# FRIGATE: Frugal Spatio-temporal Forecasting on Road Networks

Mridul Gupta
Yardi School of AI
Indian Institute of Technology, Delhi
India
Mridul.Gupta@scai.iitd.ac.in

Hariprasad Kodamana
Yardi School of AI
Indian Institute of Technology, Delhi
India
kodamana@iitd.ac.in

Sayan Ranu
Yardi School of AI
Indian Institute of Technology, Delhi
India
sayanranu@cse.iitd.ac.in

## ABSTRACT

Modelling spatio-temporal processes on road networks is a task of growing importance. While significant progress has been made on developing spatio-temporal graph neural networks (GNNs), existing works are built upon three assumptions that are not practical on real-world road networks. First, they assume sensing on every node of a road network. In reality, due to budget-constraints or sensor failures, all locations (nodes) may not be equipped with sensors. Second, they assume that sensing history is available at all installed sensors. This is unrealistic as well due to sensor failures, loss of packets during communication, etc. Finally, there is an assumption of static road networks. Connectivity within networks change due to road closures, constructions of new roads, etc. In this work, we develop FRIGATE to address all these shortcomings. FRIGATE is powered by a spatio-temporal GNN that integrates positional, topological, and temporal information into rich *inductive* node representations. The joint fusion of this diverse information is made feasible through a novel combination of *gated Lipschitz* embeddings with LSTMs. We prove that the proposed GNN architecture is provably more expressive than message-passing GNNs used in state-of-the-art algorithms. The higher expressivity of FRIGATE naturally translates to superior empirical performance conducted on real-world network-constrained traffic data. In addition, FRIGATE is robust to frugal sensor deployment, changes in road network connectivity, and temporal irregularity in sensing.

## CCS CONCEPTS

• **Applied computing** → **Transportation**; • **Computing methodologies** → **Semi-supervised learning settings**; *Cost-sensitive learning*; **Neural networks**.

## KEYWORDS

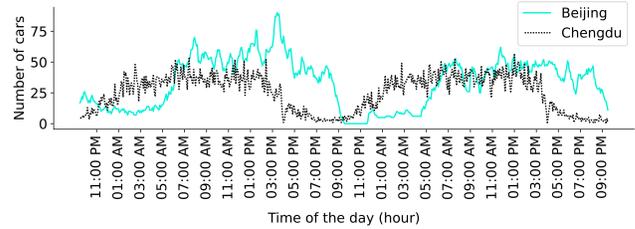spatio-temporal prediction, graph neural networks, road networks

**Figure 1: Snapshots of selected roads from Beijing and Chengdu displaying the tidal variation of traffic**

## 1 INTRODUCTION AND RELATED WORK

A road network can be modeled as a graph where nodes indicate important intersections in a region and edges correspond to streets connecting two intersections [9, 18–20, 23, 33]. Modeling the evolution of spatio-temporal events on road networks has witnessed significant interest in the recent past [6, 14, 29, 32, 35]. Specifically, each node (or edge) participates in a time-series. This time-series is a function of not only the time of the day, but also the events unfolding in other nodes of the network. The objective of the forecasting task is to model the time-series evolution in each node and predict the values in the immediate future, for e.g., the next one hour. As examples, in Fig. 1, we present the number of cars passing through two randomly selected intersections in the cities of Beijing and Chengdu. As can be seen, there is significant variation in traffic through out the day. Furthermore, the time-series vary across nodes making the forecasting problem non-trivial.

One may learn an auto-regressive model independently at each node to fit the time-series data. This strategy, however, is limited by two critical factors. First, this ignores the connectivity-induced time-series dependency among nodes. For example, if one particular intersection (node) is observing traffic congestion, it is likely for neighboring nodes connected through an outgoing edge to be affected as well. Second, the number of parameters grows linearly with the number of nodes in the graph, making it non-scalable.

### 1.1 Existing work

To address these specific needs of modeling network-dependent spatio-temporal processes, several algorithms merging models for structural data such as GNNs [10, 12, 35] or Convolutional neural networks [32] along with auto-regressive architectures have been developed. The proposed work is motivated based on some assumptions made by existing algorithms that are not realistic for real-world road networks. Table 1 summarizes these.

- **Ability to extrapolate from partial sensing:** Several of the existing algorithms assume that a sensor is placed in each node

**Table 1: Baselines**

| Algorithm | Partial sensing | Absorb network updates | Predict without temporal history |
|---|---|---|---|
| AGCRN [3] | × | × | × |
| STGODE [6] | × | × | × |
| DSTAGNN [10] | × | × | × |
| STFGNN [12] | × | × | × |
| G2S [37] | × | × | × |
| GMSDR [17] | × | × | × |
| Z-GCNETs [4] | × | × | ✓ |
| STSGCN [24] | × | × | ✓ |
| GraphWavenet [29] | × | × | ✓ |
| GMAN [36] | × | × | ✓ |
| MTGNN [28] | × | × | ✓ |
| TISV [27] | ✓ | × | ✓ |
| STGCN [35] | ✓ | ✓ | ✓ |
| DCRNN [14] | ✓ | ✓ | ✓ |
| STNN [32] | ✓ | ✓ | ✓ |
| LocaleGN [11] | ✓ | ✓ | ✓ |
| ST-GAN [26] | ✓ | ✓ | ✓ |

of the road network [3, 4, 6, 10, 12, 17, 24, 28, 29, 36, 37]. Forecasting is feasible *only* on these nodes. In reality, deploying and maintaining sensors across all intersections may be prohibitively expensive and cumbersome. Hence, it is important to also forecast accurately on nodes that do not have an explicit sensor placed on them. In this work, we show that this is indeed feasible by exploiting the network-induced dependency between nodes.

- **Ability to absorb network updates:** The connectivity within a road network may change with time. The directionality of edges may change based on time of the day, new roads may get constructed adding new nodes and edges to the network, and existing road may get removed (temporarily or permanently) to accommodate emergent needs such as street festivals, construction activities, flooding, etc. Under these circumstances, it is important to absorb small changes in the network topology without the need to retrain from scratch. Several models [3, 4, 10, 17, 24, 27–29, 36, 37] fail to absorb any updates since they are *transductive* in nature, i.e., the number of parameters in their model is a function of the network size (number of nodes or edges). In this work, we develop an inductive model, which decouples the model parameter size from that of the network. Hence, accommodating changes to network structure does not require re-training.

- **Ability to predict without temporal history or regularity:** Several of the existing algorithms can forecast on the time-series of a node *only* if the data in the past $x$ time instants are available [6, 12]. This limitation often arises from a methodology where dependency between all nodes is learned by computing similarity between their past time-series information. If past data is not available, then this similarity cannot be computed. Furthermore, some algorithms model the spatio-temporal road network as a 3-dimensional tensor, under the implicit assumption that temporal data is collected at a regular granularity (such as every five minutes) [6, 32]. In reality, sensors may fail and may therefore provide data at irregular intervals or having no temporal history in the past $x$ time instants. For a model to be deployable in real workloads, it is pertinent to be robust to sensor failures.

## 1.2 Contributions

Motivated by the above limitations, in this work, we ask the following questions: *Is it possible to design an accurate and inductive forecasting model across all nodes in a graph based on partial sensing*

*through a small subset of nodes? In addition, can the model predict on nodes with irregular time-series visibility, or in the worst case, no visibility at all?* We show that this is indeed possible through FRIGATE: FRugal and Inductive Lipschitz-GAted Spatio-TEmporal GNN. Specifically, we make the following contributions:

- **Problem formulation:** We formulate the problem of time-series forecasting on road networks. Taking a deviation from current works, the proposed formulation is cognizant of practical challenges such as limited sensor access, sensor failures, and noise in data (§ 2).

- **Novel architecture:** To enable robust forecasting, we develop a novel spatio-temporal GNN called FRIGATE. At a high-level, FRIGATE is a joint architecture composed of stacks of *siamese* GNNs and LSTMs in encoder-decoder format (§ 3). Under the hood, FRIGATE incorporates several innovations. First, to succeed in accurate forecasting under frugal sensor deployment, FRIGATE uses *Lipschitz-gated* attention over messages. Gating allows FRIGATE to receive messages from far-away neighbors without over-smoothing the node neighborhoods. In addition, Lipschitz embedding allows FRIGATE to inductively embed positional information in node representations. Second, FRIGATE is inductive and hence does not require re-training on network updates. Finally, the coupling between the GNN and LSTM is devoid of any assumption on the temporal granularity. Hence, it is robust to sensor failures or irregular data collection.

- **Empirical evaluation:** We perform extensive evaluation on real-world road network datasets from three large cities. The proposed evaluation clearly demonstrates the superiority of FRIGATE over state-of-the-art algorithms and establishes its robustness to data noise and network changes (§ 4).

## 2 PROBLEM FORMULATION

In this section, we define the concepts central to our work and formulate our problem. All key notations used in our work are summarized in Table A.1 in Appendix.

**Definition 1** (Road Network Snapshot). *A road network snapshot is represented as a directed graph $\mathcal{G} = (\mathcal{V}_t, \mathcal{E}_t, \delta, \tau_t)$, where $\mathcal{V}_t$ is the set of nodes representing road intersections at time $t$, $\mathcal{E}_t \subseteq \mathcal{V}_t \times \mathcal{V}_t$ is the set of edges representing road segments at time $t$, a distance function $\delta : \mathcal{E}_t \to \mathbb{R}$ representing the length (weight) of each road segment (edge), and the sensor readings $\tau_t = \left\{ \tau_t^v \in \mathbb{R} \mid v \in \mathcal{V}_t \right\}$ for each node at time $t \in \mathbb{N}$.*

Intuitively, a road network snapshot characterizes the state of the road network at time $t$. We use $\tau_t^v$ to denote the sensor reading at node $v$. When $\tau_t^v$ is not available, either due to non-availability of a sensor at node $v$, or due to sensor failure, we assume $\tau_t^v$ is marked with a special label. We use $\tau_t^v = \varnothing$ to denote a missing value. Furthermore, in the real world, the sensor value $\tau_t^v$ may not be recorded exactly at time $t$. We thus assume $\tau_t^v$ to be the latest value since the last snapshot; $\tau_t^v = \varnothing$ if nothing has been recorded since the last snapshot. Note that we also make the vertex and edge sets time-dependent to account for the fact that minor changes are possible on the topology over time. Examples include changing directionality of traffic on a particular road (edge), addition of new roads and intersections, temporary road closures due to street festivals, etc. We use the notation $e = (u, v)$ to denote a road

segment (edge) from node $u$ to $v$ and its length is denoted by $\delta(e)$. The length $\delta(e)$ of an edge $e$ is the *Haversine* distance from the locations represented by $u$ and $v$.

**Definition 2** (Road Network Stream). *A road network stream is the chronologically ordered set of snapshots of the road network taken at various time instances,* $\overrightarrow{\mathcal{G}} = \{\mathcal{G}_1, \mathcal{G}_2, \ldots, \mathcal{G}_T\}$*, where* $\mathcal{G}_t$ *is the road network snapshot at time instant $t$.*

We assume, that $\forall t$, $\mathcal{V}_t \approx \mathcal{V}_{t+1}$ and similarly $\mathcal{E}_t \approx \mathcal{E}_{t+1}$. This assumptions are realistic based on real-world knowledge that change events on roads are rare over both space and time.

Our goal is to learn the dynamics of the time-series at each node and forecast future values. Towards that, the modeling problem is defined as follows.

**Problem 1** (Forecasting on Road Network).

**Training:** *Given a road network stream* $\overrightarrow{\mathcal{G}} = \{\mathcal{G}_1, \ldots, \mathcal{G}_T\}$*, a forecasting horizon $\Delta$, and a timestamp $t$ where $\Delta \leq t \leq T - \Delta$, learn a function $\Psi$, parameterized by $\Theta$, to minimize the mean absolute prediction error over sensor values. Specifically,*

$$\text{minimize} \left\{ \frac{1}{|\mathcal{V}^{tr}|\Delta} \sum_{k=1}^{\Delta} \sum_{v \in \mathcal{V}_{t+k}^{tr}} \left| \Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right) - \tau_{t+k}^v \right| \right\} \quad (1)$$

*Here,* $\mathcal{V}_{t+k}^{tr} = \left\{ v \in \mathcal{V}_{t+k} \mid \tau_{t+k}^v \neq \varnothing \right\}$ *is the set of training nodes with ground truth sensor values at time $t+k$,* $\overrightarrow{\mathcal{G}}_{[1,t]} = \{\mathcal{G}_1, \cdots, \mathcal{G}_t\}$ *denotes the subset of snapshots till $t$, and* $\Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right)$ *predicts the sensor value for node $v$ at time $t+k$ when conditioned on all snapshots till now, i.e.,* $\overrightarrow{\mathcal{G}}_{[1,t]}$*. We assume that to predict on a time horizon of $\Delta$ snapshots, we must train on a history of at least $\Delta$ snapshots.*

**Inference:** *Let the last recorded snapshot be at time $t$. Given node $v \in \mathcal{V}_t$ and forecasting horizon $\Delta$, compute:*

$$\forall k \in [1, \Delta], \ \left\{ \Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right) \right\} \quad (2)$$

In addition to predicting accurately, $\Psi_\Theta$, must also satisfy the following properties:

- **Inductivity:** The number of parameters in model $\Psi_\Theta$, denoted as $|\Theta|$, should be independent of the number of nodes in the road network at any given timestamp. Inductivity enables the ability to predict on unseen nodes without retraining from scratch.
- **Permutation invariance:** Given any *permutation function* $\mathcal{P}\left(\overrightarrow{\mathcal{G}}_{[1,t]}\right)$ that randomly permutes the node set of each graph $\mathcal{G} \in \overrightarrow{\mathcal{G}}_{[1,t]}$, we require:
$\forall k \in [1, \Delta], \ \Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right) = \Psi_\Theta\left(v, t+k \mid \mathcal{P}\left(\overrightarrow{\mathcal{G}}_{[1,t]}\right)\right)$

  More simply, if a graph contains $n$ nodes, then there are $n!$ possible permutations over its node set, and hence that many adjacency matrix representations of $\mathcal{G}_t$. Permutation invariance ensures that the output is dependent only on the topology of the graph and not coupled to one of its specific adjacency matrix representations. Hence, it aids in generalizability while also ensuring that the location of change in the adjacency matrix owing to a node/edge insertion or deletion is inconsequential to the output.
- **Semi-supervised learning:** Towards our objective of frugal forecasting, $\Psi_\Theta(v, t \mid \overrightarrow{\mathcal{G}}_{[1,t]})$ should be computable even if temporal features from past snapshots are not available on $v$ as long as *some* nodes in the graph contain temporal information.

## 3 FRIGATE: PROPOSED METHODOLOGY

At an individual node level, the proposed problem is a *sequence-to-sequence* regression task, wherein we feed the time-series over past snapshots and forecast the future time-series on the target node. To model this problem, we use an *Encoder-Decoder* architecture as outlined in Fig. 2. To jointly capture the spatio-temporal dependency among nodes, the encoder is composed of a stack of *siamese* Gnns and Lstms, i.e., the weights and architecture are identical across each stack. The number of stacks corresponds to the number of historical snapshots one wishes to train on. Each stack is assigned an index depending on how far back it is from the current time $t$. Due to the siamese architecture, the number of stacks does not affect the number of parameters. In addition, the siamese design allows one to dynamically specify the stack size at inference time in a query-specific manner depending on the amount of data available. For simplicity, we will assume the number of historical snapshots to be the same as the forecasting horizon, which is denoted as $\Delta$.

In the $k^{th}$ stack, where $0 \leq k \leq \Delta$, graph snapshot $\mathcal{G}_{t-\Delta+k}$ is passed through the Gnn to learn node representations. The node representation not only characterises its own state, but also the state in its "neighborhood". The "neighborhood" that affects the time-series of the target node is automatically learned through *Lipschitz-gating*. Each stack of siamese Lstm in the encoder receives two inputs: the node embedding of the target node corresponding to its timestamp and the Lstm output from the previous stack.

The decoder has a stack of siamese Lstms as well. However, the decoder Lstms have separate weights than the encoder Lstms. The stack size is the same as the forecasting horizon $\Delta$. The output of a decoder Lstm is augmented with the *moments* [2] of the sensor value distribution in neighborhood of the target node, which injects a strong prior to the model and elevates its performance. Finally, this augmented representation is passed through an Mlp to predict the time-series value. The entire architecture is trained *end-to-end* with *mean absolute error (MAE)* loss function, as defined in Eq. 1. The $\Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right)$ term in Eq. 1 is computed as:

$$\Psi_\Theta\left(v, t+k \mid \overrightarrow{\mathcal{G}}_{[1,t]}\right) = \hat{y}_v^k \quad (3)$$

where $k \in [1, \Delta]$ and $\hat{y}_v^k$, as depicted in Fig. 2, is the predicted output generated through the $k^{th}$ Lstm in the decoder. The next sections detail these individual components.

### 3.1 Gnn Module of Frigate

In this section, we discuss the architecture of a single stack of Gnn. We use an $L$-layered message-passing Gnn to model node dependencies. In each layer, each node $v$ receives messages from its neighbors. These messages are aggregated and passed through a neural network, typically an Mlp, to construct the representation of $v$ in the next layer.

The simplest aggregation framework would be a MeanPool layer where the aggregated message at $v$ from its neighbors is simply the mean of their representations. However, this approach is too simplistic for the proposed problem. Specifically, a road network is directed, and hence drawing messages only through incoming (or outgoing) edges is not enough. On the other hand, a direction-agnostic message passing scheme by uniformly drawing messages
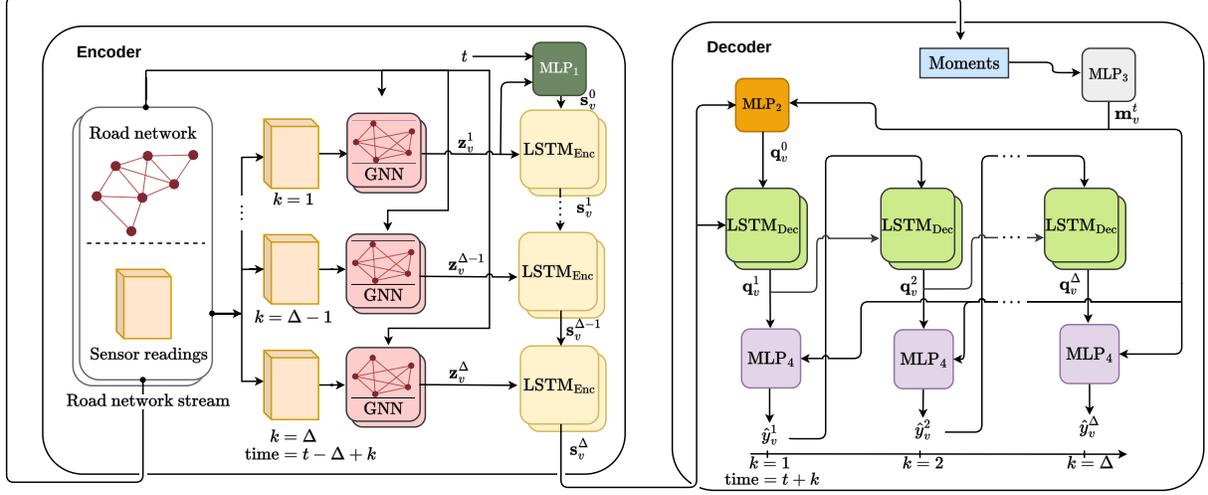
Figure 2: The architecture of FRIGATE. Each of the neural blocks that are in same color are in siamese in nature.

through all incident edges fails to capture traffic semantics. Furthermore, not all edges are equally important. An incoming edge corresponding to an arterial road is likely to have much more impact on an intersection (node) than a road from a small residential neighborhood. Hence, the importance of a road (edge) towards an intersection (node) must be learned from the data, and accordingly, its message should be weighted. To capture these semantics, we formulate our message-passing scheme as follows.

Let $\mathbf{h}_v^\ell$ denote the representation of node $v$ in layer $\ell \in [0, L]$ of the $k^{th}$ GNN. Furthermore, the outgoing and incoming neighbors of $v$ are defined as $\mathcal{N}_v^{out} = \{u \mid (v, u) \in \mathcal{E}\}$ and $\mathcal{N}_v^{in} = \{u \mid (u, v) \in \mathcal{E}\}$ respectively. Now, $h_v^{\ell+1}$ is constructed as follows:

$$\mathbf{h}_v^{\ell+1} = \sigma_1 \left( \mathbf{W}_1^\ell \left( h_v^\ell \parallel \mathbf{m}_v^{\ell,out} \parallel \mathbf{m}_v^{\ell,in} \right) \right), \text{ where} \tag{4}$$

$$\mathbf{m}_v^{\ell,out} = \text{AGGR}^{out} \left( \left\{ \mathbf{h}_u^\ell \mid u \in N_v^{out} \right\} \right) \tag{5}$$

$$\mathbf{m}_v^{\ell,in} = \text{AGGR}^{in} \left( \left\{ \mathbf{h}_u^\ell \mid u \in N_v^{in} \right\} \right) \tag{6}$$

$\parallel$ represents the *concatenation* operation, $\mathbf{W}_1^\ell \in \mathbb{R}^{3d \times d}$ is a learnable weight matrix; $d$ is the dimension of node representations. More simply, we perform two separate aggregations over the messages received from the incoming and outgoing neighbors. The aggregated vectors are concatenated and then passed through a linear transformation followed by non-linearity through an activation function $\sigma_1$. In our implementation, $\sigma_1$ is ReLU. By performing two separate aggregations over the incoming and outgoing edges and then concatenating them, we capture both directionality as well as the topology.

The aggregation functions perform a weighted summation, where the weight of a message-passing edge is learned through *sigmoid gating* over its *positional embedding*. More formally,

$$\text{AGGR}^{out} \left( \left\{ \mathbf{h}_u^\ell \mid u \in N_v^{out} \right\} \right) = \sum_{\forall u \in N_v^{out}} \beta_{v,u} \mathbf{h}_u^\ell \tag{7}$$

$$\text{AGGR}^{in} \left( \left\{ \mathbf{h}_u^\ell \mid u \in N_v^{in} \right\} \right) = \sum_{\forall u \in N_v^{in}} \beta_{u,v} \mathbf{h}_u^\ell \tag{8}$$

$$\beta_{v_i,v_j} = \frac{1}{1 + e^{-\omega_{v_i,v_j}}}, \text{where} \tag{9}$$

$$\omega_{v_i,v_j} = \mathbf{w}^\ell \cdot \mathbf{L}_{v_i,v_j} + b, \text{where} \tag{10}$$

$$\mathbf{L}_{v_i,v_j} = \sigma_2 \left( \mathbf{W}_\mathbf{L}^\ell \left( \left( \mathbf{w}_\delta^T \cdot \delta(v_i, v_j) \right) \parallel \mathbf{L}_{v_i} \parallel \mathbf{L}_{v_j} \right) \right) \tag{11}$$

Here, $\mathbf{L}_v$ denotes the positional embedding of node $v$, whose construction we will discuss in Section 3.1.1. Intuitively, the positional embedding represents the location of a node such that if two nodes are within close proximity in terms of shortest path distance, then their positional embeddings should also be similar. Since, $\mathbf{L}_{v_i,v_j}$ is a function of the positional embeddings of its two endpoints and the spatial distance between them, $\mathbf{L}_{v_i,v_j}$ may be interpreted as the positional embedding of the edge $(v_i, v_j)$. The importance of an edge is computed by passing its positional representation through an MLP (Eq. 11), which is subsequently converted to a scalar (Eq. 10). Finally, the scalar is passed through a sigmoid gate (Eq. 9) to obtain its weight. Here, $\mathbf{w}_\delta \in \mathbb{R}^{d_\delta}$, $\mathbf{W}_\mathbf{L}^\ell \in \mathbb{R}^{2d_L + d_\delta \times d_{L_e}}$ and $\mathbf{w}^\ell \in \mathbb{R}^{d_{L_e}}$ are learnable weight parameters; $d_L$ and $d_{L_e}$ are the dimensionality of the positional embeddings of nodes and edges respectively, and $d_\delta$ is the dimension of the projection created over edge distance $\delta_{v_i,v_j}$.[1] $\sigma_2$ in Eq. 11 is an activation function to apply non-linearity. $b$ in Eq 10 is the learnable bias term.

We now discuss the rationale behind this design. In a road network, the distribution of traffic over edges resemble a *power-law* (Fig. 3). While arterial roads carry a lot of traffic and therefore is a strong determinant of the state of the downstream roads, local roads, which forms the majority, have less impact. Hence, we need to learn this importance of roads from the data. Furthermore, we would like the GNN to be deep without *over-smoothing* or *over-squashing* the

---

[1] We project the edge distance to a higher dimensional representation since otherwise the significantly larger number of dimensions allocated to positional embeddings of the endpoints may dominate the single dimension allocated for edge distance.
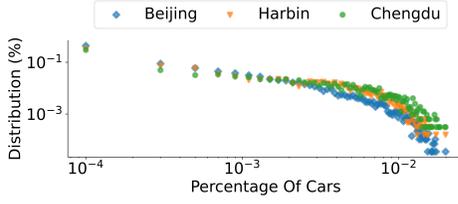
**Figure 3: The distribution of cars flowing through the edges in the cities of Beijing, Chengdu and Harbin.**

representations [5]. Sigmoid gating serves these dual requirements. First, it assigns a weight to every edge to determine its importance. Second, it enables us to go deep since as more layers are added, we receive messages only from the edges with high weights and therefore avoid over-squashing. Semantically, the roads with high weights should correspond to the arterial roads. Note that unlike Graph attention networks [25], where edges compete among each other for importance (due to normalization in SOFTMAX), we do not have that effect in a Sigmoid gating. This is more natural for road networks since the magnitude of representations on a busy intersection should be higher than an intersection with no incoming arterial roads. In addition, as we will formally prove in Section 3.4, Sigmoid gating provides higher expressive power. Finally, it is worth noting that the attention weights are directional since it is a function of the positional edge embedding where $\mathbf{L}_{v_i,v_j} \neq \mathbf{L}_{v_j,v_i}$.

**Initial embedding of nodes:** Let $t$ be the current time and $\Delta$ the number of historical snapshots we are training on. Thus, we will have $\Delta$ stacks of GNN where the $k^{th}$ GNN corresponds to snapshot of time $t - \Delta + k$, $k \in [0, \Delta]$ (Recall Fig. 2). The initial embedding of all nodes $v \in \mathcal{V}_{t-\Delta+k}$ in the $k^{th}$ GNN is defined as:

$$\mathbf{h}_v^{0,k} = \left( \mathbf{w}_\tau^T \cdot \tau_{t-\Delta+k}^v \right) \| \mathbf{L}_v \tag{12}$$

$\mathbf{w}_\tau$ is a learnable vector to map the sensor value to a higher dimensional space. $\mathbf{L}_v$ is the positional embedding of $v$. Note we do not use the stack index in the notation for $\mathbf{h}_v^\ell$ (Eq. 4) since all computations are identical except the time-dependent input in Eq. 12. The final output after the $L^{th}$ layer in the $k^{th}$ GNN is denoted as $\mathbf{z}_v^k = \mathbf{h}_v^{L,k}$.

*3.1.1 Positional Embeddings.* The simplest approach to encode the position of a node is to embed its latitude and longitude into a higher-dimensional representation. This embedding, however, would not reflect the constraints imposed by the road network. Hence, to learn network-induced positional embeddings, we use *Lipschitz embeddings*.

**Definition 3** (Lipschitz Embedding). *Let $\mathcal{A} = \{a_1, \cdots, a_m\} \subseteq \mathcal{V}$ be a randomly selected subset of nodes. We call them anchors. Each node $v$ is embedded into an $m$-dimensional feature vector $\mathbf{L}_v$ where $\mathbf{L}_v[i] = \frac{sp(a_i,v)+sp(v,a_i)}{2}$, where $sp(u,v)$ is the shortest path distance from $u$ to $v$.*

The efficacy of the Lipschitz embedding lies in how well it preserves the shortest path distances among nodes. A well-accepted convention to measure the preservation of distances between the original space and the transformed space is the notion of *distortion* [16].

**Definition 4** (Distortion). *Let $X$ be a set of points embedded in a metric space with distance function $d_X$. Given a function $f : X \to Y$*

that embeds objects in $X$ from the finite metric space $(X, d_X)$ into another finite metric space $(Y, d_Y)$. We define:

$$\text{expansion}(f) = \max_{x_1, x_2 \in X} \frac{d_Y(f(x_1), f(x_1))}{d_X(x_1, x_2)} \tag{13}$$

$$\text{contraction}(f) = \max_{x_1, x_2 \in X} \frac{d_X(x_1, x_2)}{d_Y(f(x_1), f(x_2))} \tag{14}$$

*The distortion of an embedding $f$ is defined as the product of its expansion and contraction. If the distortion is $\alpha$, this means that $\forall x, y \in X$, $\frac{1}{\alpha} d_X(x, y) \leq d_Y(f(x), f(y)) \leq d_X(x, y)$.*

In the context of a road network, $X = \mathcal{V}$, and $d_X(x_1, x_2)$ is the average two-way shortest path distance between $x_1$ and $x_2$ as defined in Def. 3. Furthermore, the mapping $f : X \to Y$ is simply the Lipschitz embedding of $X$. Now, to define $d_Y(f(x_1), f(x_2))$, along with dimensionality $m$, we use *Bourgain's Theorem* [16].

**Theorem 1** (Bourgain's Theorem [16]). *Given any finite metric space $(X, d_X)$, there exists a Lipschitz embedding of $(X, d_X)$ into $\mathbb{R}^m$ under any $L_p$ norm such that $m = O(\log^2 n)$ and the distortion of the embedding is $O(\log n)$, where $n = |X| = |\mathcal{V}|$.*

Based on Bourgain's theorem, if we can show that $d_X(u, v) = \frac{sp(u,v)+sp(v,u)}{2}$, as defined in Def. 3, is metric, then choosing $m = O(\log^2 |\mathcal{V}|)$ anchors and any $L_p$ distance in the embedded space would provide a distortion of $O(\log |\mathcal{V}|)$. We next show that $d_X(u, v)$ is indeed a metric.

**Lemma 1.** $d_X(u, v) = \frac{sp(u,v)+sp(v,u)}{2}$ *is a metric.*

PROOF. *See App. E.*

The exact algorithm to choose the anchors is described in [16]. Note that we use the same set of anchors across all snapshots. Since the topology may change across snapshots, the positional embeddings are time-varying as well. Although unlikely, it is possible for an anchor node to get deleted in a particular snapshot. In such a scenario, we denote the distance corresponding to this dimension as $\infty$ for all nodes.

### 3.2 Encoding Temporal Dependency

To encode long-term temporal dependencies over the node representations, we use an LSTM encoder. Specifically, the final-layer outputs of the GNN in the $k^{th}$ stack, denoted as $\mathbf{z}_v^k$, is fed to the $k^{th}$ LSTM stack. In addition, the $k^{th}$ LSTM also receives the hidden state of the $(k-1)^{th}$ LSTM (Recall Fig. 2). Mathematically, the output of the $k^{th}$ LSTM on node $v$, denoted as $\mathbf{s}_v^k \in \mathbb{R}^{d_{L_{enc}}}$, is computed as follows.

$$\mathbf{s}_v^k = \begin{cases} \text{LSTM}_{Enc}\left(\mathbf{s}_v^{k-1}, \mathbf{z}_v^k\right) & \text{if k>1} \\ \text{LSTM}_{Enc}\left(\text{MLP}_1\left(t, \mathbf{z}_v^k\right), \mathbf{z}_v^k\right) & \text{if k=1} \end{cases} \tag{15}$$

Since $\mathbf{s}_v^0$ is undefined for the first LSTM, i.e., $k = 1$, we make it a learnable vector through the MLP1. The output of the final LSTM stack corresponding to current time $t$ (See Fig. 2), feeds into the decoder. $d_{L_{enc}}$ is the dimension of the LSTM representations.

Lines 1–6 in Alg. 1 summarize the encoder component. Given stream $\overrightarrow{\mathcal{G}}$, current time $t$, and a target node $v$, we extract the subset $\overrightarrow{\mathcal{G}}_{[t-\Delta,t]} \subseteq \overrightarrow{\mathcal{G}}$ where $k$ denotes the number of GNN-LSTM stacks. The processing starts at $\mathcal{G}_{t-\Delta}$ where the $k^{th}$ GNN stack embeds $v$

---

**Algorithm 1** FRIGATE forward pass

---

**Require:** FRIGATE, stream $\overrightarrow{\mathcal{G}}$, target node $v$, current timestamp $t$, prediction horizon $\Delta$
**Output:** Predicted sensor values $\hat{y}_{t+1}^v \cdots \hat{y}_{t+\Delta}^v$
1: $\mathbf{s}_v^0 \leftarrow \text{MLP}_1(t, \mathbf{z}_v^1)$
2: $k \leftarrow 1$
3: **for all** $\mathcal{G} \in \overrightarrow{\mathcal{G}}_{[t-\Delta, t]}$ **do**
4: $\quad \mathbf{z}_v^k \leftarrow \text{GNN}(v)$
5: $\quad \mathbf{s}_v^k \leftarrow \text{LSTM}_{\text{Enc}}(\mathbf{s}_v^{k-1}, \mathbf{z}_v^k)$
6: $\quad k \leftarrow k+1$
7: $\mathbf{m}_v^t \leftarrow \text{MLP}_3\left(moments\left(\left\{\tau_u^{t'} \neq \varnothing | t' \in [t-\Delta, t], (u,v) \text{ or } (v,u) \in \mathcal{E}^{t'}\right\}\right)\right)$
8: $\hat{y}_0 \leftarrow \text{MLP}_2(\mathbf{s}_v^\Delta, \mathbf{m}_v^t)$
9: $\mathbf{q}_v^0 \leftarrow \mathbf{s}_v^\Delta$
10: **for all** $k \in [1 \ldots \Delta]$ **do**
11: $\quad \mathbf{q}_v^k \leftarrow \text{LSTM}_{\text{Dec}}(\mathbf{q}_v^{k-1}, \hat{y}_v^{k-1})$
12: $\quad \hat{y}_v^k \leftarrow \text{MLP}_4(\mathbf{q}_v^k, \mathbf{m}_v^t)$
13: Re-index $\hat{y}_k^v \mapsto \hat{y}_{t+k}^v$
14: **return** $\hat{y}_{t+1}^v, \ldots, \hat{y}_{t+\Delta}^v$

---

into $\mathbf{z}_v^k$. Next, $\mathbf{z}_v^k$ is passed to the $k^{th}$ LSTM. This completes one stack of computation. Iteratively, we move to the $(k+1)^{th}$ stack and the same operations are repeated till $k = \Delta$, after which $\mathbf{s}_v^\Delta$ (Eq. 15) is fed to the decoder.

## 3.3 Decoder

The decoder is composed of a stack of $\Delta$ (forecasting horizon) siamese LSTMs. Assuming $t$ to be the current time, the output of the $k^{th}$ LSTM in the decoder corresponds to the predicted sensor reading at time $t+k$. Each LSTM receives two inputs: (1) the predicted sensor value in the previous timestamp denoted as $\hat{y}_v^{k-1} \in \mathbb{R}$ and (2) the hidden state of the previous LSTM, denoted as $\mathbf{q}_v^{k-1} \in \mathbb{R}^{d_{dec}}$. Thus, the output of the $k^{th}$ LSTM is expressed as:

$$\mathbf{q}_v^k = \begin{cases} \text{LSTM}_{Dec}\left(\hat{y}_v^{k-1}, \mathbf{q}_v^{k-1}\right) & \text{if k>1} \\ \text{LSTM}_{Dec}\left(\text{MLP}_2\left(\mathbf{s}_v^\Delta, \mathbf{m}_v^t\right), \mathbf{s}_v^\Delta\right) & \text{if k=1} \end{cases} \quad (16)$$

We have a special case for $k = 1$, since both $\hat{y}_v^0$ and $\mathbf{q}_v^0$ are undefined. Since we assume a setting of partial sensing across nodes, $\tau_v^t$ may not be available and hence cannot be used to substitute $\hat{y}_v^0$. To mitigate this situation, $\mathbf{q}_v^{k-1}$ is replaced with the output $\mathbf{s}_v^\Delta$ of the last LSTM in the encoder and $\hat{y}_v^0$ is estimated through the $\text{MLP}_4$. This MLP takes as input $\mathbf{s}_v^\Delta$ and a representation of the *moments* [2] of the observed sensor values in the neighborhood of $v$, denoted as $\mathbf{m}_v^t$. Formally,

$$\mathbf{m}_v^t = \text{MLP}_3\left(moments\left(\left\{\tau_u^{t'} \neq \varnothing | t' \in [t-\Delta, t], (u,v) \text{ or } (v,u) \in \mathcal{E}^{t'}\right\}\right)\right)$$

Finally, the predicted sensor value $\forall k \in [1, \Delta]$, $\hat{y}_v^k$ is computed as:

$$\hat{y}_v^k = \text{MLP}_4\left(\mathbf{q}_v^k, \mathbf{m}_v^t\right) \quad (17)$$

Note that instead of directly predicting the sensor value from the LSTM hidden state $\mathbf{q}_v^k$, we augment this information with statistical information $\mathbf{m}_v^t$ on the time-series at $v$. This provides a strong inductive bias to the model and elevates its performance, which we will substantiate empirically during our ablation study in Section 4. Lines 7–12 in Alg. 1 summarize the computations in the decoder.

## 3.4 Theoretical Characterization of FRIGATE

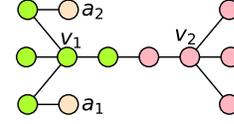**Fact 1.** *FRIGATE is inductive and permuation-invariant.*



**Figure 4: Sample graph to illustrate expressivity of GNN in FRIGATE.**

DISCUSSION. As outlined in Sections 3.1-3.3 and summarized in Table C.1 in Appendix, the parameter-size is independent of the network size, number of snapshots being used for training, forecasting horizon, and the temporal granularity across snapshots. Furthermore, since FRIGATE uses sigmoid-weighted sum-pool, it is permutation invariant to node set. Consequently, FRIGATE can inherently adapt to changes in topology or forecast on unseen nodes without the need to re-train from scratch. As outlined in Table 1 and discussed in Sec. 1.1, majority of existing works on spatio-temporal network-constrained forecasting do not satisfy the needs of being inductive and permutation-invariant. □

We now focus on the expressive power of FRIGATE. Message-passing GNNs that aggregate messages from their neighbors are no more powerful than 1-WL [30]. Hence, given nodes $v$ and $u$, if their $\ell$-hop neighborhoods are isomorphic, they will get identical embeddings. Thus, the position of a node plays no role. Methodologies such as DCRNN [14], STGCN [35], or STGODE [6], that build on top of these message-passing GNNs will consequently inherit the same limitation. FRIGATE does not suffer from this limitation.

**Lemma 2.** *FRIGATE can distinguish between isomorphic neighborhoods based on positioning.*

PROOF. Consider nodes $v_1$ and $v_2$ in Fig. 4. If we use a 1-layered message-passing GNN, then their embeddings would be identical since the 1-hop neighborhoods (color coded in green and pink) are isomorphic. Hence, DCRNN, STGCN, or STGODE would not be able to distinguish between them. FRIGATE augments initial node features with Lipschitz embeddings. Assuming $a_1$ and $a_2$ to be the anchor nodes, $v_1$ would have a Lipschitz embedding of $[2, 2]$ vs. $[5, 5]$ for $v_2$. Hence, FRIGATE would generate different embeddings for $v_1$ and $v_2$. □

As we will show in our ablation study, removing positional embeddings have a significant impact on the performance.

**Lemma 3.** *FRIGATE is at least as powerful as 1-WL.*

PROOF. *See App. F.*

**Corollary 1.** FRIGATE *is strictly more expressive than DCRNN, STGCN, and STGODE.*

PROOF. This follows naturally by combining Lemmas 2 and 3 since FRIGATE retains the 1-WL expressivity of DCRNN, STGCN, and STGODE, while also being capable of distinguishing between isomorphic topologies through positional embeddings. □

We note that while positional embeddings have been used in the context of GNNs [21, 34], they do not provide 1-WL expressivity since messages are exchanged only among anchor nodes. The methodology proposed in this work is therefore unique.

**Table 2: (a) Data statistics and (b) model size.**

**(a) Details about the datasets**

| Dataset | #Nodes | #Edges | Geographical Area (km$^2$) | Mean | Std | Default Seen Percent | #Trips |
|---|---|---|---|---|---|---|---|
| Beijing [15] | 28465 | 64478 | 16,411 | 3.22 | 9.09 | 5% | 785709 |
| Chengdu [1] | 3193 | 7269 | 14,378 | 4.54 | 7.19 | 50% | 1448940 |
| Harbin [13] | 6235 | 15205 | 53,068 | 77.85 | 124.11 | 30% | 659141 |

**(b) Number of parameters**

| Model | #Parameters |
|---|---|
| Frigate | 198985 |
| DCRNN | 372353 |
| STNN | 313906 |
| GraphWavent | 247660 |
| STGODE | 558582 |
| LocaleGN | 333636 |

## 4 EXPERIMENTS

In this section, we benchmark Frigate and establish:

- **Prediction Accuracy:** Frigate is superior compared to baselines for the task of spatio-temporal forecasting on road networks.
- **Ablation study:** Through extensive ablation studies, we illuminate the significance of each component of Frigate and provide insights into the critical role they play in the overall performance.
- **Robustness:** We test the limits of our model to frugality in sensing, graph structure modifications at inference time and resilience to non-uniform temporal granularity across snapshots.

The codebase of Frigate and datasets are available at https://github.com/idea-iitd/Frigate.

### 4.1 Datasets

We use three real-world datasets collected by tracking the GPS trajectory of taxis. As summarized in Table 2a, the three datasets correspond to the cities of Beijing, Chengdu and Harbin [8]. We map-match [31] the raw trajectories to their corresponding road networks obtained from Openstreetmap [22]. The traffic is aggregated into buckets of 5 minutes and the sensor value corresponds to the number of vehicles passing through each node in a 5-minute interval. From the entire node-set, we select a subset of nodes uniformly at random as the "seen" nodes with installed sensors. Inference and loss calculation is performed only on the unseen nodes. The default percentage of nodes we retain as "seen" is shown in Table 2a.

### 4.2 Experimental Setup

- **Computational engine:** We use a system running on Intel Xeon 6248 processor with 96 cores and 1 NVIDIA A100 GPU with 40GB memory for our experiments.
- **Forecasting horizon:** In all the experiments, we choose to predict one hour of traffic information from one hour of historical traffic information. Concretely, since in the datasets the duration between two snapshots is 5 minutes, we set $\Delta = 12$. Later, we also simulate what happens if the duration between snapshots is irregular, but even then we have past one hour of data and make predictions of traffic condition up to one hour in the future.
- **Evaluation metric:** We use mean absolute error (MAE) as the primary metric across all methods (lower is better). We also use sMAPE and RMSE for our key results in Table 3. In addition, we report a 95% confidence interval around the mean by using bootstrapping on the distribution of the metric being used on the "unseen" nodes.
- **Training setting:** We use a 70%-20%-10% train-val-test split for training. This split is on the time dimension. So, to be explicit, during training 70% of the total time series on seen nodes are used. Validation and test are both on different parts of the data both from node perspective and timestep perspective. We stop the training of

the model if it does not improve the validation loss for more than 15 epochs. Further, we have varied percentage of training nodes in various experiments to check the fidelity of the FRIGATE, as indicated in the results.

- **Baselines:** We consider DCRNN [14], STGCN [35], LocaleGN [11] and STNN [32] as our baselines. We adapt GraphWavenet [29] and STGODE [6] to our setting by removing node embedding matrices or Dynamic Time Warping based graph computation so that they support the triple objectives outlined in Table 1. For all six baselines, we obtain the official code-base released by the authors.
- **Parameter settings:** We use 16 anchor nodes to calculate Lipschitz embeddings for all graphs. We use 10 layers of Gnn in Frigate.

### 4.3 Inference Accuracy

In Table 3, we present the MAE obtained by Frigate and all other baselines on all three datasets. We observe that Frigate consistently outperforms all baselines. On average, the MAE in Frigate is more than 25% lower than the closest baseline. Among the baselines, LocaleGN and STGODE perform the best, followed by DCRNN, followed by GraphWavenet, followed by STNN and then STGCN. We further note that STGCN fails to train on Harbin and Beijing even after 48 hours (the largest dataset evaluated in STGCN [35] was of 1026 nodes and it only considered weekday traffic). Likewise, GraphWavenet also fails to train on Harbin and Beijing (the largest dataset evaluated in GraphWavenet [29] was of 325 nodes). Now, to better contextualize the results, let us compare the MAE with the standard deviation of sensor values reported in Table 2a. We observe a clear correlation of MAE with the standard deviation, which explains why all techniques perform comparatively poorer in Harbin. Due to the substantially inferior performance of STGCN and its inability to scale on large datasets, subsequent experiments only consider DCRNN, STNN, STGODE, and LocaleGN as baselines.

To further diagnose the performance pattern among the benchmarked techniques, we segregate the nodes into three buckets based on the number of cars going through them and plot the distribution of errors within these three buckets. Fig. 5 presents the results. Here, "High", "Medium" and "Low" corresponds to the top-33 percentile, 33-66 percentile, and bottom-33 percentile, respectively. We derive three key insights from this experiments. First, across all techniques, the MAE reduces as we move towards nodes handling lower traffic. This is natural, since these nodes have low variation in terms of traffic volume. In contrast, high-traffic nodes undergo more fluctuation depending on the time of the day, weekday vs. weekends, etc. Second, we note that Frigate performs better than all other models on high frequency nodes, good on medium frequency nodes, and worse than the rest in low frequency nodes. Overall, it may be argued that doing better in high and medium-frequency nodes
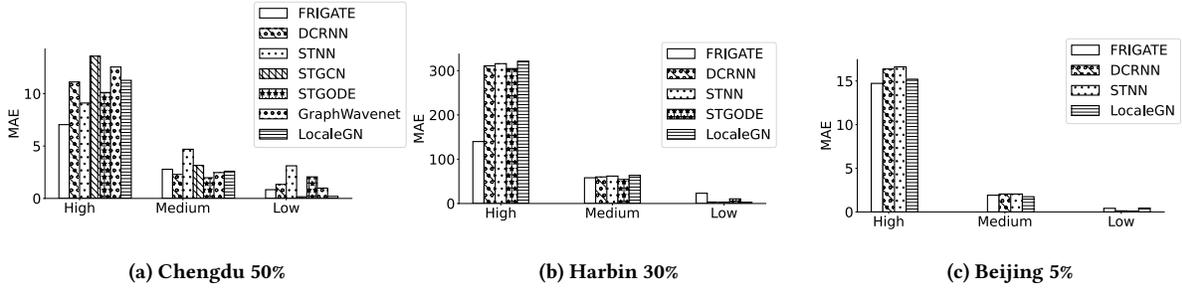
(a) Chengdu 50%                    (b) Harbin 30%                    (c) Beijing 5%

**Figure 5: Comparison of the distribution of error over nodes grouped by frequency of cars moving through them.**

are more important since they handle a large majority of the traffic, where FRIGATE mostly outperforms other techniques. Finally, one of the key reasons of under-performance in DCRNN, STNN, and STGODE is that they over-smooth the high-frequency nodes into inferring lower values as the majority of the nodes are of low-frequency (recall power-law distribution from Fig. 5). Positional embeddings and deep layering through gating enables FRIGATE to avoid over-smoothing. As we established in Section 3.4, these design choices make FRIGATE provable more expressive. We further investigate the impact of these two components in our ablation study in Section 4.5.

## 4.4 Impact of Volume of "Seen" Nodes

As in any machine learning task, we expect the accuracy to improve with larger volume of seen nodes. This increases training data, as well as induces a higher likelihood of an unseen node being close to a seen node. Towards that objective, we vary the number of nodes that are "seen" by the models from 10% to 90% of the total number of nodes in the respective road networks and measure MAE against forecasting horizon. Fig. 6 presents the results. We observe that FRIGATE performs significantly better because of more informative priors and inherent inductivity. Furthermore, the baselines

**Table 3: Performance comparison of different approaches for traffic volume prediction. OOT stands for Out of Time and OOM stands for Out of Memory. The lowest MAE (best accuracy), lowest sMAPE and lowest RMSE are highlighted in bold for each dataset.**

| Model | Dataset | MAE | sMAPE | RMSE |
|---|---|---|---|---|
| FRIGATE | Chengdu | **3.547 ± 0.20** | **131.17 ± 3.39** | **5.93 ± 0.31** |
| STNN | | 5.633 ± 0.18 | 199.86 ± 0.00 | 9.50 ± 0.50 |
| DCRNN | | 4.893 ± 0.26 | 138.17 ± 2.07 | 8.10 ± 0.48 |
| LocaleGN | | 4.597 ± 0.22 | 192.80 ± 0.19 | 8.57 ± 0.33 |
| STGODE | | 4.693 ± 0.19 | 147.12 ± 2.25 | 7.72 ± 0.39 |
| STGCN | | 5.712 ± 0.33 | 198.38 ± 0.00 | 9.47 ± 0.50 |
| GraphWavenet | | 5.308 ± 0.30 | 147.77 ± 1.81 | 8.89 ± 0.43 |
| FRIGATE | Harbin | **73.559 ± 2.04** | **93.00 ± 1.89** | **105.01 ± 3.51** |
| STNN | | 126.305 ± 4.54 | 199.68 ± 0.02 | 206.91 ± 6.93 |
| DCRNN | | 124.109 ± 5.06 | 185.73 ± 0.68 | 204.75 ± 6.36 |
| LocaleGN | | 128.674 ± 5.10 | 169.54 ± 1.43 | 207.36 ± 6.58 |
| STGODE | | 122.493 ± 4.97 | 152.28 ± 1.27 | 198.81 ± 6.74 |
| STGCN | | OOT | OOT | OOT |
| GraphWavenet | | OOM | OOM | OOM |
| FRIGATE | Beijing | **5.651 ± 0.11** | 169.73 ± 0.45 | 11.87 ± 0.30 |
| STNN | | 6.215 ± 0.12 | 199.99 ± 0.01 | 12.82 ± 0.27 |
| DCRNN | | 6.122 ± 0.14 | 195.56 ± 0.15 | 12.63 ± 0.24 |
| LocaleGN | | 5.759 ± 0.13 | **172.73 ± 0.33** | 12.10 ± 0.03 |
| STGODE | | OOM | OOM | OOM |
| STGCN | | OOT | OOT | OOT |
| GraphWavenet | | OOM | OOM | OOM |

need considerably more data to reach the same performance. For Chengdu, DCRNN trained on 50% of the graph beats our model trained on only 10% of the graph while on Harbin, none of the models ever even surpass our model trained at only 10% of the graph. And as the data increases, the gap between our model and the baselines increases. This economic use of data without losing the ability to generalize is a desirable property that FRIGATE contains.

Another interesting trend we note is that the performance slightly deteriorates in the baselines from 70% to 90%. It is hard to pinpoint the exact reason. We hypothesize two factors as possible reasons. First, with higher volumes of training data, there is a stronger chance of over-fitting. In addition, we note that the confidence interval expands as the percentage of seen nodes increases since the sample size of test nodes decreases. This trend, which is consistent with statistical theory, means that at low volumes of test sets, the results have higher variability. FRIGATE does not suffer from this trend, which means it does not overfit. It is well-known in machine learning theory, that tendency to overfit is correlated to the number of parameters. In this regard, we draw attention to Table 2b, which shows that FRIGATE has almost 50% and 33% smaller parameter set than DCRNN and STNN, respectively. Also, FRIGATE utilizes moments as a robust inductive bias, reducing overfitting risks.
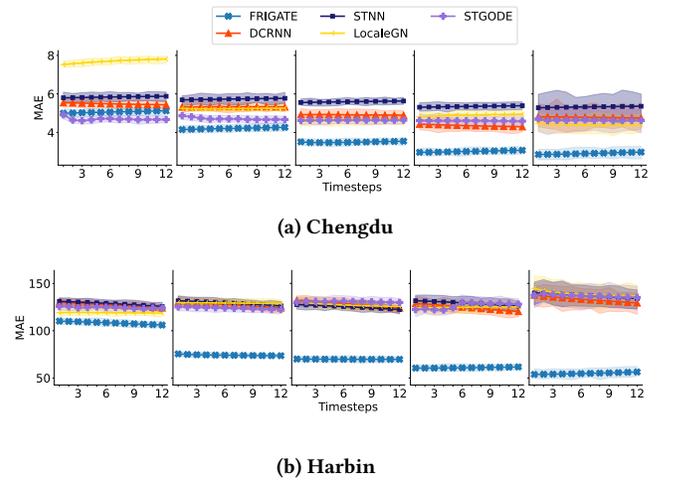


(a) Chengdu



(b) Harbin

**Figure 6: From left to right, we progressively increase the volume of seen nodes and measure MAE in (a) Chengdu and (b) Harbin. The envelopes signify the 95% confidence bound of the MAEs.**
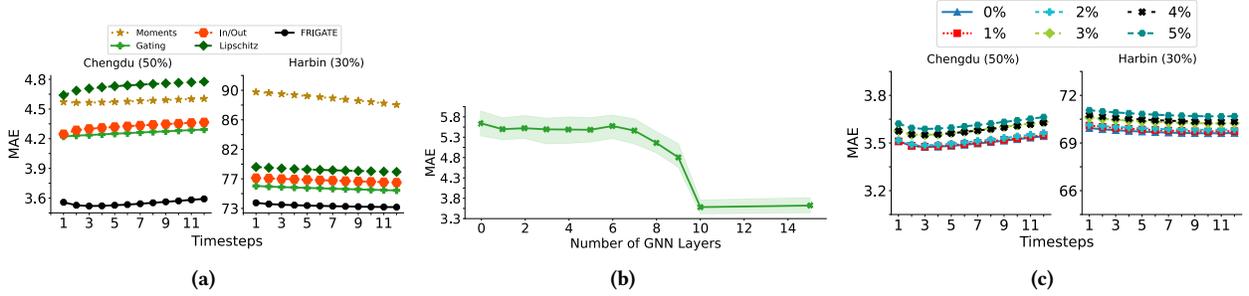
**Figure 7: (a) Ablation study. (b) Impact of Gnn layer variations on Chengdu 50% (c) Model resilience to inference time changes on road network topology.**

## 4.5 Ablation Study

In our ablation study, we systematically turn off various components of Frigate and measure the impact on MAE. Fig. 7a presents the results on Chengdu and Harbin.

**Gating:** Here, we replace the gated convolution with a GraphSAGE layer [7], leaving rest of the architecture intact. We see a clear increase in MAE indicating the importance of gated convolution layer.

**Positional embeddings:** To understand its utility on performance, we remove Lipschitz embeddings as features of the node, and thus remove it as an input to the calculation of gating. We see a significant drop in performance, with Chengdu being more pronounced where Lipschitz has the highest drop. This result empirically demonstrates the value of positional embeddings in time-series forecasting on road networks.

**In-Out aggregation:** In Frigate, we separately aggregate messages from incoming and outgoing neighbors. Here, we measure the impact of aggregating into a single vector. We observe an increase in MAE, showcasing the need for direction-based aggregation.

**Moments:** We remove the moments from the model and keep everything else in the architecture the same. We see a big drop in performance on both datasets, with Harbin being more pronounced. This establishes the importance of having an informative prior.

**Number of GNN layers:** To understand the effect of Gnn layers on performance of Frigate, we vary this parameter while training on Chengdu 50% dataset. As visible in Fig. 7b, the performance saturates at 10, which we use as our default value.

## 4.6 Robustness and Resilience

In this section, we analyze the robustness and resilience of Frigate to changes to network topology and irregular topological sensing. In these experiment, we train Frigate on the original datasets, and then slightly perturb the topology or temporal granularity. We evaluate inference performance on this perturbed dataset. Note that we do not re-train the model after perturbations.

**Topology change:** To simulate the real-world situation where the road network might change due to road blocks or opening of new roads, we first select the volume of perturbations to be introduced. Assuming this to be $X\%$, we change the network by randomly dropping a maximum of $\frac{X}{2}\%$ of the edges. In addition, we create an equal number edges among unconnected nodes whose distance is within a threshold. The distances for these edges are sampled, with replacement, from the original distribution of edge distance. We then vary $X$ and measure the impact on MAE. As visible in Fig. 7c,

**Table 4: Impact of irregular temporal granularity on MAE.**

| Time series | Dataset | MAE |
|---|---|---|
| Regular (original) | Chengdu (50%) | $3.547 \pm 0.23$ |
| Irregular (perturbed) | Chengdu (50%) | $3.582 \pm 0.17$ |
| Regular (original) | Harbin (50%) | $69.700 \pm 3.21$ |
| Irregular (perturbed) | Harbin (50%) | $69.834 \pm 2.78$ |

Frigate is resilient to changes to the road network with minimal drop in accuracy.

**Irregular temporal sensing:** In this experiment, we try to break our model by changing the time step granularity at test time. We randomly drop 33.33% of snapshots in $\overrightarrow{\mathcal{G}}_{[t-\Delta,t]}$, effectively reducing the time sequence length and also changing the granularity at which the data is captured. The results are tabulated in table 4, where we observe negligible increase in MAE. The results indicate that the model is resilient to changes in the granularity of snapshots.

## 5 CONCLUSION

In this paper, we have proposed a new spatio-temporal Gnn, called Frigate, to forecast network-constrained time-series processes on road networks. We show through experiments on real-world traffic datasets that Frigate significantly outperforms existing baselines. In addition, Frigate is robust to several practical challenges such as partial sensing across nodes, absorb minor updates to road network topology and resilience to irregular temporal sensing. The core competency of Frigate originates from its novel design that incorporates Lipschitz embedding to encode position, sigmoid gating to learn message importance and enable pathways for long-range communication across nodes, directional aggregation of messages, and strong priors in the form of moments of the time-series distribution. In addition, Frigate is built from the ground-up keeping inductivity as a core objective. On the whole, Frigate takes us closer to deployable technology for practical workloads.

## ACKNOWLEDGMENTS

# REFERENCES

[1] [n. d.]. Gaia Open Dataset. gaia.didichuxing.com.
[2] [n. d.]. Moment (Mathematics). https://en.wikipedia.org/wiki/Moment_(mathematics). Accessed: 2023-02-01.
[3] Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. 2020. Adaptive graph convolutional recurrent network for traffic forecasting. *Advances in neural information processing systems* 33 (2020), 17804–17815.
[4] Yuzhou Chen, Ignacio Segovia, and Yulia R Gel. 2021. Z-GCNETs: time zigzags at graph convolutional networks for time series forecasting. In *International Conference on Machine Learning*. PMLR, 1684–1694.
[5] Vijay Prakash Dwivedi, Ladislav Rampášek, Mikhail Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu, and Dominique Beaini. 2022. Long Range Graph Benchmark. In *NeurIPS*.
[6] Zheng Fang, Qingqing Long, Guojie Song, and Kunqing Xie. 2021. Spatial-temporal graph ode networks for traffic flow forecasting. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 364–373.
[7] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems* 30 (2017).
[8] Jayant Jain, Vrittika Bagadia, Sahil Manchanda, and Sayan Ranu. 2021. NeuroMLR: Robust & Reliable Route Recommendation on Road Networks. *Advances in Neural Information Processing Systems* 34 (2021), 22070–22082.
[9] Manas Joshi, Arshdeep Singh, Sayan Ranu, Amitabha Bagchi, Priyank Karia, and Puneet Kala. 2021. Batching and matching for food delivery in dynamic road networks. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2099–2104.
[10] Shiyong Lan, Yitong Ma, Weikang Huang, Wenwu Wang, Hongyu Yang, and Pyang Li. 2022. DSTAGNN: Dynamic Spatial-Temporal Aware Graph Neural Network for Traffic Flow Forecasting. In *International Conference on Machine Learning*. PMLR, 11906–11917.
[11] Mingxi Li, Yihong Tang, and Wei Ma. 2022. Few-Sample Traffic Prediction With Graph Networks Using Locale as Relational Inductive Biases. *IEEE Transactions on Intelligent Transportation Systems* (2022).
[12] Mengzhang Li and Zhanxing Zhu. 2021. Spatial-temporal fusion graph neural networks for traffic flow forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 4189–4196.
[13] Xiucheng Li, Gao Cong, Aixin Sun, and Yun Cheng. 2019. Learning travel time distributions with deep generative model. In *The World Wide Web Conference*. 1017–1027.
[14] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. 2018. Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting. In *International Conference on Learning Representations (ICLR '18)*.
[15] Jing Lian and Lin Zhang. 2018. One-month beijing taxi GPS trajectory dataset with taxi IDs and vehicle status. In *Proceedings of the First Workshop on Data Acquisition To Analysis*. 3–4.
[16] Nathan Linial, Eran London, and Yuri Rabinovich. 1995. The geometry of graphs and some of its algorithmic applications. *Combinatorica* 15 (1995), 215–245.
[17] Dachuan Liu, Jin Wang, Shuo Shang, and Peng Han. 2022. Msdr: Multi-step dependency relation networks for spatial temporal forecasting. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1042–1050.
[18] Sourav Medya, Sayan Ranu, Jithin Vachery, and Ambuj Singh. 2018. Noticeable network delay minimization via node upgrades. *Proceedings of the VLDB Endowment* 11, 9 (2018), 988–1001.
[19] Shubhadip Mitra, Sayan Ranu, Vinay Kolar, Aditya Telang, Arnab Bhattacharya, Ravi Kokku, and Sriram Raghavan. 2015. Trajectory aware macro-cell planning for mobile users. In *2015 IEEE Conference on Computer Communications (INFOCOM)*.

[20] IEEE, 792–800.
Shubhadip Mitra, Priya Saraf, Richa Sharma, Arnab Bhattacharya, Sayan Ranu, and Harsh Bhandari. 2017. Netclus: A scalable framework for locating top-k sites for placement of trajectory-aware services. In *2017 IEEE 33rd International Conference on Data Engineering (ICDE)*. IEEE, 87–90.
[21] Sunil Nishad, Shubhangi Agarwal, Arnab Bhattacharya, and Sayan Ranu. 2021. GraphReach: Position-Aware Graph Neural Network using Reachability Estimations. *IJCAI* (2021).
[22] OpenStreetMap contributors. 2017. Planet dump retrieved from https://planet.osm.org . https://www.openstreetmap.org.
[23] Shiladitya Pande, Sayan Ranu, and Arnab Bhattacharya. 2017. SkyGraph: Retrieving regions of interest using skyline subgraph queries. *Proceedings of the VLDB Endowment* 10, 11 (2017), 1382–1393.
[24] Chao Song, Youfang Lin, Shengnan Guo, and Huaiyu Wan. 2020. Spatial-temporal synchronous graph convolutional networks: A new framework for spatial-temporal network data forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 914–921.
[25] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. *ICLR* (2018).
[26] Pengkun Wang, Chaochao Zhu, Xu Wang, Zhengyang Zhou, Guang Wang, and Yang Wang. 2022. Inferring intersection traffic patterns with sparse video surveillance information: An st-gan method. *IEEE Transactions on Vehicular Technology* 71, 9 (2022), 9840–9852.
[27] Yang Wang, Yiwei Xiao, Xike Xie, Ruoyu Chen, and Hengchang Liu. 2018. Real-time Traffic Pattern Analysis and Inference with Sparse Video Surveillance Information.. In *IJCAI*. 3571–3577.
[28] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. 2020. Connecting the dots: Multivariate time series forecasting with graph neural networks. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 753–763.
[29] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. 2019. Graph wavenet for deep spatial-temporal graph modeling. *International Joint Conference on Artificial Intelligence* 28 (2019).
[30] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How powerful are graph neural networks? *ICLR* (2019).
[31] Can Yang and Gyozo Gidofalvi. 2018. Fast map matching, an algorithm integrating hidden Markov model with precomputation. *International Journal of Geographical Information Science* 32, 3 (2018), 547–570.
[32] Song Yang, Jiamou Liu, and Kaiqi Zhao. 2021. Space Meets Time: Local Space-time Neural Network For Traffic Flow Forecasting. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 817–826.
[33] Pranali Yawalkar and Sayan Ranu. 2019. Route recommendations on road networks for arbitrary user preference functions. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE, 602–613.
[34] Jiaxuan You, Rex Ying, and Jure Leskovec. 2019. Position-aware Graph Neural Networks. In *ICML*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), Vol. 97. 7134–7143.
[35] Bing Yu, Haoteng Yin, and Zhanxing Zhu. 2018. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *International Joint Conference on Artificial Intelligenc* 27 (2018).
[36] Chuanpan Zheng, Xiaoliang Fan, Cheng Wang, and Jianzhong Qi. 2020. Gman: A graph multi-attention network for traffic prediction. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 1234–1241.
[37] Zhengyang Zhou, Kuo Yang, Wei Sun, Binwu Wang, Min Zhou, Yunan Zong, and Yang Wang. 2023. Towards Learning in Grey Spatiotemporal Systems: A Prophet to Non-consecutive Spatiotemporal Dynamics. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*. SIAM, 190–198.

**Table B.1: Inference times**

| Model | Time(s) |
| --- | --- |
| DCRNN | 0.0081 |
| STGODE | 0.0164 |
| Frigate | 0.0207 |
| STNN | 0.0375 |
| LocaleGN | 0.0492 |

## A  NOTATIONS

**Table A.1: Notations used in the paper**

| Symbol | Meaning |
| --- | --- |
| $\overrightarrow{\mathcal{G}}$ | The road network stream |
| $\mathcal{G}_t$ | Road network snapshot at time $t$ |
| $\overrightarrow{\mathcal{G}}_{[i,j]},\ i \leq j$ | $\{\mathcal{G}_i, \mathcal{G}_{i+1}, \ldots, \mathcal{G}_j\} \subseteq \overrightarrow{\mathcal{G}}$ |
| $\mathcal{V}_t$ | Vertex set at timestep $t$ |
| $\mathcal{V}^{tr}$ | Set of vertices to be trained on |
| $\mathcal{E}_t$ | Edge set at timestep $t$ |
| $\delta$ | *Haversine* distance function |
| $\tau_t$ | Sensor readings $\forall$ nodes at time $t$ |
| $T$ | Number of snapshots of the graph |
| $\Delta$ | Prediction horizon: how many timesteps in future to predict |
| $\Psi_\Theta$ | The model with parameters $\Theta$ |
| $\mathcal{P}$ | Permutation function |
| $\mathcal{N}(v)$ | Set of nodes in the neighborhood of node $v$ |
| $\mathbf{L}_v$ | Positional embedding of node $v$ |
| $\mathbf{L}_{v_i,v_j}$ | Positional embedding of edge $(v_i, v_j)$ |
| $\mathbf{h}_v^\ell$ | Intermediate representation of node $v$ at layer $\ell$ in Gnn |
| $\mathbf{z}_v^k$ | Final representation of node $v$ in $k^{th}$ stack output by Gnn |
| $\mathbf{s}_v^k$ | Hidden state + cell state of LSTM$_{\text{Enc}}$ at step $k$ |
| $\mathbf{m}_v^t$ | Representation of moments of $\mathcal{N}(v)$ at time $t$ |
| $\mathbf{q}_v^k$ | Hidden state + cell state of LSTM$_{\text{Dec}}$ at step $k$ |
| $\hat{y}_v^k$ | Prediction at $k^{th}$ step from LSTM$_{\text{Dec}}$ |
| $\cdot\|\cdot$ | Concatenation operator |
| $sp(u, v)$ | Shortest path between $u$ and $v$ |
| $\|\cdot\|_2$ | L$_2$ norm |

## B  INFERENCE TIME

The inference times of the three models are shown in Table B.1. Frigate is faster than STNN but slower than DCRNN. However, we note that even though Frigate is slower than DCRNN it is still fast enough to be deployed for real-world workloads.

## C  PARAMETER SIZE

**Table C.1: Dimensionality of the parameters used in Frigate.**

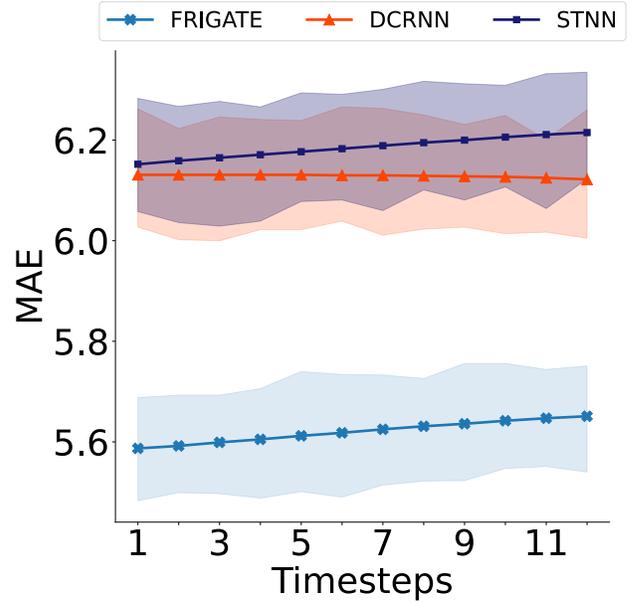| Parameter | Shape |
| --- | --- |
| $\mathbf{W}_1^\ell$ | $\mathbb{R}^{3d \times d}$ |
| $\mathbf{w}_\delta^T$ | $\mathbb{R}^{d_\delta}$ |
| $\mathbf{W}_L^\ell$ | $\mathbb{R}^{2d_L + d\delta \times d_{L_e}}$ |
| $\mathbf{w}^\ell$ | $\mathbb{R}^{d_{L_e}}$ |
| $\mathbf{w}_\tau^T$ | $\mathbb{R}^{d_\tau}$ |
| Mlp$_1$ | $\mathbb{R}^{d_{Lenc}}$ |
| Mlp$_2$ | $\mathbb{R}^{d_{dec}}$ |
| Mlp$_3$ | $\mathbb{R}^{d_{moments}}$ |
| Mlp$_4$ | $\mathbb{R}^{d_{dec} + d_{moments}}$ |

## D  OTHER RESULTS



**Figure D.1: Prediction on Beijing (5%) dataset**

Fig. D.1 shows the MAE of the three models on Beijing (5%) dataset against the prediction horizon.

# E  PROOF OF LEMMA 1

**Lemma 1.** $d_X(u,v) = \frac{sp(u,v)+sp(v,u)}{2}$ is a metric.

PROOF. We need to show **(1)** Symmetry: $d_X(u,v) = d_X(v,u)$, **(2)** Non-negativity: $d_X(u,v) \geq 0$, **(3)** $d_X(u,v) = 0$ iff $u = v$ and **(4)** Triangle inequality: $d_X(u,v) \leq d_X(u,w) + d_X(w,v)$.

We omit the proofs of first three properties since they are trivial. We use proof-by-contradiction to establish triangle inequality. Let us assume $d_X(u,v) > d_X(u,w) + d_X(w,v)$        (18)

$$or, sp(u,v) + sp(v,u) > sp(u,w) + sp(w,u) + sp(w,v) + sp(v,w)$$

From the definition of shortest paths, $sp(u,v) \leq sp(u,w) + sp(w,v)$ and $sp(v,u) \leq sp(v,w) + sp(w,u)$. Hence, Eq. 18 is a contradiction. □

# F  PROOF OF LEMMA 3

**Lemma 3.** FRIGATE is at least as powerful as 1-WL.

PROOF. GIN [30] is as powerful as 1-WL [30]. This power is induced by the SUM-POOL aggregation since sum-pool is *injective* function, i.e., two separate aggregation over messages would be identical if and only if the input messages are identical [2]. FRIGATE can also model sum-pool, and in the more general case, an injective message aggregation, whenever the sigmoid gate over all edges is some value $x \neq 0$ (Recall Eq. 9). This happens if $\mathbf{w}^\ell$ in Eq. 10 is a zero vector, the bias $b$ in Eq. 10 is non-zero and the submatrix in $\mathbf{W}_1^1$ of Eq. 4 applying linear transformation to Lipschitz embeddings of nodes is 0. □

---

[2]As in GIN [30], we assume countable features on nodes.