

# LSTOD: LATENT SPATIAL-TEMPORAL ORIGIN-DESTINATION PREDICTION MODEL AND ITS APPLICATIONS IN RIDE-SHARING PLATFORMS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Origin-Destination (OD) flow data is an important instrument in transportation studies. Precise prediction of customer demands from each original location to a destination given a series of previous snapshots helps ride-sharing platforms to better understand their market mechanism. However, most existing prediction methods ignore the network structure of OD flow data and fail to utilize the topological dependencies among related OD pairs. In this paper, we propose a latent spatial-temporal origin-destination (LSTOD) model, with a novel convolutional neural network (CNN) filter to learn the spatial features of OD pairs from a graph perspective and an attention structure to capture their long-term periodicity. Experiments on a real customer request dataset with available OD information from a ride-sharing platform demonstrate the advantage of LSTOD in achieving at least 6.5% improvement in prediction accuracy over the second best model.

## 1 INTRODUCTION

Spatial-temporal prediction of large-scale network-based OD flow data plays an important role in traffic flow control, urban routes planning, infrastructure construction, and the policy design of ride-sharing platforms, among others. On ride-sharing platforms, customers keep sending requests with origins and destinations at each moment. Knowing the exact original location and destination of each future trip allows platforms to prepare sufficient supplies in advance to optimize resource utilization and improve users' experience. Given the destinations of prospective demands, platforms can predict the number of drivers transferring from busy to idle status. Prediction of dynamic demand flow data helps ride-sharing platforms to design better order dispatch and fleet management policies for achieving the demand-supply equilibrium as well as decreased passenger waiting times and increased driver serving rates.

Many efforts have been devoted to developing traffic flow prediction models in the past few decades. Before the rise of deep learning, traditional statistical and machine learning approaches dominate this field (Li et al., 2012; Lippi et al., 2013; Moreira-Matias et al., 2013; Shekhar & Williams, 2007; Idé & Sugiyama, 2011; Zheng & Ni, 2013). Most of these models are linear and thus ignore some important non-linear correlations among the OD flows. Some other methods (Kwon & Murphy, 2000; Yang et al., 2013) further use additional manually extracted external features, but they fail to automatically extract the spatial representation of OD data. Moreover, they roughly combine the spatial and temporal features when fitting the prediction model instead of dynamically modelling their interactions.

The development of deep learning technologies brings a significant improvement of OD flow prediction by extracting non-linear latent structures that cannot be easily covered by feature engineering. (Xingjian et al., 2015; Ke et al., 2017; Zhou et al., 2018). Zhang et al. (2016; 2017) modeled the whole city as an entire image and employed residual neural network to capture temporal closeness. Ma et al. (2017) and Yu et al. (2017) also learned traffic as images but they used LSTM instead to obtain the temporal dependency. Yao et al. (2018b) proposed a Deep Multi-View Spatial-Temporal Network (DMVST-Net) framework to model both spatial and temporal relationships. However, using standard convolution filters suffers from the problem that some OD flows covered by a receptive field of regular CNNs are not spatially important.

The aim of this paper is to introduce a hierarchical Latent Spatial-Temporal Origin-Destination (LSTOD) prediction model to jointly extract the complex spatial-temporal features of OD data by using some well-designed CNN-based architectures. Instead of modelling the dynamic OD networks as a sequence of images and applying standard convolution filters to capture their spatial information, we introduce a novel Spatial Adjacent Convolution Network (SACN) that uses an irregular convolution filter to cover most related OD flows for a target one. The OD flows connected by common starting and/or ending vertexes, which may fall into different regions of the flow map, can be spatially correlated and topologically connected. Moreover, for most ride-sharing platforms, a passenger is more likely to send a new request from the location where his/her last trip ends in. To learn such sequential dependency, we introduce a temporal gated CNN (TGCNN) (Yu et al., 2018) and integrate it with SACN by using the sandwich-structured ST-conv block in order to collectively catch the evolutionary mechanism of dynamic OD flow systems. A periodically shifted attention mechanism is also used to capture the shift in the long-term periodicity. Finally, the combined short-term and long-term representations are fed into the final prediction layer to complete the architecture. Our contributions are summarized as follow:

- To the best of our knowledge, it is the first time that we propose purely convolutional structures to learn both short-term and long-term spatio-temporal features simultaneously from dynamic origin-destination flow data.
- We propose a novel SACN architecture to capture the graph-based semantic connections and functional similarities among correlated OD flows by modeling each OD flow map as an adjacency matrix.
- We design a periodically shift attention mechanism to model the long-term periodicity when using convolutional architecture TGCNN in learning temporal features.
- Experimental results on two real customer demand data sets obtained from a ride-sharing platform demonstrate that LSTOD outperforms many state-of-the-art methods in OD flow prediction, with 6.5% to 15.0% improvement of testing RMSE.

## 2 DEFINITIONS AND PROBLEM STATEMENT

For a given urban area, we observe a sequence of adjacency matrices representing the OD flow maps defined on a fixed vertex set  $V$ , which indicates the  $N$  selected sub-regions from this area. We let  $\mathbb{V} = \{v_1, v_2, \dots, v_N\}$  denote the vertex set with  $v_i$  being the  $i$ -th sub-region. The shape of each grid  $v_i$  could be either rectangles, hexagons or irregular sub-regions. We define the dynamic OD flow maps as  $\{\mathbf{O}_{d,t}\}$ , where  $d = 1, \dots, D$  and  $t = 1, \dots, T$  represent the day and time indexes, respectively. For each snapshot  $\mathbf{O}_{d,t} = (o_{d,t}^{ij}) \in R^{N \times N}$ , the edge weight  $o_{d,t}^{ij}$  at row  $i$  and column  $j$  denotes the flow amount from node  $v_i$  to node  $v_j$  at time  $t$  of day  $d$ . A larger edge weight  $o_{d,t}^{ij}$  is equivalent to a strong connection between nodes  $v_i$  and  $v_j$ . The  $\mathbf{O}_{d,t}$ s' are asymmetric since all the included OD flows are directed. Specifically,  $o_{d,t}^{ij} = 0$  if there is no demand from  $v_i$  to  $v_j$  within the  $t$ -th time interval of day  $d$ .

The goal of our prediction problem is to predict the snapshot  $\mathbf{O}_{d,t+j} \in R^{N \times N}$  in the future time window  $(t+j)$  of day  $d$  given previously observed data, including both short-term and long-term historical information. The short-term input data consists of the last  $p_1$  time-stamps from  $t+1-p_1$  to  $t$ , denoted by  $\mathbb{X}_1 = \{\mathbf{O}_{d,t+1-p_1}, \mathbf{O}_{d,t+1-p_1+1}, \dots, \mathbf{O}_{d,t}\}$ . The long-term input data is made up of  $q$  time series  $\{\mathbf{O}_{d-\varphi,t+j-(p_2-1)/2}, \dots, \mathbf{O}_{d-\varphi,t+j+(p_2-1)/2}\}$  of length  $p_2$  for each previous day  $(d-\varphi)$  with the predicted time index  $(t+j)$  in the middle for  $\varphi = 1, \dots, q$ . We let  $\mathbb{X}_2 = \{\mathbf{O}_{d-q,t+j-(p_2-1)/2}, \dots, \mathbf{O}_{d-q,t+j+(p_2-1)/2}, \dots, \mathbf{O}_{d-1,t+j-(p_2-1)/2}, \dots, \mathbf{O}_{d-1,t+j+(p_2-1)/2}\}$  denote the entire long-term data. Increasing  $p_1$  and  $p_2$  leads to higher prediction accuracy, but more training time.

We reformulate the set of short-term OD networks  $\mathbb{X}_1$  into a 4D tensor  $\mathbf{X}_1 \in R^{N \times N \times p_1 \times 1}$  and concatenate the long-term snapshots  $\mathbb{X}_2$  into a 5D tensor  $\mathbf{X}_2 = (\mathbf{X}_{2,d-1}, \dots, \mathbf{X}_{2,d-q}) \in R^{q \times N \times N \times p_2 \times 1}$ . Each  $\mathbf{X}_{2,d-\varphi} \in R^{N \times N \times p_2 \times 1}$  for day  $d-\varphi$  is a 4D tensor for  $\varphi = 1, \dots, q$ . Therefore, we can finally define our latent prediction problem as follows:

$$o_{d,t+j} = F(\mathbf{X}_1, \mathbf{X}_2), \quad (1)$$

where  $F(\cdot, \cdot)$  represents the LSTOD model, which captures the network structures of OD flow data as well as the temporal dependencies in multiple scales. A notation table is attached in the appendix.

### 3 LSTOD FRAMEWORK

In this section, we describe the details of our proposed LSTOD prediction model. See Figure 1 for the architecture of LSTOD. The four major novelties and functionalities of LSTOD model include

- an end-to-end framework LSTOD constructed by all kinds of CNN modules to process dynamic OD flow maps and build spatio-temporal prediction models;
- a novel multi-layer architecture SACN to extract the network patterns of OD flow maps by propagating through edge connections, which can not be covered by traditional CNNs;
- a special module ST-Conv block used to combine SACN and gated temporal convolutions to coherently learning the essential spatio-temporal representations;
- a periodically shifted attention mechanism which is well designed for the purely convolutional ST-Conv blocks to efficiently utilize the long-term information by measuring its similarities with short-term data.

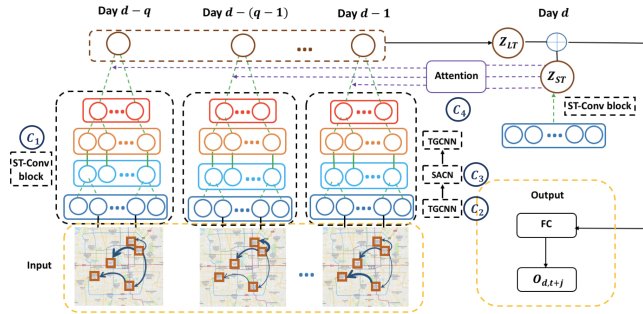


Figure 1: The architecture of LSTOD consisting of (C1) ST-Conv blocks, (C2) TGCNN, (C3) SACN, and (C4) attention.

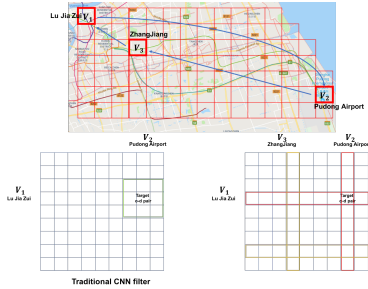


Figure 2: An example to show how standard CNN fails to capture the network structure of OD flows

#### 3.1 SPATIAL ADJACENT CONVOLUTION NETWORK

Before we formally introduce SACN, we first discuss why directly applying standard CNNs to the OD flow map  $O_{d,t}$  may disregard the connections between neighboring OD flows in the graph space. Figure 2 demonstrates that it fails to capture enough semantic information using the real-world example of ride demands. For the OD flow starting from  $v_1$  to  $v_2$ , as illustrated in the upper sub-figure, the most related OD flows should be those with either origin or destination being  $v_1$  or  $v_2$  in the past few timestamps. A certain part of the travel requests from  $v_1$  to  $v_2$  can be matched with some historical finished trips from a third-party location to  $V_1$  by the same group of people, for example a trip from  $v_3$  to  $v_1$ . However, as the lower-left sub-figure illustrates, some of the OD flows covered by a single CNN filter (the green square) such as the four corners of the kernel window may be topologically far away from the target one in the graph.

More generally, let's consider a target OD flow  $o_{d,t}^{ij}$  in a map  $O_{d,t}$ . Most of the components covered by a standard CNN of size  $q \times q$  with  $o_{d,t}^{ij}$  in the middle such as  $o_{d,t}^{(i+1)(j+1)}$  are less correlated than some out of the kernel window, for example  $o_{d,t}^{ki}$  when  $|k-j| > (q+1)/2$ . Moreover, if we change the order of the  $N$  vertexes in  $O_{d,t}$ , then the network structure is unchanged, but a different set of OD flows is covered by the receptive field of the same size with the center being  $o_{d,t}^{ij}$ .

We propose a novel CNN-based architecture SACN using a global-view receptive field to include all connected edges in the graph and exclude the topologically unrelated ones considered by standard CNN filters. As shown in the lower right sub-figure of Figure 2, OD flows with either origin or destination being  $v_i$  or  $v_j$ , covered by the red and yellow kernel windows, are considered to be the

most related ones for  $o_{d,t}^{ij}$  in row  $i$  and column  $j$ . Formally, we use SACN to extract the latent topological structure inside the OD flow network  $\mathbf{O}_{d,t}$  at each timestamp  $(d, t)$ . For an  $L$ -layer SACN architecture, the  $l$ -th layer takes  $M^{l-1}$  edge features obtained from the previous  $(l-1)$ -th layer as input and feeds the  $M^l$ -dimensional output to the next layer. The input of a general SACN layer  $l$  is a 3D tensor,  $\mathbf{A}_{d,t}(l) = \{A_{d,t}^{ijn}(l)\} \in R^{N \times N \times M^l}$ , which includes the  $M^l$  features of each of the  $N^2$  OD flows, and the output is another 3D tensor  $\mathbf{A}_{d,t}(l+1) = \{A_{d,t}^{ijn}(l+1)\}$  of size  $N \times N \times M^{l+1}$ . The learned representation of a target edge is defined as the weighted sum of those from the same row or column in the adjacency matrix, and those from the row or column in the transposed adjacency matrix. The  $n$ -th output of the  $l$ -th layer-wise SACN propagation for the OD flow from  $v_i$  to  $v_j$  of snapshot  $\mathbf{O}_{d,t}$ , denoted as  $A_{d,t}^{ijn}(l+1)$ , is written as

$$\mathcal{F}\left\{\sum_{m=1}^{M^l} \sum_{k=1}^N r_1^{kmn}(l) A_{d,t}^{ikm}(l) + c_1^{kmn}(l) A_{d,t}^{kjm}(l) + r_2^{kmn}(l) A_{d,t}^{kim}(l) + c_2^{kmn}(l) A_{d,t}^{jkm}(l)\right\}, \quad (2)$$

where  $n \in \{1, \dots, M^l\}$  and  $\{r_1^{kmn}(l)\}, \{r_2^{kmn}(l)\}, \{c_1^{kmn}(l)\},$  and  $\{c_2^{kmn}(l)\} \in R^{N \times M^l \times M^{l+1}}$  include all the related parameters to be learnt for the  $l$ -th SACN layer. The  $\mathcal{F}(\cdot)$  represents an elementwise activation function, such as  $\text{ReLU}(x) = \max(0, x)$ . The first part of (2) works by summing up the feature values of OD flows having either the same origin or destination with the target OD flow. The second part covers another set of OD flows that either start at  $v_j$  or end at  $v_i$ . Similar to standard CNN architectures, OD flows more related to the target one are more highly weighted by a multi-layer SACN operator. For notational simplification, we use  $A(\theta)_{*L}$  to represent a  $L$ -layer SACN operator and denote the overall output as  $\mathbf{S}_{d,t} \in R^{N \times N \times M^{L+1}}$ . Figure 4 in the appendix more clearly illustrates the architecture of SACN.

### 3.2 TEMPORAL GATED CNN

We use temporal gated CNN (TGCNN) (Yu et al., 2018) instead of RNN-based architectures such as LSTMs to capture the temporal representation, which makes our LSTOD a pure convolutional architecture. RNNs suffer from the problem of lower training efficiency, gradient instability, and time-consuming convergence. Moreover, the high dimension of the spatial representations captured by SACN and a potential long temporal sequence length make RNNs notoriously difficult to train. The CNNs is more flexible in handling various data structures and allows parallel computations to increase training speed.

TGCNN consists of two parts including one being a 3D convolution kernel applied to the spatial representations of all the  $N^2$  OD flows along the time axis and the other being a gated linear unit (GLU) as the gate mechanism. Given an input spatial-temporal sequential data  $\mathbf{Y} \in R^{N \times N \times r \times m_0}$ , the detailed architecture of a one-layer temporal gated CNN is formally defined as follows:

$$G(\gamma) *_{\tau} \mathbf{Y} = \mathbf{P} \odot \sigma(\mathbf{Q}) \in R^{N \times N \times (r-K+1) \times m_1}, \quad (3)$$

where  $\odot$  denotes the element-wise Hadamard product and  $\gamma$  denotes the set of parameters to be learnt. The output  $\mathbf{Q}$  with an element-wise sigmoid function  $\sigma(\cdot)$  work together as a gate to evaluate the importance of each element in  $\mathbf{P}$  and assign a weight before moving to the following layer.

### 3.3 ST-CONV BLOCKS

We use the spatial-temporal convolutional block (ST-conv block) to jointly capture the spatial-temporal features of OD flow data, which has a ‘sandwich’-structure architecture with an  $L$ -layer SACN operator in the middle connecting the two TGCNNs. The use of ST-Conv blocks have two major advantages. First, the block can be stacked or extended based on the dimension and characteristic of the spatio-temporal input data. Second, a temporal operation is applied before extracting the spatial information, which greatly reduces its computation complexity and memory consumption.

Both the input and output of a single ST-Conv block are 4D tensors. We denote the input as  $\mathbf{Y}^0 \in R^{N \times N \times r \times c_0}$  of  $r$  continuous snapshots and  $c_0$  features, which can be the short-term OD flow data by setting  $r = p_1$  and  $c_0 = 1$ . The mathematical definition of the ST-Conv block is given by.

$$\mathbf{Y}^1 = G_1(\gamma_1) *_{\tau} [A(\theta_0)_{*L} \{G_0(\gamma_0) *_{\tau} \mathbf{Y}^0\}], \quad (4)$$

where  $G_1(\cdot)$  and  $G_0(\cdot)$  are the two temporal gated CNN layers and  $A(\theta)*_L$  is an  $L$ -layer SACN operator. The  $(\theta_0, \gamma_0, \gamma_1)$  is the set of all parameters to be learnt. The  $m_1$  3D convolutional filters of kernel size  $1 \times 1 \times K_0$  and  $1 \times 1 \times K_1$  are used by the two TGCNN  $G_0(\gamma_0)*_\tau$  and  $G_1(\gamma_1)*_\tau$ , respectively. The  $L$ -layer SACN is applied to each of the  $\{r - (K_0 - 1)\}$  3D output of size  $N \times N \times m_1$  obtained from TGCNN  $G_0(\gamma_0)*_\tau$ , and then fed into the other TGCNN operator  $G_1(\gamma_1)*_\tau$ . One ST-Conv block shortens the temporal length of input  $\mathbf{Y}^0$  by  $(K_0 + K_1 - 2)$ , and the dimension of the output  $\mathbf{Y}^1$  becomes  $N \times N \times \{r - (K_0 + K_1 - 2)\} \times m_1$ . Accordingly, a set of  $n_{ST} = (r - 1)/(K_0 + K_1 - 2)$  ST-Conv blocks reduces the sequential length from  $r$  to 1. We can then flatten the spatial-temporal representation into a 3D tensor by squeezing out the temporal dimension.

The short-term spatial-temporal representation  $\mathbf{Z}_{ST} \in R^{N \times N \times c_{ST}}$  is obtained by continuously applying  $(p_1 - 1)/(K_{ST}^0 + K_{ST}^1 - 2)$  ST-Conv blocks to the short-term OD flow data  $\mathbf{X}_1 \in R^{N \times N \times p_1 \times 1}$ . The kernel sizes of the two TGCNNs in all ST-Conv blocks are fixed to be  $1 \times 1 \times K_{ST}^0$  and  $1 \times 1 \times K_{ST}^1$ , respectively. The  $c_{ST}$  filters are used by all the  $L$ -layer SACN and TGCNN layers. The detailed propagation of the  $n$ -th ST-Conv block is defined as

$$\mathbf{Z}_{ST}(n + 1) = G_1(\gamma_{ST}^1) *_\tau [A(\theta_{ST}) * _L \{G_0(\gamma_{ST}^0) *_\tau \mathbf{Z}_{ST}(n)\}], \quad (5)$$

Moreover,  $\mathbf{Z}_{ST}(1) = \mathbf{X}_1$ , and  $\mathbf{Z}_{ST} = \mathbf{Z}_{ST}(n_{ST} + 1) \in R^{N \times N \times c_{ST}}$  is the output of the last ST-Conv block and is denoted as the final learned short-term spatial-temporal representations.

### 3.4 PERIODICALLY SHIFTED ATTENTION MECHANISM

In addition to capturing the the spatial-temporal features from short-term OD flow data  $\tilde{\mathbf{X}}_1$ , we also take into account the long-term temporal periodicity due to the potential day-wise cycling patterns insides the OD flow data, decided by customer’s travelling schedule and the city’s traffic conditions. Directly applying ST-Conv blocks to an extremely long OD sequence which covers previous few days or weeks is computationally expensive. Only a small set of timestamps from each previous day is necessary to capture the long-term periodicity. As mentioned, we pick  $p_2$  time intervals for each day  $d - \varphi$  when predicting the time window  $(d, t + j)$  considering the non-strict long-term periodicity. This slight time shifting may be caused by unstable traffic peaks, holidays and extreme weather conditions among different days.

Inspired by the recently widely used attention mechanisms (Xu et al., 2015; Yao et al., 2018a; Liang et al., 2018) in spatial-temporal prediction problems, we propose a modified periodically shifted attention to work for the CNN-based ST-Conv blocks here. Different from Yao et al. (2018a) that measures the similarity between hidden units of LSTMs, the attention here is built on the intermediate outputs of TGCNNs where the concatenations are then fed into a new set of ST-Conv blocks. For each day  $(d - \varphi)$ , we apply  $(p_2 - n_{LT}^0)/(2K_{LT}^0 - 2)$  ST-Conv blocks to the day-level  $p_2$ -length sequential OD flow data indexed by  $\{\mathbf{O}_{d-\varphi, t+j-(p_2+1)/2}; \dots; \mathbf{O}_{d-\varphi, t+j+(p_2+1)/2}\}$  to reduce the sequence length from  $p_2$  to  $n_{LT}^0$ . We let the two TGCNNs in all the  $(p_2 - n_{LT}^0)/(2K_{LT}^0 - 2)$  ST-Conv blocks have the same filter size  $1 \times 1 \times K_{LT}^0$ . The propagation rule of the  $n$ -th ST-Conv blocks is defined as:

$$\mathbf{Z}_{d-\varphi}^{n+1} = G_1(\gamma_{LT}^{01}) *_\tau [A(\theta_{LT}^0) * _L \{G_0(\gamma_{LT}^{00}) *_\tau \mathbf{Z}_{d-\varphi}^n\}] \quad (6)$$

with  $\mathbf{Z}_{d-\varphi}^n \in R^{N \times N \times \{p_2 - 2(n-1)(K_{LT}^0 - 1)\} \times \tilde{c}}$  and  $\mathbf{Z}_{d-\varphi}^{n+1} \in R^{N \times N \times \{p_2 - 2n(K_{LT}^0 - 1)\} \times \tilde{c}}$  being the input and output, respectively. Specifically,  $\mathbf{Z}_{d-\varphi}^1 = \mathbf{X}_{2, d-\varphi}$  is the original OD flow data at day  $d - \varphi$ . All SACN and TGCNN layers use  $\tilde{c}$  convolutional filters and  $(\theta_{LT}^0, \gamma_{LT}^{00}, \gamma_{LT}^{01})$  is the parameter set.

We denote the day-level features of day  $(d - \varphi)$  captured by  $(p_2 - n_{LT}^0)/(2K_{LT}^0 - 2)$  ST-Conv blocks as  $\tilde{\mathbf{Z}}_{d-\varphi} \in R^{N \times N \times n_{LT}^0 \times \tilde{c}}$ , where  $\tilde{z}_{d-\varphi, \phi}^{ij} \in R^{\tilde{c} \times 1}$  denotes the  $\phi$ -th element along the time axis for the OD flow from  $v_i$  to  $v_j$ . We let  $z_{ST}^{ij} \in R^{c_{ST} \times 1}$  be the learned short-term representation at the OD flow from  $v_i$  to  $v_j$ . Then, a day-level output  $z_{d-\varphi}^{ij}$  can be obtained by summing up all the  $n_{LT}^0 \tilde{z}_{d-\varphi, \phi}^{ij}$ ’s by the weights which measure their similarities with  $z_{ST}^{ij}$ :

$$z_{d-\varphi}^{ij} = \sum_{\phi=1}^{n_{LT}^0} \beta_{d-\varphi, \phi}^{ij} \tilde{z}_{d-\varphi, \phi}^{ij}, \quad (7)$$

where  $\beta_{d-\varphi,\phi}^{ij}$  is the weight function of quantifying the similarity between  $\tilde{z}_{d-\varphi,\phi}^{ij}$  and  $z_{ST}^{ij}$  based on a score function  $\text{score}(\tilde{z}_{d-\varphi,\phi}^{ij}, z_{ST}^{ij})$ , which is defined as:

$$\beta_{d-\varphi,\phi}^{ij} = \frac{\exp(\text{score}(\tilde{z}_{d-\varphi,\phi}^{ij}, z_{ST}^{ij}))}{\sum_{\phi'} \exp(\text{score}(\tilde{z}_{d-\varphi,\phi'}^{ij}, z_{ST}^{ij}))}. \quad (8)$$

Moreover,  $\text{score}(\tilde{z}_{d-\varphi,\phi}^{ij}, z_{ST}^{ij})$  is defined as

$$v_\phi^T \tanh(\mathbf{W}_1 \tilde{z}_{d-\varphi,\phi}^{ij} + \mathbf{W}_2 z_{ST}^{ij} + b_s), \quad (9)$$

where  $\mathbf{W}_1 \in R^{\tilde{c} \times \tilde{c}}$ ,  $\mathbf{W}_2 \in R^{\tilde{c} \times c_{ST}}$ , and  $v_\phi \in R^{\tilde{c} \times 1}$  are learned projection matrices, and  $b_s$  is the added bias term. We let  $\mathbf{Z}_{d-\varphi} = (z_{d-\varphi}^{ij}) \in R^{N \times N \times \tilde{c}}$  denote the day-level output and then concatenate the  $q$   $\mathbf{Z}_{d-\varphi}$ 's along a new additional axis in the third dimension as

$$\mathbf{Z}_{LT}^0 = \text{Concat}_{\varphi=q}^1 \mathbf{Z}_{d-\varphi} \quad (10)$$

to build a new day-wise time series  $\mathbf{Z}_{LT}^0 \in R^{N \times N \times q \times \tilde{c}}$  of length  $q$ . Finally, we apply another set of  $(q-1)/(2K_{LT}^1-2)$  ST-Conv blocks to the day-wise sequence data generated by (10) to capture the long-term spatial-temporal representations, which is denoted by  $\mathbf{Z}_{LT} \in R^{N \times N \times c_{LT}}$ , where  $c_{LT}$  is the number of feature channels.

We concatenate the short-term and long-term spatial-temporal representations  $\mathbf{Z}_{ST}$  and  $\mathbf{Z}_{LT}$  together along the feature axis as  $\mathbf{Z} = \mathbf{Z}_{ST} \oplus \mathbf{Z}_{LT} \in R^{N \times N \times \mathcal{C}}$ , where  $\mathcal{C} = c_{ST} + c_{LT}$ . Then,  $\mathbf{Z}$  is modified to a 2D tensor  $\mathbf{Z} \in R^{N^2 \times \mathcal{C}}$  by flattening the first two dimensions while keeping the third one. We apply a fully connected layer to the  $\mathcal{C}$  feature channels together with an element-wise non-linear sigmoid function to get the final predictions for all the  $N^2$  OD flows.

### 3.5 DATA PROCESSING AND TRAINING

We normalize the original OD flow data in the training set to  $(0, 1)$  by Max-Min normalization and use 'sigmoid' activation for the final prediction layer to ensure that all predictions fall into  $(0, 1)$ . The upper and lower bounds are saved and used to denormalize the predictions of testing data to get the actual flow volumes.

We use  $L_2$  loss to build the objective loss during the training. The model is optimized via Back-propagation Through Time (BPTT) and Adam (Kingma & Ba, 2014). The whole architecture of our model is realized using Tensorflow (Abadi et al., 2016) and Keras (Chollet et al., 2015). All experiments were run on a cluster with one NVIDIA 12G-memory Titan GPU.

## 4 EXPERIMENT

In this section, we compare the proposed LSTOD model with some state-of-the-art approaches for latent traffic flow predictions. All compared methods are classified into traditional statistical methods and deep-learning based approaches. We use the demand flow data collected by a ride-sharing platform to examine the finite sample performance of OD flow predictions for each method.

### 4.1 DATASET DESCRIPTION

We employ a large-scale demand dataset obtained from a large-scale ride-sharing platform to do all the experiments. The dataset contains all customer requests received by the platform from 04/01/2018 to 06/30/2018 in two big cities A and B. Within each urban area,  $N = 50$  hexagonal regions with the largest customer demands (cover more than 80% of total demands) are selected with radius being 2 km,  $N = 50$  of which with the largest customer demands are selected to build the vertex set  $V$ . In total, 2500 OD flows are generated based on the 50 sub-regions.

We split the whole dataset into two parts. The data from 04/01/2018 to 06/16/2018 is used for model training, while the other part from 06/17/2018 to 06/30/2018 (14 days) serves as the testing set. The first two and half months of OD flow data is further divided in half to the training and

validation sets. The size ratio between the two sets is around 4:1. We let 30 min be the length of each timestamp and the value of the OD flow from  $v_i$  to  $v_j$  is the cumulative number of customer requests. We make predictions for all the  $50^2$  OD flows in the incoming 1st, 2nd, and 3rd 30 minutes (i.e.  $t + 1, t + 2, t + 3$ ) by each compared method, given the historical data with varied  $(p_1, p_2)$  combinations. For those model settings incorporating long-term information, we trace back  $q = 3$  days to capture the time periodicity. We use Rooted Mean Square Error to evaluate the performance of each method:

$$\text{RMSE} = \sqrt{\frac{1}{N^2 * |\mathcal{T}_0|} \sum_{i=1}^N \sum_{j=1}^N \sum_{(d,t) \in \mathcal{T}_0} (o_{d,t}^{ij} - \hat{o}_{d,t}^{ij})^2}, \quad (11)$$

$o_{d,t}^{ij}$  and  $\hat{o}_{d,t}^{ij}$  are the true value and prediction at the OD flow from vertex  $v_i$  to vertex  $v_j$  at time  $(d, t)$ , respectively.  $\mathcal{T}_0$  is the set containing all the predicted time points in the testing data.

## 4.2 COMPARED METHODS

All state-of-the-art methods to be compared are listed as follows, some of which are modified to work for the OD flow data: (i) **Historical average (HA)**: HA predicts the demand amount at each OD flow by the average value of previous 5 days, (ii) **Autoregressive integrated moving average (ARIMA)**, (iii) **Support Vector Machine Regression (SVMR)**, (iv) **Latent Space Model for Road Networks (LSM-RN)** (Deng et al., 2016), (v) **Dense + BiLSTM (DLSTM)** (Alché & de La Fortelle, 2017) and (vi) **Spatiotemporal Recurrent Convolutional Networks (SRCN)** (Yu et al., 2017). We only consider latent models in this paper, that is, no external covariates are allowed, while only the historical OD flow data is used to extract the hidden spatial-temporal features.

## 4.3 PREPROCESSING AND PARAMETERS

We tune the hyperparameters of each compared model to obtain the optimal prediction performance. Specifically, we get  $(p^*, d^*, q^*) = (3, 0, 3)$  for ARIMA and  $k^* = 15, \gamma^* = 2^{-5}, \lambda^* = 10$  for LSM-RN. The optimal kernel size of the spatial based CNN kernel is  $11 \times 11$  in SRCN model.

For fair comparison, we set the length of short-term OD flow sequence to be  $p_1 = 9$  (i.e., previous 4.5 hours),  $q = 3$  for long-term data which covers the three most recent days, and the length of each day-level time series  $p_2 = 5$  to capture the periodicity shifting (one hour before and after the predicted time index). More analysis of how variant  $(p_1, p_2)$  combinations may affect the prediction performance of LSTOD will be studied latter.

A two-layer architecture is used by all the deep-learning based methods to extract the spatial patterns inside the OD flow data ( $L = 2$  for both short-term and long-term SACN). We set the filter size of all deep learning layers in both spatial and temporal space to be 64, including the SACNs and TGCNNs in our LSTOD model with  $\tilde{c} = c_{ST} = c_{LT} = 64$ .

## 4.4 RESULTS

**Comparison with state-of-the-art methods.** Table 1 summarizes the finite sample performance for all the competitive methods and our LSTOD model in terms of the prediction RMSE on the testing data. For city A, LSTOD outperforms all other methods on the testing data with the lowest RMSE (2.44/2.59/2.69), achieving (6.51%/6.83%/7.24%) improvement over the second best method 'SRCN'. This demonstrates the advantages of using our spatial-temporal architecture and long-term periodicity mechanism in modelling the dynamic evolution of OD flow networks. The improvement increases as the predicting scope increases since our model captures the long-term periodicity. The 'Dense + BiLSTM' outperforms traditional approaches by more precisely learning the temporal dependency using deep learning architecture, but it fails to model the underlying graph structure of OD flow data. Both 'ARIMA' and 'LSM-RN' perform poorly, even much worse than HA, indicating that they cannot capture enough short-term spatial-temporal features to get the evolution trend of OD flow data. The LSTOD performs even better on city B compared to SRCN (11.34% - 15% improvement) since the long-term periodical pattern in city B is more significant compared with that in city A.

Table 1: Comparison with State-of-the-art methods

Method	City A			City B		
	30 min	60 min	90 min	30 min	60 min	90 min
HA		4.64			3.53	
ARIMA	5.64	6.01	6.49	8.82	9.47	10.82
LSVR	3.53	3.95	4.06	4.13	4.85	5.12
LSM-RN	5.73	6.36	6.74	6.02	6.75	7.36
DLSTM	3.08	3.59	3.99	3.88	4.09	4.59
SRCN	2.61	2.78	2.90	2.91	3.08	3.20
<b>LSTOD</b>	<b>2.44</b>	<b>2.59</b>	<b>2.69</b>	<b>2.58</b>	<b>2.67</b>	<b>2.72</b>

**ACN VS standard local CNN.** In this experiment, we will show that our proposed SACN outperforms standard CNNs in capturing the hidden network structure of the OD flow data. Given the model setting that  $N = 50$  sub-regions of city A are used to build the dynamic OD flow matrices, the number of pixels being covered by SACN at each single snapshot is  $50 \times 4 = 200$ . For fair comparison, the largest receptive field of standard CNN should be no bigger than a  $15 \times 15$  window, which includes 225 elements each time. We consider five different kernel sizes including  $5 \times 5$ ,  $8 \times 8$ ,  $11 \times 11$ ,  $14 \times 14$ , and  $15 \times 15$ . We replace SCAN in our model by standard CNN in order to fairly compare its performance. All hyper-parameters are fixed but only the kernel size of CNNs being changed. Moreover, we only consider the baseline short-term mode of LSTOD model while ignoring the long-term information. As Figure 3 illustrates, standard CNN achieves the best performance with the smallest RMSE = 2.64 on testing data with the filter size being  $11 \times 11$ , which is still higher than that using SACN with RMSE = 2.54. Specifically, RMSE increases when the receptive field is larger than  $11 \times 11$  since the spatial correlations among the most related OD flows (sharing common origin or destination nodes) are smoothed with the increase in the filter size  $((8 \times 2 - 1)/64 > (14 \times 2 - 1)/196)$ . This experiment shows that treating the dynamic demand matrix as an image, and applying standard CNN filters does not capture enough spatial correlations among related OD flows without considering their topological connections from the perspective of graphs. For more details, please refer to Figure 3.

**Comparison with variants of LSTOD.** Table 2 shows the finite sample performance of our proposed model LSTOD and its different variants based on the demand data from city A. We can see that the LSTOD model without using attention mechanism (RMSE = 2.49) outperforms the baseline setting only using short-term data (RMSE = 2.54). This shows the necessity of modeling the seasonal temporal patterns. Moreover, the complete model using the attention mechanism (RMSE = 2.44) outperforms the one without using it (RMSE = 2.49). It indicates that the periodically shifted attention can capture the shifting of the day-wise periodicity and extract more seasonal patterns to improve prediction accuracy. We do some more experiments to show how different hyperparameter configurations influence the model performance. For more details, please refer to the appendix.

Method	RMSE		
	30 min	60 min	90 min
ACN + GCNN	2.54	2.71	2.83
LSTOD (no attention)	2.49	2.63	2.72
<b>LSTOD (complete)</b>	<b>2.44</b>	<b>2.59</b>	<b>2.69</b>

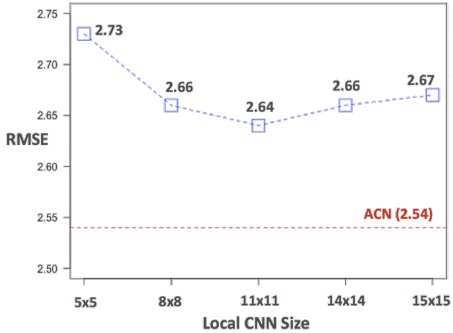


Figure 3: RMSE on testing data with respect to ACN and standard CNN using different kernel sizes.

Table 2: Evaluation of LSTOD and its variants



## REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation*, 265-283, 2016.
- Florent Alché and Arnaud de La Fortelle. An lstm network for highway trajectory prediction. In *20th International Conference on Intelligent Transportation Systems (ITSC)*, 353-359. IEEE, 2017.
- François Chollet et al. Keras, 2015.
- Dingxiong Deng, Cyrus Shahabi, Ugur Demiryurek, Linhong Zhu, Rose Yu, and Yan Liu. Latent space model for road networks to predict time-varying traffic. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016.
- Tsuyoshi Idé and Masashi Sugiyama. Trajectory regression on road networks. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- Jintao Ke, Hongyu Zheng, Hai Yang, and Xiqun Michael Chen. Short-term forecasting of passenger demand under on-demand ride services: A spatio-temporal deep learning approach. *Transportation Research Part C: Emerging Technologies*, 85:591–608, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jaimyoung Kwon and Kevin Murphy. Modeling freeway traffic with coupled hmms. Technical report, Technical report, Univ. California, Berkeley, 2000.
- Xiaolong Li, Gang Pan, Zhaohui Wu, Guande Qi, Shijian Li, Daqing Zhang, Wangsheng Zhang, and Zonghui Wang. Prediction of urban human mobility using large-scale taxi traces and its applications. *Frontiers of Computer Science*, 6(1):111–121, 2012.
- Yuxuan Liang, Songyu Ke, Junbo Zhang, Xiuwen Yi, and Yu Zheng. Geoman: Multi-level attention networks for geo-sensory time series prediction. In *IJCAI*, 3428-3434, 2018.
- Marco Lippi, Matteo Bertini, and Paolo Frasconi. Short-term traffic flow forecasting: An experimental comparison of time-series analysis and supervised learning. *IEEE Transactions on Intelligent Transportation Systems*, 14(2):871–882, 2013.
- Xiaolei Ma, Zhuang Dai, Zhengbing He, Jihui Ma, Yong Wang, and Yunpeng Wang. Learning traffic as images: a deep convolutional neural network for large-scale transportation network speed prediction. *Sensors*, 17(4):818, 2017.
- Luis Moreira-Matias, Joao Gama, Michel Ferreira, Joao Mendes-Moreira, and Luis Damas. Predicting taxi-passenger demand using streaming data. *IEEE Transactions on Intelligent Transportation Systems*, 14(3):1393–1402, 2013.
- Shashank Shekhar and Billy M Williams. Adaptive seasonal time series models for forecasting short-term traffic flow. *Transportation Research Record*, 2024(1):116–125, 2007.
- SHI Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, 802-810, 2015.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, 2048-2057, 2015.
- Bin Yang, Chenjuan Guo, and Christian S Jensen. Travel cost inference from sparse, spatio-temporally correlated time series using markov models. *Proceedings of the VLDB Endowment*, 6(9): 769–780, 2013.

- Huaxiu Yao, Xianfeng Tang, Hua Wei, Guanjie Zheng, Yanwei Yu, and Zhenhui Li. Modeling spatial-temporal dynamics for traffic prediction. *arXiv preprint arXiv:1803.01254*, 2018a.
- Huaxiu Yao, Fei Wu, Jintao Ke, Xianfeng Tang, Yitian Jia, Siyu Lu, Pinghua Gong, and Jieping Ye. Deep multi-view spatial-temporal network for taxi demand prediction. *arXiv preprint arXiv:1802.08714*, 2018b.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 3634-3640. AAAI Press, 2018.
- Haiyang Yu, Zhihai Wu, Shuqin Wang, Yunpeng Wang, and Xiaolei Ma. Spatiotemporal recurrent convolutional networks for traffic prediction in transportation networks. *Sensors*, 17(7):1501, 2017.
- Junbo Zhang, Yu Zheng, Dekang Qi, Ruiyuan Li, and Xiuwen Yi. Dnn-based prediction model for spatio-temporal data. In *Proceedings of the 24th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 92. ACM, 2016.
- Junbo Zhang, Yu Zheng, and Dekang Qi. Deep spatio-temporal residual networks for citywide crowd flows prediction. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Jiangchuan Zheng and Lionel M Ni. Time-dependent trajectory regression on road networks via multi-task learning. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.
- Xian Zhou, Yanyan Shen, Yanmin Zhu, and Linpeng Huang. Predicting multi-step citywide passenger demands using attention-based neural networks. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, 736-744. ACM, 2018.

## A NOTATION TABLE

$\mathbf{O}_{d,t}$	OD flow maps of day $d$ , time $t$
$o_{d,t}^{ij}$	Flow amount from vertex $v_i$ to $v_j$ at day $d$ , time $t$
$\mathbf{X}_1, \mathbf{X}_2$	Short-term and long-term spatial-temporal OD flow data
$\mathbf{X}_{2,d-\varphi}$	long-term spatial-temporal OD flow data of day $d - \varphi$
$\mathbf{A}_{d,t}(l+1)$	Output of $l$ -th SACN layer when input being $\mathbf{O}_{d,t}$
$A(\theta)*_L$	$L$ -layer SACN operator
$G(\gamma)*_\tau$	One-layer TGCNN operator
$\mathbf{Z}_{ST}(n+1)$	Short-term output of the $l$ -th ST-Conv block
$\mathbf{Z}_{ST}$	Final captured short-term spatial-temporal representations
$\mathbf{Z}_{d-\varphi}^{n+1}$	Output of $n$ -th ST-Conv block of day $d - \varphi$
$\tilde{\mathbf{Z}}_{d-\varphi}$	features of day $(d - \varphi)$ captured by ST-Conv blocks
$\tilde{z}_{d-\varphi,\phi}^{ij}$	$\phi$ -th element along the time axis of $\tilde{\mathbf{Z}}_{d-\varphi}$ from $v_i$ to $v_j$
$\beta_{d-\varphi,\phi}^{ij}$	similarity between $\tilde{z}_{d-\varphi,\phi}^{ij}$ and $z_{ST}^{ij}$
$\mathbf{Z}_{d-\varphi}$	day-level output of day $d - \varphi$
$\mathbf{Z}_{LT}$	long-term spatial-temporal representations
$\mathbf{Z}$	$\mathbf{Z}_{ST} \oplus \mathbf{Z}_{LT}$ , combined spatial-temporal representations

## B ILLUSTRATION OF SACN

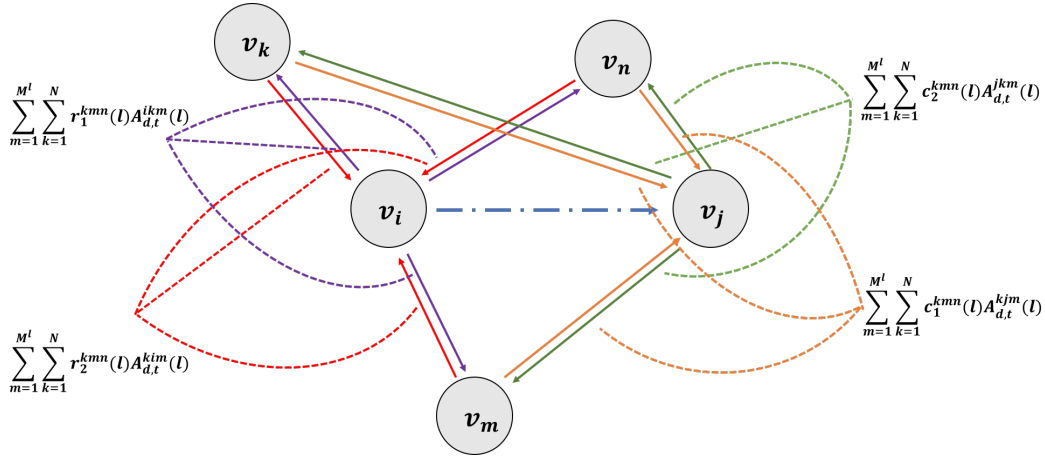


Figure 4: Working mechanism of spatial adjacent convolution network (SACN) for a target OD flow from  $v_i$  to  $v_j$

## C TRAINING DETAILS

Batch normalization is used in the SACN component. The batch size in our experiment was set to 10, corresponding to 10 randomly sampled timestamps and all the  $50^2$  OD flows in each snapshot. The initial learning rate is set to be  $10^{-4}$  with a decay rate  $10^{-6}$ . We use early stopping for all the deep learning-based methods where the training process is terminated when the RMSE over validation set has not been improved for 10 successive epochs. The maximal number of epochs allowed is 100.

## D COMPARISON OF DIFFERENT HYPERPARAMETER CONFIGURATIONS

In this section, we want to explore how some important hyperparameters of input OD flow data, for example  $p_1$  and  $p_2$ , may affect the performance of our LSTOD model.

Figure 5 (b) compares RMSE on testing data by STOD model with different data settings. Varied combinations of the short-term sequence length  $p_1$  and the long-term day-level sequence length  $p_2$  are studied. We can see that the best performance is achieved as  $(p_1, p_2) = (7, 5)$  with RMSE = 2.41. Specifically, settings with different  $p_1$ 's under  $p_2 = 5$  consistently outperform those under  $p_2 = 7$ . It may demonstrate that the shift can usually be captured within a short time range, while a longer time sequence may smooth the significance. Table 3 provides the detailed prediction results for each data setting.

Table 3: Comparison of STOD under different  $p_1, p_2$  combinations

$p_2$	$(K_{LT}^0, K_{LT}^1)$	$p_1$	$(K_{ST}^0, K_{ST}^1)$	RMSE
5	(2, 2)	5	(2, 2)	2.45
		<b>7</b>	<b>(2, 3)</b>	<b>2.41</b>
		9	(3, 3)	2.42
		11	(3, 4)	2.43
		13	(4, 4)	2.43
7	(3, 2)	5	(2, 2)	2.45
		7	(2, 3)	2.44
		9	(3, 3)	2.44
		11	(3, 4)	2.44
		13	(4, 4)	2.49

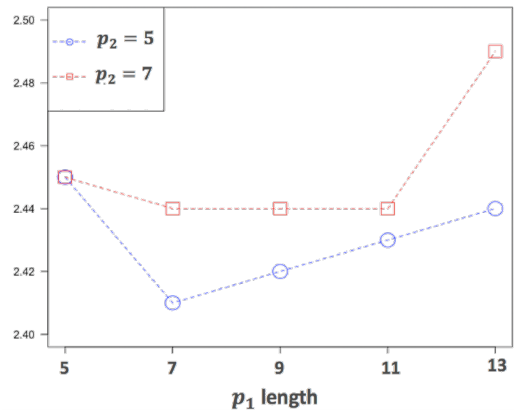


Figure 5: RMSE on testing data with respect to STOD with different  $p_1$  and  $p_2$  combinations.