

# Transportation Research Record

## A Pseudo-3D Convolutional Neural Network based Framework for Short-term Mixed Passenger Flow Prediction in Large-scale Public Transit

--Manuscript Draft--

<b>Full Title:</b>	A Pseudo-3D Convolutional Neural Network based Framework for Short-term Mixed Passenger Flow Prediction in Large-scale Public Transit
<b>Abstract:</b>	<p>Globally, many metropolitan areas tend to be promoting the multi-modal public transit which interconnects different modes to move large quantities of commuters in the urban area. Therefore, the accurate short-term passenger flow prediction can contribute to improve the reliability, responsiveness and the service quality of the multi-modal public transit. In this paper, we propose an end-to-end deep learning framework, called ST-Pseudo3D Net to collectively predict different types of public transit passenger flow at each region of city in the near future. The framework employs the deep Pseudo-3D residual architecture to model the network-wide spatial-temporal correlations among different types of passenger flow. Based on smart card data collected from Singapore's multi-modal public transit, considering metro passenger flow, bus passenger flow as well as the transfer passenger flow between these metro and bus, totally we construct 7 different types of passenger flow to be collectively predicted, where the experiments results demonstrate that the proposed model synthetically outperforms other baselines. To the best knowledge of the authors, this is the first attempt to investigate the integrated prediction for multi-modal passenger flow leveraging deep learning techniques.</p>
<b>Manuscript Classifications:</b>	Data and Information Technology; Artificial Intelligence and Advanced Computing Applications ABJ70; Artificial Intelligence; Machine Learning (Artificial Intelligence); Neural Networks; Transformative Transit Data AP000; Smartcard; Transportation Issues in Major Cities and Urban Mobility ABE30; Urban Mobility; Public Transportation; Urban Transportation Data and Information Systems ABJ30; Multimodal Analysis; General
<b>Manuscript Number:</b>	
<b>Article Type:</b>	Presentation
<b>Order of Authors:</b>	Siyu Hao
	Dingyi Zhuang
	De Zhao
	Der-Horng Lee

1 **A Pseudo-3D Convolutional Neural Network based**  
2 **Framework for Short-term Mixed Passenger Flow**  
3 **Prediction in Large-scale Public Transit**

4 Siyu Hao  
5 (Corresponding author)  
6 Department of Civil & Environmental Engineering,  
7 National University of Singapore  
8 Singapore 117576, Singapore  
9 Email: siyuhao@u.nus.edu

10 Dingyi Zhuang  
11 Department of Civil Engineering and Applied Mechanics  
12 McGill University  
13 Montreal, Quebec H3A 0C3, Canada  
14 Email: dingyi.zhuang@mail.mcgill.ca

15 De Zhao  
16 Department of Civil & Environmental Engineering  
17 National University of Singapore  
18 Singapore 117576, Singapore  
19 Email: ceezde@nus.edu.sg

20 Der-Horng Lee  
21 Department of Civil & Environmental Engineering,  
22 National University of Singapore  
23 Singapore 117576, Singapore  
24 Email: dhl@nus.edu.sg

25 Word count: 4516 words + 4 table × 250 + 5 figures = 5516 words

26 Submission Date: 1st Aug 2019

**1 ABSTRACT**

2 Globally, many metropolitan areas tend to be promoting the multi-modal public transit which in-  
3 terconnects different modes to move large quantities of commuters in the urban area. Therefore,  
4 the accurate short-term passenger flow prediction can contribute to improve the reliability, re-  
5 sponsiveness and the service quality of the multi-modal public transit. In this paper, we propose  
6 an end-to-end deep learning framework, called *ST-Pseudo3D Net* to collectively predict different  
7 types of public transit passenger flow at each region of city in the near future. The framework  
8 employs the deep Pseudo-3D residual architecture to model the network-wide spatial-temporal  
9 correlations among different types of passenger flow. Based on smart card data collected from  
10 Singapore's multi-modal public transit, considering metro passenger flow, bus passenger flow as  
11 well as the transfer passenger flow between these metro and bus, totally we construct 7 different  
12 types of passenger flow to be collectively predicted, where the experiments results demonstrate  
13 that the proposed model synthetically outperforms other baselines. To the best knowledge of the  
14 authors, this is the first attempt to investigate the integrated prediction for multi-modal passenger  
15 flow leveraging deep learning techniques.

16

17 **Keywords:** public transit, deep learning, 3DCNN, passenger flow, multi-modal

## 1 INTRODUCTION

2 The core concept of a multi-modal public transit system is to provide passengers with reliable, con-  
3 venient, highly connective and integrated services using different effective transportation modes.  
4 Normally, an integrated multi-modal public transport system mainly consists of an urban metro  
5 system and a bus system, where these two subsystems interconnect collaboratively to guarantee  
6 that the urban is functioning smoothly and properly. However, many public transport systems in  
7 global metropolitan areas are still facing a lot of challenges, especially when the urban dynamics  
8 are becoming more and more active. The accurate short-term passenger flow prediction is one of  
9 the critical solutions that can provide significant improvements at the operation and management  
10 level. Specifically, the operators might pointedly adjust the real-time operation scheme and allocate  
11 the resources on demand according to the prediction for the passenger demand in the near future.  
12 Additionally, if any anomaly pattern was detected from the prediction, then relevant agencies could  
13 timely implement the responses, which contributes is of great significance to public safety.

14 In fact, there are many related research that are focusing on predicting short-term passen-  
15 ger/traffic flow in the past years, where the methodologies that have been utilized range from some  
16 traditional statistical method to a batch of advanced deep learning models. Especially, leverag-  
17 ing deep learning techniques to solve transportation related problems has become the recent main  
18 stream, which has also demonstrated that the deep learning based techniques do have a number of  
19 superiorities in modelling transportation data.

20 Among these related research, the prediction for traffic flow or traffic speed on road using  
21 vehicle trajectory data or traffic flow data accounts for a relatively larger proportion. For instance,  
22 Ma et al. (1) proposed a deep learning framework to predict the traffic congestion using GPS  
23 trajectory data from taxi. The framework is formed by combining Restricted Boltzmann Machine  
24 (RBM) and Recurrent Neural Network (RNN), which shows effectiveness in the predicting traffic  
25 congestion in large-scale network. Lv et al. (2) firstly utilized a stacked autoencoder (SAE) based  
26 model to forecast the short-term traffic flow, through which the complex spatial-temporal features  
27 of traffic data could be well captured and abstracted. Polson and Sokolov (3) developed a hybrid  
28 deep learning architecture for short-term traffic prediction using road sensor data. This framework  
29 contains a liner vector autoregressive model for detecting the spatial-temporal correlations among  
30 predictors, of which the findings will be fed into another neural network with a deep structure to  
31 further model the nonlinear dependencies. Cui et al. (4), Zhao et al. (5), Fu et al. (6), Yu et al.  
32 (7) and Wang et al. (8) investigated the use of Long Short Term Memory (LSTM) based models  
33 to predict short-term traffic flow/speed, as LSTM based models can effectively learn relatively  
34 long-range temporal dependencies.

35 Compared with the short-term prediction for link/network traffic based on trajectory data,  
36 leveraging deep learning techniques to forecast passenger flow/demand in public transit system  
37 has received less attention from researchers. For bus system, Liu and Chen (9) combined SAE and  
38 Deep Neural Network (DNN) forming a hierarchical hybrid framework to predict short-term bus  
39 passenger flow. Bai et al. (10) proposed a multi-pattern deep fusion model to predict short-term  
40 bus passenger flow, where the Affinity Propagation (AP) algorithm is applied to identify different  
41 patterns from passenger flow and then the Deep Belief Network (DBN) is fused in each pattern  
42 to obtain the abstract representation. When it comes to metro system, Ma et al. (11) developed a  
43 hybrid architecture to predict large-scale metro ridership by taking advantages of both CNN and  
44 Bi-directional LSTM. Liu et al. (12) proposed a deep learning based framework that incorporates  
45 both LSTM module and manually designed features which is able to model the spatial-temporal

1 characteristics of metro passenger flow. Furthermore, Zhang et al. (13) designed a deep residual  
2 CNN based framework to predict crowd flows (inflow and outflow) at each region in the city area,  
3 which has been validated on taxi trajectory data from Beijing and bike sharing data from New York  
4 City.

5 As reviewed and discussed above, although the existing research on the related topics have  
6 already gain some achievements in both theory and practice, some limitations and unexplored  
7 gaps are still observed. According to our investigation, for the short-term passenger/traffic flow  
8 prediction related tasks, all the proposed frameworks can only make prediction for a single type of  
9 flow (e.g. metro passenger, bus passenger, taxi), without considering the mixed passenger flows in  
10 the multi-modal transport system. The motivations of collectively predicting the mixed passenger  
11 flow in a large-scale multi-modal public transit are summarized as follows:

- 12 • Many multi-modal public transit systems around the world bring together both buses and  
13 metro to move people from place to place in the city area. Therefore, the investigating  
14 merely on a single mode (e.g. bus or metro) could not comprehensively characterize the  
15 overall pattern of the multi-modal public transit.
- 16 • Normally, the metro system in the urban area forms the backbone network, which is  
17 mainly served for mass and long distance transit while the bus system can further com-  
18 plement the metro system and cover the area devoid of metro service as a local auxiliary  
19 feeder solution. Thus, incorporating the bus passenger prediction into metro passenger  
20 prediction would not only broaden the scope of the network to larger-scale, but also in-  
21 crease the density of the network that allow the prediction to be made in a finer resolution.
- 22 • The highly integration of different modes of transportation in a multi-modal public transit  
23 is imperative for attracting passengers and providing them with seamless services. Many  
24 commuters use more than one mode of transportation to complete the journey, where  
25 the whole journey consists of several different boarding, alighting and transfer activities.  
26 Therefore, if different types of passenger flow (e.g. bus boarding demand, metro alight-  
27 ing demand, etc.) at each region in the city could be collectively predicted by a single  
28 integrated framework, it would provide comprehensive decision supports for the relevant  
29 agencies and operators.
- 30 • Many cities are promoting the use of public transport to embrace the sustainability and  
31 safety. We have witnessed that the share of usage of public transport have been increasing  
32 over the past years in many cities. In some cases, such as Hongkong, Singapore, London,  
33 the trips carried by public transport account for a larger proportion. Therefore, compared  
34 with other types of traffic flow, we believe that the mixed passenger flows in large-scale  
35 multi-modal public transit might be a better reflection of the citywide mobility dynamics.

36 Taking advantages of Singapore's distance-based AFC system, the smart card data record  
37 both passengers' boarding and alighting activities and more importantly it integrates both bus and  
38 metro trips in the same data frame, which allows us to put metro passenger flow and bus passenger  
39 flow together to conduct analysis. In this paper, we propose an end-to-end Pseudo-3D Convolu-  
40 tional Neural Network based framework that can collectively predict totally 7 different types of  
41 passenger flow at each region of the city area in the near future. The key contributions of this  
42 research summarized as follows:

- 1 • We propose a Pseudo-3D Convolutional Neural Network (Pseudo-3DCNN) based model  
2 to predict the public transport passenger flow in a network-wide region level.
- 3 • We take metro passenger flow, bus passenger flow as well as the transfer flow between  
4 metro system and bus system together into consideration instead of merely predicting a  
5 single type of passenger flow.
- 6 • The feasibility and effectiveness of our proposed model has been demonstrated on real-  
7 world data collected in Singapore.

8 The remainder of this paper is organized as follows. Section 3 elaborates all the prelim-  
9 inaries and methodologies. Section 4 describes the experiment settings and analyzes the results.  
10 Lastly, section 5 summarizes the main findings and discusses the future directions.

## 11 METHODOLOGY

12 In this section, we first give a detailed statement for all relevant definitions and concepts in this  
13 mixed passenger flow prediction task, and then we introduce the core methodologies in this re-  
14 search.

### 15 *Prediction Scope*

16 In this research, we predict the mixed passenger flow at each region of the city, which means  
17 that the prediction is performed in region level( we term as network-wide region level) instead of  
18 stop/station level. We first uniformly partition the entire city area into grid-like form with shape  
19 of  $I \times J$ , hence each of the grid refers to a region in this study. Then, we assign the bus stops  
20 and metro stations to corresponding regions based on the longitude and latitude coordinates of the  
21 stops/stations. The assignment procedure is illustrated as follows:

$$G_{i,j} = \{s_{bus} \in BS | s_{bus}^c \in (i,j)\} \cup \{s_{metro} \in MS | s_{metro}^c \in (i,j)\} \quad (1)$$

22 Where  $BS$  and  $MS$  denote the set of all the bus stops and metro stations, respectively.  $s_{bus}$   
23 refers to an element of set  $BS$ , a particular bus stop, and  $s_{bus}^c$  is the coordinate of  $bs$ . Similarly,  
24  $s_{metro}$  represents a metro station in the set of  $BM$ , of which the coordinate is denoted by  $s_{metro}^c$ .  
25 Therefore,  $G_{i,j}$  is the collection of all bus stops and metro stations that are located in grid  $(i,j)$ .

### 26 *Mixed Passenger Flow*

27 Totally, we construct 7 different types of passenger flow/demand by simultaneously taking both  
28 metro and bus system into consideration. The first 4 types of passenger flows are relatively straight-  
29 forward, which are *bus boarding* demand, *bus alighting* demand, *metro boarding* demand and  
30 *metro alighting* demand, respectively.

31 We first obtain 4 sets that contain the corresponding feasible trips during time interval  $t$ ,  
32 which can be defined as:

$$D_t^{bus-board} = \{\forall d \in D | d^{bus-board-time} \in t\} \quad (2)$$

$$D_t^{bus-alight} = \{\forall d \in D | d^{bus-alight-time} \in t\} \quad (3)$$

$$D_t^{metro-board} = \{\forall d \in D | d^{metro-board-time} \in t\} \quad (4)$$

35

$$D_t^{metro-align} = \{\forall d \in D | d^{metro-align-time} \in t\} \quad (5)$$

1 Where  $D$  is the collection of all the trips.  $d$  refers to a record in  $D$  representing a complete trip made  
 2 by a passenger and this trip might consists of a sequence of subtrips (e.g. a bus trip followed by a  
 3 metro trip). Specifically, in a trip  $d$ , if there is a bus boarding activity that takes place during time  
 4 interval  $t$ , then  $d$  will be added to set  $D_t^{bus-board}$  (see Equation 2) or if there is a metro alighting  
 5 activity that takes place during time interval  $t'$ , then  $d$  will be also added to set  $D_{t'}^{metro-align}$  (see  
 6 Equation 5). The same rules applies for the other sets. Subsequently, the *bus boarding* demand, *bus*  
 7 *alighting* demand, *metro boarding* demand and *metro alighting* demand can be obtained according  
 8 to the following logic:

$$x_{t,i,j}^{bus-board} = Card(\{\forall d_t \in D_t^{bus-board} | d_t^{bus-board-stop} \in G_{i,j}\}) \quad (6)$$

9

$$x_{t,i,j}^{bus-align} = Card(\{\forall d_t \in D_t^{bus-align} | d_t^{bus-align-stop} \in G_{i,j}\}) \quad (7)$$

10

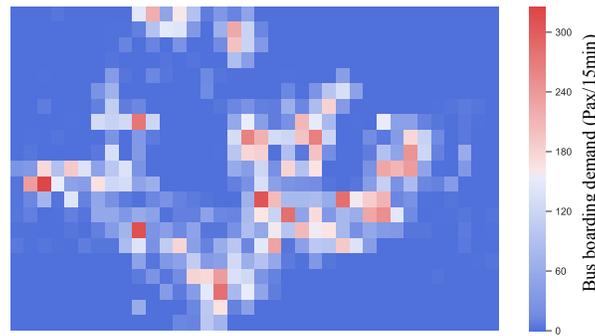
$$x_{t,i,j}^{metro-board} = Card(\{\forall d_t \in D_t^{metro-board} | d_t^{metro-board-station} \in G_{i,j}\}) \quad (8)$$

11

$$x_{t,i,j}^{metro-align} = Card(\{\forall d_t \in D_t^{metro-align} | d_t^{metro-align-station} \in G_{i,j}\}) \quad (9)$$

12

Take bus passenger for instance (see Equation 6 and Equation 7),  $D_t^{bus-board}$  represents the  
 13 set of all the trips that have a bus boarding activity during time interval  $t$  and similarly  $D_t^{bus-align}$   
 14 denotes the collection of all trips that have an alighting activity during time interval  $t$ .  $d_t$  is a  
 15 satisfied trip in  $D_t^{bus-board}$ , of which the bus boarding stop is denoted by  $d_t^{bus-board-stop}$ . Hence,  
 16  $\{d_t \in D_t | d_t^{bus-board-stop} \in G_{i,j}\}$  is the set consisting of all the trips that have a boarding activity in  
 17 grid  $(i, j)$  during time interval  $t$ . Lastly, we take the cardinality of this set to obtain the exact bus  
 18 boarding demand (total number of bus boarding passengers) in grid  $(i, j)$  during time interval  $t$ ,  
 19 which is denoted by  $x_{t,i,j}^{bus-board}$ . Thus,  $x_{t,i,j}^{bus-board}$  is a 2D flow matrix, shown in in Fig 1. The same  
 20 rule applies for the other types of passenger flow (see Equation 7, 8 and 9).



**FIGURE 1 : Flow matrix (bus boarding)**

21

In fact, the transfer activities are very common in a mature multi-modal public transit,  
 22 where there are lots of commuters that would take more than 1 mode of transportation to complete  
 23 the journey. Therefore, in addition to the 4 types of passenger flow mentioned above, we construct  
 24 another 3 different types of passenger flow by considering the transfer between different modes,

1 which are *bus-to-metro* demand, *metro-to-bus* demand and *bus-to-bus* demand, respectively. The  
 2 transfer related flows can be obtained by performing the following operations:

$$x_{t,i,j}^{bus-metro} = Card(\{\forall d_t \in \{D_t^{bus-align} \cap D_t^{metro-board}\} \\ |d_t^{bus-align-station} \in G_{i,j} \wedge d_t^{metro-board-station} \in G_{i,j}\}) \quad (10)$$

$$x_{t,i,j}^{metro-bus} = Card(\{\forall d_t \in \{D_t^{metro-align} \cap D_t^{bus-board}\} \\ |d_t^{metro-align-station} \in G_{i,j} \wedge d_t^{bus-board-station} \in G_{i,j}\}) \quad (11)$$

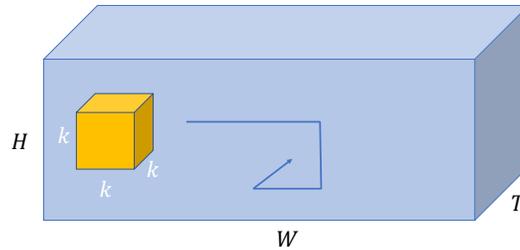
$$x_{t,i,j}^{bus-bus} = Card(\{\forall d_t \in \{D_t^{bus-align} \cap D_t^{bus-board}\} \\ |d_t^{bus-align-station} \in G_{i,j} \wedge d_t^{bus-board-station} \in G_{i,j}\}) \quad (12)$$

5 Where  $x_{t,i,j}^{bus-metro}$  refers to the *bus-to-metro* demand in grid  $(i, j)$  during time  $t$ . To be  
 6 specific,  $x_{t,i,j}^{bus-metro}$  is the total number of passengers who alight from a bus in grid  $(i, j)$  during time  
 7 interval  $t$  and then make a transfer to board a metro in the same grid during the same time interval.  
 8 Similarly,  $x_{t,i,j}^{metro-bus}$  and  $x_{t,i,j}^{bus-bus}$  represent the corresponding *metro-to-bus* demand and *bus-to-bus*  
 9 demand, respectively.

10 At a specific time interval  $t$ , the mixed passenger flow observation in all regions  $I \times J$   
 11 could be structured as a tensor  $X_t \in \mathbb{R}^{7 \times I \times J}$  by concatenating different types of passenger flow  
 12 matrix along the channel dimension, where  $(X_t)_{0,i,j} = x_{t,i,j}^{bus-board}$ ,  $(X_t)_{1,i,j} = x_{t,i,j}^{bus-align}$ ,  $(X_t)_{2,i,j} =$   
 13  $x_{t,i,j}^{metro-board}$ , ...,  $(X_t)_{6,i,j} = x_{t,i,j}^{bus-bus}$ . The objective of this framework is to simultaneously predict  
 14 different types of passenger flow at each region in the near future according to the mixed flow  
 15 tensors from the last few short-term periods. The input of the model is a tensor by concatenating a  
 16 sequence of historical observations,  $[X_{t-T}, \dots, X_{t-1}, X_t] \in \mathbb{R}^{7 \times T \times I \times J}$ , where  $T$  refers to the number of  
 17 historical time intervals that are taken into consideration. Besides, the predicted output is denoted  
 18 by  $\hat{Y} \in \mathbb{R}^{7 \times T' \times I \times J}$ , where  $T'$  is the number of predicted steps ahead.

### 19 3D Convolutional Neural Network

20 Convolutional Neural Network (CNN) have been widely applied for many tasks and shown a cou-  
 21 ple of superiorities, owing to its great ability of modelling spatial structures. In transportation  
 22 related prediction tasks, not only the spatial dependencies should be effectively captured, but also  
 23 the temporal correlations are of great importance. Unlike traditional 2DCNN that the convolution  
 24 operations are performed only spatially, 3DCNN can perform the convolution operations spatial-  
 25 temporally, which have exhibited significant effectiveness in some video/motion prediction tasks  
 26 (14); (15). Fig 2 demonstrates how 3D convolution performs on the tensor.



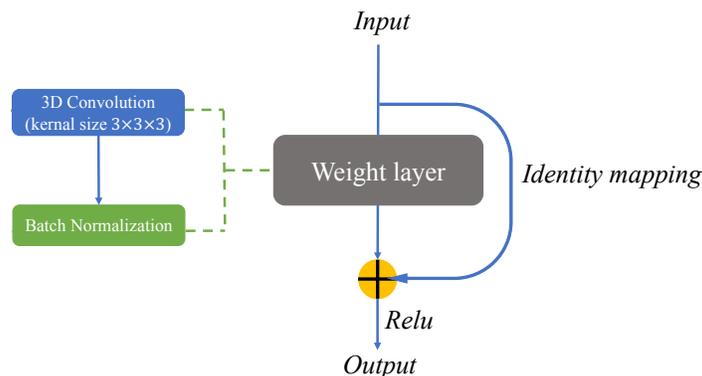
**FIGURE 2 : The illustration of 3D convolution**

1 As illustrated in Fig 2, 3DCNN applies a 3 dimensional filter to the tensor and the filter  
 2 moves in 3 dimensions to aggregate the nearby spatial-temporal features.

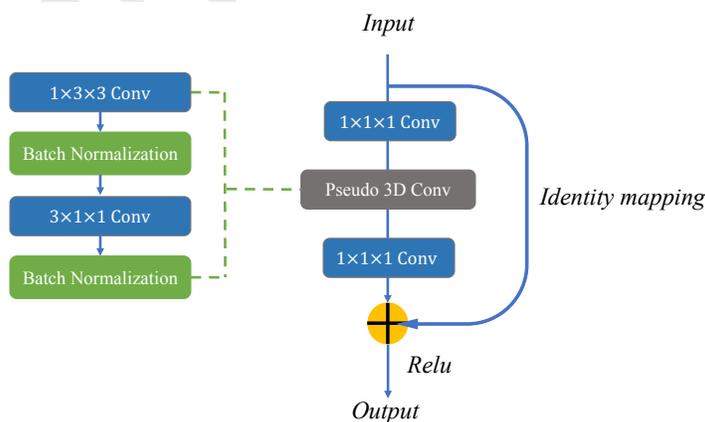
### 3 Pseudo-3DCNN Residual Learning

#### 4 *3D Residual Learning*

5 The local spatial patterns can be continuously captured through the sliding convolution operations.  
 6 Since many commuters are likely to take public transport for relatively longer-distance trips, con-  
 7 sequently the non-local or long-range spatial-temporal patterns might be more critical. Normally  
 8 there are several different ways to capture the long-range dependencies, such as downsampling  
 9 (e.g. pooling), applying fully connected layers as well as stacking multiple convolutional layers.  
 10 Given that applying downsampling strategies (e.g. pooling) might lead to the loss of information  
 11 and using fully connected layers would not only significantly increase the number of parameters  
 12 but also miss the spatial information. Hence, inspired by the deep residual learning mechanism  
 13 (16), we stack multiple residual units forming a deeper network in order to model the long-range  
 14 citywide dependencies. Fig 3a illustrates the structure of a residual unit.



(a) The structure of a 3DCNN residual unit



(b) The structure of a Pseudo-3D residual unit

**FIGURE 3 : The structure of the residual unit**

15 The identity mapping in the residual unit, defined in Equation 13, is to directly transmit  
 1 the activation from previous layers to deeper layers which contributes to ease the vanishing and

2 exploding gradient problems when training a deep network.

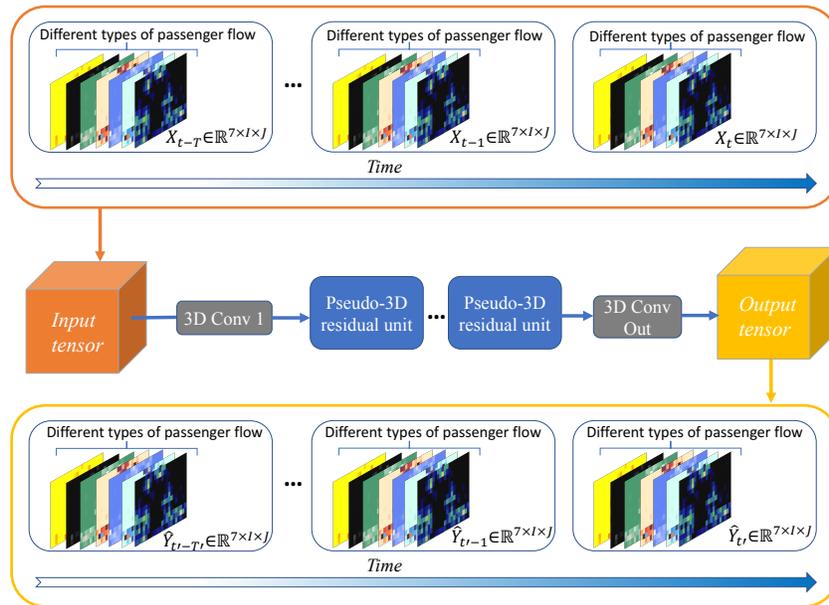
$$X^{l+1} = X^l + F(X^l) \quad (13)$$

3 Where  $X^{l+1}$  and  $X^l$  are the output and input of the  $l$ -th residual unit. Besides, in this research,  $F$   
4 refers to the operation of 1 layer 3DCNN.

### 5 *Pseudo-3D Residual Learning*

6 However, compared with 2DCNN, training 3DCNN is very computationally expensive meanwhile  
7 the model size also has a quadratic growth (17). Thus, a relatively light weight framework with  
8 great prediction performance is of better efficiency and scalability, especially when being applied  
9 or deployed in practice. Fortunately, Qiu et al. (17) proposed a Pseudo-3DCNN residual architec-  
10 ture that can learn spatial-temporal video representation with greater efficiency by simplifying the  
11 normal 3DCNN which is computational expensive. Hence, we adopt the concept of the Pseudo-  
12 3DCNN to construct our prediction framework. In pseudo-3DCNN (17), a normal 3DCNN oper-  
13 ation with kernel size  $3 \times 3 \times 3$  could be decomposed to a combination of a 2D spatial CNN with  
14 kernel size  $1 \times 3 \times 3$  and a 1D temporal CNN with kernel size  $3 \times 1 \times 1$ . Therefore, we first replace  
15 the normal 3DCNN in our original residual units (Fig 3a) with two consecutive convolution oper-  
16 ations, which are Spatial-CNN ( $1 \times 3 \times 3$ ) and Temporal-CNN ( $3 \times 1 \times 1$ ), respectively. Besides,  
17 we add another two  $1 \times 1 \times 1$  convolutions at both ends of the unit for reducing and rescaling the  
18 tensor dimensions, which helps to further reduce the computational costs. The revised residual  
19 unit is illustrated in Fig 3b.

20 Fig 4 presents the schematic architecture of our proposed model which we term as *ST-*  
21 *Pseudo3D Net*.



**FIGURE 4 : The schematic architecture of the proposed framework**

22 In the framework, the input tensor  $X \in \mathbb{R}^{7 \times T \times I \times J}$  is first transformed by *conv-1* (3DCNN  
23 with kernel size  $3 \times 3 \times 3$ ) to generate the internal feature maps which will be subsequently fed into

1 the stacked residual units followed by a final output 3DCNN layer to obtain the predicted result  
 2  $\hat{Y} \in \mathbb{R}^{7 \times T' \times I \times J}$ . The entire computation flow could be written as:

$$a^l = \text{ReLU}(\text{BN}(W^l * X + b^l)), \quad l = 1 \quad (14)$$

$$a^l = a^{l-1} + F(a^{l-1}; \theta^l), \quad \forall l = 2, \dots, L-1 \quad (15)$$

$$\hat{Y} = W^L * a^{L-1} + b^L, \quad l = L \quad (16)$$

5 Equation 14 represents the computation in the first 3DCNN layer  $\text{conv} - 1$ , where  $*$  denotes the  
 6 3D convolution operation,  $a^l$  is the output of layer  $l$ . Both  $W^l$  and  $b^l$  are trainable parameters at  
 7 layer  $l$ .  $\text{BN}$  refers to the batch normalization operation (18) which helps to accelerate the training  
 8 by reducing internal covariate shift and  $\text{ReLU}$  is the non-linear activation function,  $\text{ReLU}(x) =$   
 9  $\max(0, x)$ . Equation 15 outlines the computation flow through the stacked residual units, where  $a^l$   
 10 and  $a^{l-1}$  are the activation from layer  $l$  and  $l-1$ , respectively.  $F$  is the residual function, which  
 11 is a sequence of convolution operations in our research (see Fig 3b). Finally, the output layer  
 12 is illustrated in Equation 16, where  $*$  denotes the 3D convolution operation,  $W^L$  and  $b^L$  are both  
 13 learnable parameters. Notably, we just use the direct linear output from  $W^L * a^{L-1} + b^L$  as the final  
 14 predicted output without applying any non-linear activation at the last layer.

## 15 EXPERIMENTS AND RESULTS

16 In this section, we first introduce the experimental settings and then use Singapore as the case to  
 17 evaluate the performance of the proposed framework.

### 18 Data Preparation

19 The data set used for training and evaluation is collected from Singapore's AFC system (19th  
 20 March 2012 to 25th March 2012). Singapore adopts the distance-based AFC system and integrates  
 21 both bus trips and metro trips in the same data frame which allows us to conduct research on the  
 22 multi-modal related issues. The exact study area of this particular research starts from 103.678  
 23 to 104.014 in longitude and from 1.255 to 1.448 in latitude, which is uniformly partitioned into  
 24  $36 \times 21$  grids. Thus, the size of each grid equals to  $0.0093 \text{ lng} \times 0.992 \text{ lat}$ , which is about  $1.03 \text{ km} \times$   
 25  $1.02 \text{ km}$ . Additionally, in the experiment, we use the historical observations in the last 90 minutes  
 26 by a 15-min interval to predict the passenger flows in the next 1 hour by a 15-min. Notably, the  
 27 different types of passenger flow might have significant differences in magnitude. For example,  
 28 the metro boarding demand tend to be much larger than the bus-to-bus transfer demand. Therefore,  
 29 we first apply the Min-Max normalization (illustrated in Equation 17) to scale the different types  
 30 of passenger flow into the same range ( $[0, 1]$ ) in order to improve training efficiency.

$$z = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (17)$$

### 31 Training Settings

32 The randomly split is applied on the data set where 85% of the data are used for training while  
 33 another 15% are used for evaluation. The batch size is 32 and the maximum training epochs is  
 1 fixed to 500. Additionally, the early-stopping mechanism is introduced to continuously monitor  
 2 the loss on the validation set in order to avoid overfitting issue. Root Mean Square Error (RMSE)  
 3 is the evaluation metric for the proposed model, given by:

$$z = \sqrt{\frac{1}{N} \sum_{n=1}^N (Y^{(n)} - \hat{Y}^{(n)})^2} \quad (18)$$

4 The model is optimized by Adam optimization algorithm (19) with learning rate of 0.001.  
 5 All experiments are conducted on our Google Cloud Platform, of which the detail configuration  
 6 are listed in Table 1.

**TABLE 1 : Experiment environment**

Item	Configuration
OS	GNU/Linux server
CPU version	Intel (R) Xeon (R) CPU @2.20GHz
RAM	15GB
GPU	2 × NVIDIA Tesla T4
Pytorch version	1.1.0
CUDA version	9.0

## 7 Results Analysis

8 We compare our proposed model *ST-Pseudo3D Net* with a group of state-of-the-art benchmarks,  
 9 including DNN, RNN, LSTM (20), GRU (21), 2DCNN with residual connection (Res2DCNN)  
 10 and 3DCNN with residual connection (Res3DCNN). In order to show the results in a more metic-  
 1 ulous manner, we decompose the overall prediction results into several groups with respect to the  
 2 type of the passenger flow, which are presented in Table 2-4. The best score in each passenger flow  
 3 type is highlighted in bold.

**TABLE 2 : Results comparison for bus passenger flow**

Model	RMSE	
	<i>Bus boarding demand</i>	<i>Bus alighting demand</i>
DNN	12.51	12.80
RNN	14.33	14.71
LSTM	13.65	14.35
GRU	13.84	14.51
Res2DCNN	10.56	11.93
Res3DCNN	9.57	10.73
<b>ST-Pseudo3D</b>	<b>9.15</b>	<b>10.55</b>

**TABLE 3 : Results comparison for metro passenger flow**

Model	RMSE	
	<i>Metro boarding demand</i>	<i>Metro alighting demand</i>
DNN	15.93	<b>19.62</b>
RNN	16.49	20.84
LSTM	<b>15.46</b>	20.10
GRU	15.84	20.29
Res2DCNN	17.91	21.82
Res3DCNN	15.82	20.26
ST-Pseudo3D	15.98	20.33

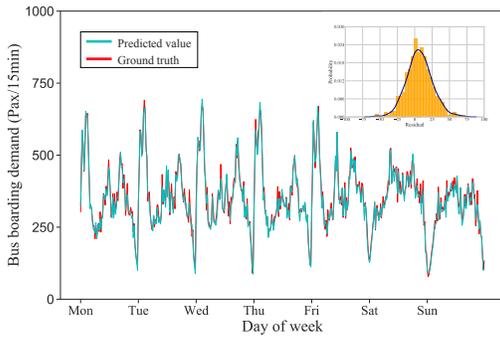
**TABLE 4 : Results comparison for transfer passenger flow**

Model	RMSE		
	<i>Bus-to-bus demand</i>	<i>Bus-to-metro demand</i>	<i>Metro-to-bus demand</i>
DNN	7.87	10.24	9.93
RNN	7.99	10.98	10.46
LSTM	7.76	10.69	10.07
GRU	7.81	10.78	10.26
Res2DCNN	7.00	9.68	9.03
Res3DCNN	6.41	<b>7.81</b>	7.45
ST-Pseudo3D	<b>6.33</b>	7.89	<b>7.38</b>

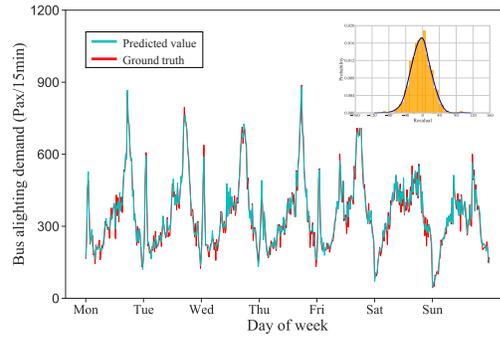
4 As shown in Table 2 - 4, our proposed model *ST-Pseudo3D Net* outperforms other baselines  
5 in most of the above-listed prediction tasks. Specifically, for bus passenger flow prediction, CNN  
6 based models perform better than DNN and other RNN based models, owing to the effectiveness in  
7 extracting spatial features, where the proposed *ST-Pseudo Net* achieve the best results for both *bus*  
8 *boarding demand* prediction and *bus alighting demand* prediction, which are relatively 4.4% up to  
9 36.1% better than other baselines in predicting *bus boarding demand* and 1.7% up to 28.3% better  
10 than other baselines in predicting *bus alighting demand* at each region of the city. When it comes  
11 to metro passenger flow, we observe that the DNN and RNN based models generally perform bet-  
12 ter than CNN based models, where our proposed model slightly lags behind LSTM by 3.3% and  
13 DNN by 3.5% in predicting *metro boarding demand* and *metro alighting demand*, respectively.  
14 However, for transfer passenger flow, the CNN based models, especially 3DCNN based models  
15 show significant superiorities again, where 3DCNN based models are relatively 8.4% up to 19.7%,  
16 18.5% up to 28.9% and 18.3% up to 29.4% better than other baselines in predicting *bus-to-bus*  
17 *demand*, *bus-to-metro demand*, *metro-to-bus demand*, respectively. Notably, although Res3DCNN  
18 almost perform as good as the proposed ST-Pseudo3D, the computation cost and memory de-  
19 mand of Res3DNN are much higher than ST-Pseudo3D. Specifically, in the experiment, we keep  
20 Res3DCNN and ST-Pseudo3D with the same depth (stacking 16 residual units) and the total num-  
21 ber of trainable parameters of ST-Pseudo3D is 46.1% less than Res3DCNN. Therefore, we believe  
22 that the proposed framework is superior to Res3DCNN in terms of efficiency and practicability.

### 23 Case Analysis

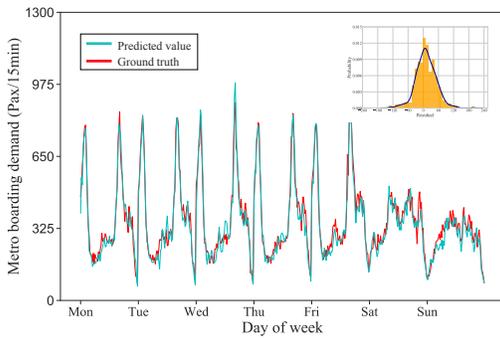
24 In order to show the prediction results of the proposed framework intuitively, a specific region is  
25 selected for detailed performance demonstration. The selected region is located in Clementi area  
26 which is one of the major residential towns in Singapore. The selected region covers 1 metro  
1 station (Clementi station), 1 integrated bus interchange hub and several bus stops, which makes it  
2 an ideal spot for results analysis. Both the true and predicted demand are plotted in Fig 5 according  
3 to different types of passenger flow.



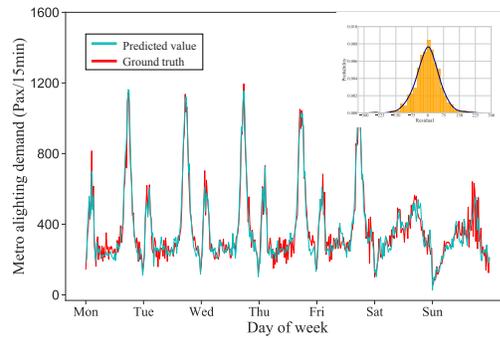
(a) Bus boarding demand comparison



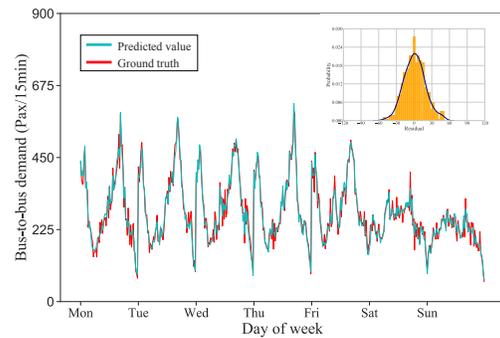
(b) Bus alighting demand comparison



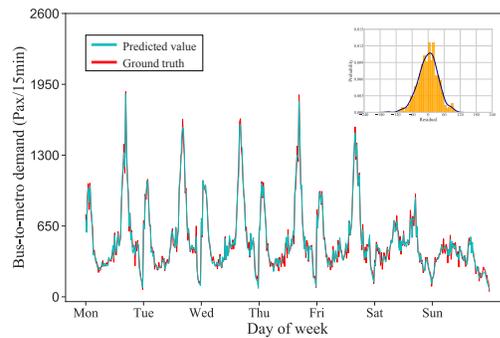
(c) Metro boarding demand comparison



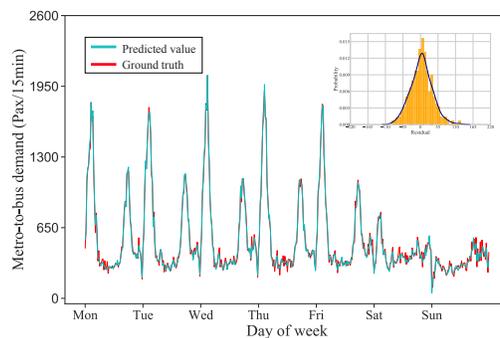
(d) Metro alighting demand comparison



(e) Bus-to-bus demand comparison



(f) Bus-to-metro demand comparison



(g) Metro-to-bus demand comparison

**FIGURE 5 : Comparisons of the ground truth and predicted passenger flow**

4 As illustrated in Fig 5a - 5g, the proposed model can make relatively reliable prediction  
5 for different types of passenger flow. Furthermore, since different prediction tasks have different  
6 numerical scales, so a small *RMSE* or *MAE* value does not necessarily mean that the model is  
7 good enough. Therefore, we further analyze the residual of the prediction, which are plotted in  
8 the upper right corner of each figure. We observe that the residuals of our prediction for different  
9 types of passenger flow basically follow normal distributions with zero mean, which indicates that  
10 there are no significant non-random patterns in the residuals.

## 11 CONCLUSION

12 In this paper, we propose an end-to-end deep learning based framework to collectively predict  
13 network-wide multi-modal passenger flows in a large-scale public transit. Considering metro pas-  
14 senger flow, bus passenger flow as well as transfer flow between metro and bus system, we to-  
15 tally construct 7 different types of passenger flow to be collectively predicted. The framework is  
16 based on deep 3DCNN residual architecture which can effectively extract spatial-temporal fea-  
17 tures. Moreover, inspired by Pseudo-3D, we replace the traditional 3D convolution operation  
18 in the residual units with two consecutive convolution operations (a Spatial-CNN followed by a  
19 Temporal-CNN), which significantly reduce the number of parameters and the computational cost.  
20 Additionally, the model is validated by real-world data collected from Singapore's public transit  
21 and the experiment results demonstrate that the proposed framework is superior to other baselines  
22 in terms of the accuracy and practicability.

23 In the future, we plan to train the model on some larger data sets in order to further improve  
24 its generalization ability. In addition, combining with the theories from graph learning domain, we  
25 will also try to explore the movement of passenger flows in the network.

## 1 ACKNOWLEDGEMENT

2 We would like to show our gratitude to Land Transport Authority of Singapore for providing the  
3 data.

#### 4 REFERENCES

- 5 [1] Ma, X., H. Yu, Y. Wang, and Y. Wang, Large-scale transportation network congestion evolu-  
6 tion prediction using deep learning theory. *PloS one*, Vol. 10, No. 3, 2015, p. e0119044.
- 7 [2] Lv, Y., Y. Duan, W. Kang, Z. Li, F.-Y. Wang, et al., Traffic flow prediction with big data:  
8 A deep learning approach. *IEEE Trans. Intelligent Transportation Systems*, Vol. 16, No. 2,  
9 2015, pp. 865–873.
- 10 [3] Polson, N. G. and V. O. Sokolov, Deep learning for short-term traffic flow prediction. *Trans-*  
11 *portation Research Part C: Emerging Technologies*, Vol. 79, 2017, pp. 1–17.
- 12 [4] Cui, Z., R. Ke, and Y. Wang, Deep Bidirectional and Unidirectional LSTM Recurrent Neural  
13 Network for Network-wide Traffic Speed Prediction. *CoRR*, Vol. abs/1801.02143, 2018.
- 14 [5] Zhao, Z., W. Chen, X. Wu, P. C. Chen, and J. Liu, LSTM network: a deep learning approach  
15 for short-term traffic forecast. *IET Intelligent Transport Systems*, Vol. 11, No. 2, 2017, pp.  
16 68–75.
- 17 [6] Fu, R., Z. Zhang, and L. Li, Using LSTM and GRU neural network methods for traffic  
18 flow prediction. In *2016 31st Youth Academic Annual Conference of Chinese Association of*  
19 *Automation (YAC)*, IEEE, 2016, pp. 324–328.
- 20 [7] Yu, H., Z. Wu, S. Wang, Y. Wang, and X. Ma, Spatiotemporal recurrent convolutional net-  
21 works for traffic prediction in transportation networks. *Sensors*, Vol. 17, No. 7, 2017, p. 1501.
- 22 [8] Wang, J., R. Chen, and Z. He, Traffic speed prediction for urban transportation network: A  
23 path based deep learning approach. *Transportation Research Part C: Emerging Technologies*,  
24 Vol. 100, 2019, pp. 372–385.
- 25 [9] Liu, L. and R.-C. Chen, A novel passenger flow prediction model using deep learning meth-  
26 ods. *Transportation Research Part C: Emerging Technologies*, Vol. 84, 2017, pp. 74–91.
- 27 [10] Bai, Y., Z. Sun, B. Zeng, J. Deng, and C. Li, A multi-pattern deep fusion model for short-term  
28 bus passenger flow forecasting. *Applied Soft Computing*, Vol. 58, 2017, pp. 669–680.
- 29 [11] Ma, X., J. Zhang, B. Du, C. Ding, and L. Sun, Parallel Architecture of Convolutional Bi-  
30 Directional LSTM Neural Networks for Network-Wide Metro Ridership Prediction. *IEEE*  
31 *Transactions on Intelligent Transportation Systems*, 2018.
- 32 [12] Liu, Y., Z. Liu, and R. Jia, DeepPF: A deep learning based architecture for metro passenger  
33 flow prediction. *Transportation Research Part C: Emerging Technologies*, Vol. 101, 2019, pp.  
34 18–34.
- 35 [13] Zhang, J., Y. Zheng, and D. Qi, Deep Spatio-Temporal Residual Networks for Citywide  
36 Crowd Flows Prediction. In *AAAI*, 2017, pp. 1655–1661.
- 1 [14] Ji, S., W. Xu, M. Yang, and K. Yu, 3D convolutional neural networks for human action  
2 recognition. *IEEE transactions on pattern analysis and machine intelligence*, Vol. 35, No. 1,  
3 2012, pp. 221–231.

- 4 [15] Tran, D., L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, Learning spatiotemporal fea-  
5 tures with 3d convolutional networks. In *Proceedings of the IEEE international conference*  
6 *on computer vision*, 2015, pp. 4489–4497.
- 7 [16] He, K., X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition. In *Pro-*  
8 *ceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–  
9 778.
- 10 [17] Qiu, Z., T. Yao, and T. Mei, Learning spatio-temporal representation with pseudo-3d residual  
11 networks. In *proceedings of the IEEE International Conference on Computer Vision*, 2017,  
12 pp. 5533–5541.
- 13 [18] Ioffe, S. and C. Szegedy, Batch normalization: Accelerating deep network training by reduc-  
14 ing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- 15 [19] Kingma, D. P. and J. Ba, Adam: A method for stochastic optimization. *arXiv preprint*  
16 *arXiv:1412.6980*, 2014.
- 17 [20] Hochreiter, S. and J. Schmidhuber, Long short-term memory. *Neural computation*, Vol. 9,  
18 No. 8, 1997, pp. 1735–1780.
- 402 [21] Cho, K., B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and  
403 Y. Bengio, Learning phrase representations using RNN encoder-decoder for statistical ma-  
404 chine translation. *arXiv preprint arXiv:1406.1078*, 2014.