 8.38% | 6.93% | 12.08% | 6.60% |
| HK | ST-GCN | 42.27% | **24.47**% | **29.25**% | 45.24% |
| | T-GCN | 23.91% | 23.09% | 22.07% | **22.86**% |
| | A3T-GCN | **72.37**% | 52.71% | **56.86**% | 174.56% |

system in smart city," *IEEE Internet of Things Journal*, vol. 8, no. 8, pp. 6416–6424, 2021.

[40] K. F. Chu, A. Y. S. Lam, and V. O. K. Li, "Deep multi-scale convolutional lstm network for travel demand and origin-destination predictions," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 8, pp. 3219–3232, 2020.

[41] A. Monti, A. Bertugli, S. Calderara, and R. Cucchiara, "Dag-net: Double attentive graph neural network for trajectory forecasting," *arXiv preprint arXiv:2005.12661*, 2020.

[42] A. Mohamed, K. Qian, M. Elhoseiny, and C. Claudel, "Social-stgcnn: A social spatio-temporal graph convolutional neural network for human trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 424–14 432.

[43] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, "Geometric deep learning: going beyond euclidean data," *IEEE Signal Processing Magazine*, vol. 34, no. 4, pp. 18–42, 2017.

[44] X. Ma, Z. Tao, Y. Wang, H. Yu, and Y. Wang, "Long short-term memory neural network for traffic speed prediction using remote microwave sensor data," *Transportation Research Part C: Emerging Technologies*, vol. 54, pp. 187–197, 2015.

[45] J. Dai, W. Zhu, and X. Luo, "A targeted universal attack on graph convolutional network," *arXiv preprint arXiv:2011.14365*, 2020.

[46] D. Zügner and S. Günnemann, "Adversarial attacks on graph neural networks via meta learning," *arXiv preprint arXiv:1902.08412*, 2019.

[47] H. Dai, H. Li, T. Tian, X. Huang, L. Wang, J. Zhu, and L. Song, "Adversarial attack on graph structured data," *arXiv preprint arXiv:1806.02371*, 2018.

[48] J. Ma, S. Ding, and Q. Mei, "Towards more practical adversarial attacks on graph neural networks," *Advances in neural information processing systems*, 2020.

[49] J. Ma, J. Deng, and Q. Mei, "Near-black-box adversarial attacks on graph neural networks as an influence maximization problem," in *ICLR Conference OpenReview*, 2021.

[50] K. Xu, H. Chen, S. Liu, P.-Y. Chen, T.-W. Weng, M. Hong, and X. Lin, "Topology attack and defense for graph neural networks: An optimization perspective," *arXiv preprint arXiv:1906.04214*, 2019.

[51] X. Xu, X. Du, and Q. Zeng, "Attacking graph-based classification without changing existing connections," in *Annual Computer Security Applications Conference*, ser. ACSAC '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 951–962.

[52] F. Liu, L. M. Moreno, and L. Sun, "One vertex attack on graph neural networks-based spatiotemporal forecasting," in *ICLR Conference OpenReview*, 2021.

[53] B. Finkelshtein, C. Baskin, E. Zheltonozhskii, and U. Alon, "Single-node attack for fooling graph neural networks," *arXiv preprint arXiv:2011.03574*, 2020.

[54] B. Wang, T. Zhou, M. Lin, P. Zhou, A. Li, M. Pang, C. Fu, H. Li, and Y. Chen, "Efficient evasion attacks to graph neural networks via influence function," *arXiv preprint arXiv:2009.00203*, 2020.

[55] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *2016 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.

[56] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*, 2017, pp. 39–57.

[57] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2574–2582.

[58] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, 2017, pp. 506–519.

[59] W. Brendel, J. Rauber, and M. Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," *arXiv preprint arXiv:1712.04248*, 2017.

[60] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, 2017, pp. 15–26.

[61] C.-C. Tu, P. Ting, P.-Y. Chen, S. Liu, H. Zhang, J. Yi, C.-J. Hsieh, and S.-M. Cheng, "Autozoom: Autoencoder-based zeroth order optimization method for attacking black-box neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 742–749.

[62] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin, "Black-box adversarial attacks with limited queries and information," *arXiv preprint arXiv:1804.08598*, 2018.

[63] N. Akhtar and A. Mian, "Threat of adversarial attacks on deep learning in computer vision: A survey," *Ieee Access*, vol. 6, pp. 14 410–14 430, 2018.

[64] J. C. Spall, "An overview of the simultaneous perturbation method for efficient optimization," *Johns Hopkins apl technical digest*, vol. 19, no. 4, pp. 482–492, 1998.

[65] J. Uesato, B. O'Donoghue, P. Kohli, and A. van den Oord, "Adversarial risk and the dangers of evaluating against weak attacks," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. Stockholmsmässan, Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 5025–5034.

[66] S. Martello, D. Pisinger, and P. Toth, "Dynamic programming and strong bounds for the 0-1 knapsack problem," *Management science*, vol. 45, no. 3, pp. 414–424, 1999.

[67] B. Zhang and Z. Hua, "A unified method for a class of convex separable nonlinear knapsack problems," *European Journal of Operational Research*, vol. 191, no. 1, pp. 1–6, 2008.

[68] Y. Rong, W. Huang, T. Xu, and J. Huang, "Dropedge: Towards deep graph convolutional networks on node classification," in *ICLR Conference OpenReview*, 2020.

[69] X. L. Lingwei Chen and D. Wu, "Enhancing robustness of graph convolutional networks via dropping graph connections." [Online]. Available: https://faculty.ist.psu.edu/wu/papers/DropCONN.pdf

[70] N. K. Kaur, U. Kaur, and D. D. Singh, "K-medoid clustering algorithm-a review," *International Journal of Computer Application and Technology (IJCAT)*, vol. 1, no. 1, pp. 2349–1841, 2014.

[71] M. L. Tam and W. H. Lam, "Application of automatic vehicle identification technology for real-time journey time estimation," *Information Fusion*, vol. 12, no. 1, pp. 11–19, 2011.

[72] J. Zhang, Y. Wang, S. Li, and S. Shi, "An architecture for iot-enabled smart transportation security system: A geospatial approach," *IEEE Internet of Things Journal*, vol. 8, no. 8, pp. 6205–6213, 2021.

[73] S. Yang, W. Ma, X. Pi, and S. Qian, "A deep learning approach to real-time parking occupancy prediction in transportation networks incorporating multiple spatio-temporal data sources," *Transportation Research Part C: Emerging Technologies*, vol. 107, pp. 248–265, 2019.

[74] J. Zhang, H. Che, F. Chen, W. Ma, and Z. He, "Short-term origin-destination demand prediction in urban rail transit systems: A channel-wise attentive split-convolutional neural network method," *Transportation Research Part C: Emerging Technologies*, vol. 124, p. 102928, 2021.

**Wei Ma** received bachelor's degrees in Civil Engineering and Mathematics from Tsinghua University, China, master degrees in Machine Learning and Civil and Environmental Engineering, and PhD degree in Civil and Environmental Engineering from Carnegie Mellon University, USA. He is currently an assistant professor with the Department of Civil and Environmental Engineering at the Hong Kong Polytechnic University (PolyU). His research focuses on intersection of machine learning, data mining, and transportation network modeling, with applications for smart and sustainable mobility systems. He has received awards for research excellence and his contributions to the area, including 2020 Mao Yisheng Outstanding Dissertation Award, and best paper award (theoretical track) at INFORMS Data Mining and Decision Analytics Workshop.

**Lyuyi ZHU** is an undergraduate student from College of Civil Engineering and Architecture, Zhejiang University, Hangzhou, China. He will join the School of Data Science, City University of Hong Kong as a PhD student. His research interest includes machine learning, optimization and numerical method.

**Kairui Feng** is currently a Ph.D. student from Department of Civil and Environmental Engineering, Princeton University, New Jersey, USA. He received bachelor's degrees in Civil Engineering and Mathematics from Tsinghua University, China. His research interest includes infrastructure system modeling/optimization and climate change using data-driven and numerical approaches.

**Ziyuan Pu** is currently a Lecturer (Assistant Professor) at Monash University. He received B.S. degree in transportation engineering in 2010 at Southeast University, China. He received M.S. and Ph.D. degree in civil and environmental engineering in 2015 and 2020, respectively, at the University of Washington, US. His research interest includes transportation data science, smart transportation infrastructures, connected and autonomous vehicles (CAV), and urban computing.