



Applicability of Long Short-Term Memory Traffic Volume Imputation Model to Drive Connected Corridor Simulation

Abhilasha Saroj^{a,*}, Angshuman Guin^a, Michael Hunter^a

^a School of Civil and Environmental Engineering, Georgia Institute of Technology, 788 Atlantic Drive NW, Atlanta, GA 30332; *Corresponding Author

Received: 12 Dec. 2020 • Revised: 28 Dec. 2020 • Accepted: 29 Dec. 2020 • Published: 29 Dec. 2020

© 2020 Mason Publishing Group (a division of George Mason University libraries); Sponsor: [ATPIO](#)



ABSTRACT

For effective implementation of connected corridor applications, it is imperative to study the characteristics of the high-resolution connected corridor data streams leveraged in smart city applications. In a previous effort, a smart city application – real-time corridor data-driven traffic simulation model, i.e., Digital Twin – is developed. Investigation of the corridor field volume data revealed the presence of data gaps. To address these gaps, deep Long Short-Term Memory (LSTM) Recurrent Neural Network univariate and multivariate volume imputation models are developed. In this paper, the impact of the developed model imputations on the digital twin generated travel times are investigated. Simulation runs are conducted for typical and atypical traffic, for three volume input cases: base (original volumes), univariate model imputations, and multivariate model imputations. For the given methodology it was seen that: 1) the travel times generated using multivariate imputations are the closest to that generated using base data, 2) the impact of imputations on travel times is focused on congested routes, and 3) the impact on travel time is minimal despite input volume overestimation on routes that have the capacity to accommodate higher volumes. These findings demonstrate the need to prioritize data streams based on the given application and underlying corridor conditions.

Keywords: smart cities, connected corridor, long short-term memory, real-time simulation, traffic volume imputation

1. INTRODUCTION

Smart cities across the world utilize smart corridor testbeds to explore technology implementations [1, 2, 3, 4, 5]. Often, a smart corridor is equipped with communications technologies [3], enabling the transfer of significant data between vehicles, the infrastructure, and corridor management centers. These data can be in different forms, such as connected vehicle data providing high resolution instantaneous vehicle specific data and signal phase and timing data, vehicle counts from in-road or roadside detectors, probe vehicle data such as that from INRIX [6], HERE [7], etc., to name a few. Smart corridor applications seek to convert these data into actionable information, to improve corridor performance. However, the presence of data gaps in the data streams can impair such efforts. Thus, it is imperative to develop data imputation methodologies as well as to understand the impact of such imputation on the application performance.

In a previous effort the authors developed a smart corridor application, a real-time data-driven traffic simulation model, i.e., Digital Twin, for the North Avenue Smart Corridor in Atlanta, Georgia [8, 9]. The Digital Twin, driven using high frequency volume and signal data, is capable of dynamically providing corridor traffic and environmental performance measures [8, 9]. However, investigation of the corridor data streams revealed the presence of data gaps. To address the volume data gaps, bi-directional Long Short-Term Memory (LSTM) Recurrent Neural Networks (RNNs) univariate and

multivariate imputation models were developed. Experiments were conducted to investigate the LSTM RNN model performance under typical and atypical day conditions [10]. Of specific interest was exploring if supplementing historic data from the given data stream with the most recent data from similar data streams (multivariate model) provides superior predictions over utilizing only historic data (univariate model). Experimental results indicated the potential for a multivariate LSTM RNN model to provide reasonable imputations on typical and atypical days [10].

In this paper the reasonableness of the data estimations is explored by identifying limitations and evaluating the appropriateness of the imputation models to drive the Digital Twin. A simulation experiment is conducted using the Digital Twin to investigate the impact of the univariate and multivariate model imputations on generated travel times for selected corridor routes.

2. LSTM RNN MODEL DEVELOPMENT

RNNs are a variant of Neural Networks, capable of utilizing the “memory” of previous event data to predict the next values in a sequence [11]. However, RNNs can suffer from the vanishing gradient problem in backpropagation implementation [12]. This may be a drawback when accounting for long-term dependency in a sequence is crucial to prediction accuracy. LSTM RNN [13] seeks to address this issue through the inclusion of a

‘memory’ cell component along with gates to regulate the memory cell value [14]. A variant of LSTM RNN is bidirectional LSTM RNN (BLSTM), where output mapping may learn from both past and future information [15].

In the previous study [10], deep bidirectional LSTM RNNs were used to develop univariate and multivariate volume time series prediction models for six selected detectors on three approaches on the North Avenue corridor, a 2.3-mile long actuated corridor, as shown in Figure 1. Each of these approaches has two lanes, referenced as L_1 and L_2. The multivariate models are trained using the historic data of the detector experiencing data loss as well as data from a corridor detector drawn from a cluster of detectors that have been identified to have a similar time series data pattern using cluster analysis. For a comprehensive literature review on time series similarity measures, traffic data imputation methodologies, and the LSTM RNN model development process, the reader is referred to Saroj [10].

3. EXPERIMENT DESIGN

A simulation experiment is designed to study the impact of the previously developed univariate and multivariate prediction models on simulation generated performance measures for a typical weekday, Monday, March 18th, 2019, and a weekday with atypical traffic conditions, Monday, May 27th, 2019, (Memorial Day). For each of these days the PM peak hours (3 PM to

6 PM) are simulated for three traffic volume sets, input at the three corridor approaches: 1) base traffic condition (original volume), 2) univariate model imputations, and 3) multivariate model imputations. A discussion of the base traffic volumes and signal timings for each experiment day may be found in [9]. The second and third volumes cases assume a three-hour data gap in the base traffic data at the three study approaches, utilizing the imputed volumes for these locations. Volumes during simulation run-time are imputed (i.e., predicted) as would occur in a real-time event, that is, the simulation model and algorithms are only fed data up to the equivalent wall clock time, i.e., the actual time in the field. Imputations are then based on the current wall clock (real-time) data and previous (historical) data. For each of the three data cases, for each of the two traffic days, ten replicate simulation trials are run to evaluate the impact on travel times on the nine corridor routes, i.e., the three studied side-street approaches and the six mainline routes (Figure 1).

4. RESULTS AND DISCUSSION

The developed univariate and multivariate models for each of the six detectors are used to predict volumes from 3 PM to 6 PM. Table 1 presents the performance error measures for the model predictions.

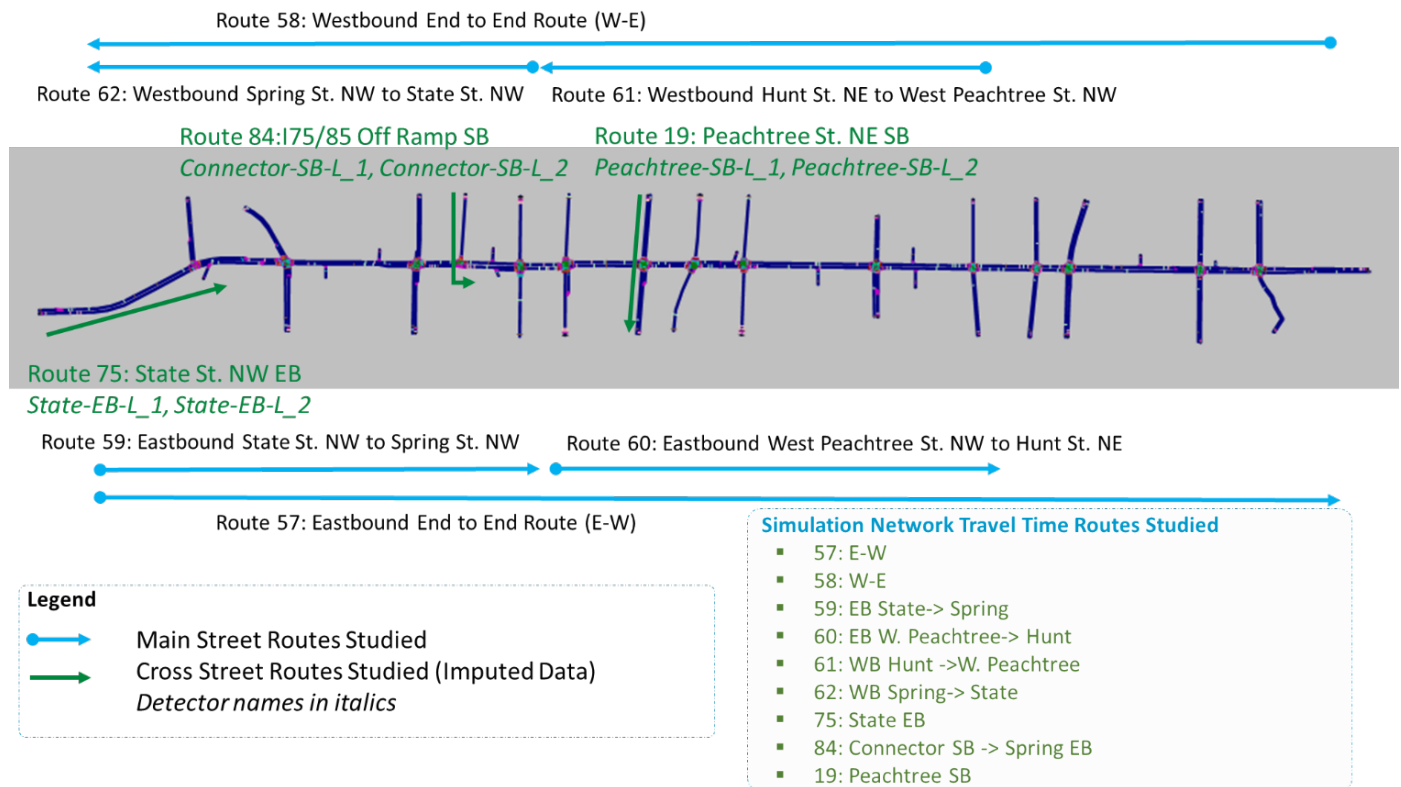


Figure 1. LSTM RNN model developed for three North Ave. corridor approach locations and the nine studied routes [10]

Table 1. Error measures for LSTM RNN model predictions for 3 PM to 6 PM.

| Detector | Model Type | March 18 th (Typical Day) | | | May 27 th (Atypical Day) | | |
|------------------|-----------------|--------------------------------------|------|----------|-------------------------------------|------|----------|
| | | MAE | RMSE | Std. Dev | MAE | RMSE | Std. Dev |
| State-EB-L_1 | Univariate | 5.2 | 6.3 | 6.3 | 20.6 | 22.7 | 9.6 |
| State-EB-L_1 | *Multivariate** | 5.0 | 6.1 | 6.1 | 5.2 | 6.1 | 4.0 |
| State-EB-L_2 | *Univariate** | 4.8 | 6.0 | 6.0 | 32.0 | 33.1 | 8.4 |
| State-EB-L_2 | Multivariate** | 5.4 | 7.24 | 7.2 | 16.4 | 18.2 | 7.8 |
| Connector-SB-L_1 | Univariate | 32.1 | 38.3 | 21.5 | 40.5 | 42.9 | 16.1 |
| Connector-SB-L_1 | *Multivariate** | 19.4 | 26.2 | 20.7 | 7.0 | 8.2 | 8.0 |
| Connector-SB-L_2 | *Univariate** | 8.7 | 11.1 | 10.7 | 12.2 | 15.0 | 13.0 |
| Connector-SB-L_2 | Multivariate | 10.0 | 12.8 | 12.8 | 23.1 | 25.1 | 9.8 |
| Peachtree-SB-L_1 | *Univariate** | 6.4 | 8.6 | 7.4 | 9.2 | 11.1 | 8.1 |
| Peachtree-SB-L_1 | Multivariate | 7.2 | 8.6 | 7.8 | 9.8 | 11.3 | 7.9 |
| Peachtree-SB-L_2 | *Univariate** | 6.9 | 8.4 | 7.5 | 4.6 | 6.0 | 5.6 |
| Peachtree-SB-L_2 | Multivariate | 8.3 | 10.6 | 8.0 | 12.5 | 13.3 | 5.1 |

Notes:

- * An asterisk indicates lower error values among the two model types on typical day predictions
- ** Two asterisks indicate lower values among the two model types on atypical day predictions
- MAE: mean absolute error, RMSE: root mean square error, Std. Dev: Standard Deviation of Errors

It is observed that the multivariate and univariate predictions tend to be similar on the typical day, with univariate errors often lower than that of multivariate, consistent with previous research findings [10]. On the atypical day, it is expected that the multivariate model will provide improved imputation values compared to the univariate model, as the historical data is not consistent with current conditions. This is seen to be true for

most detectors. However, at Connector-SB-L_2 and Peachtree-SB-L_2 the univariate prediction errors are observed to be much lower than multivariate prediction errors. For example, for Connector-SB-L_2 Figure 2 shows the plots for observed traffic volumes from midnight to 6 PM for each day along with the univariate and multivariate model predictions from 3 PM to 6 PM.

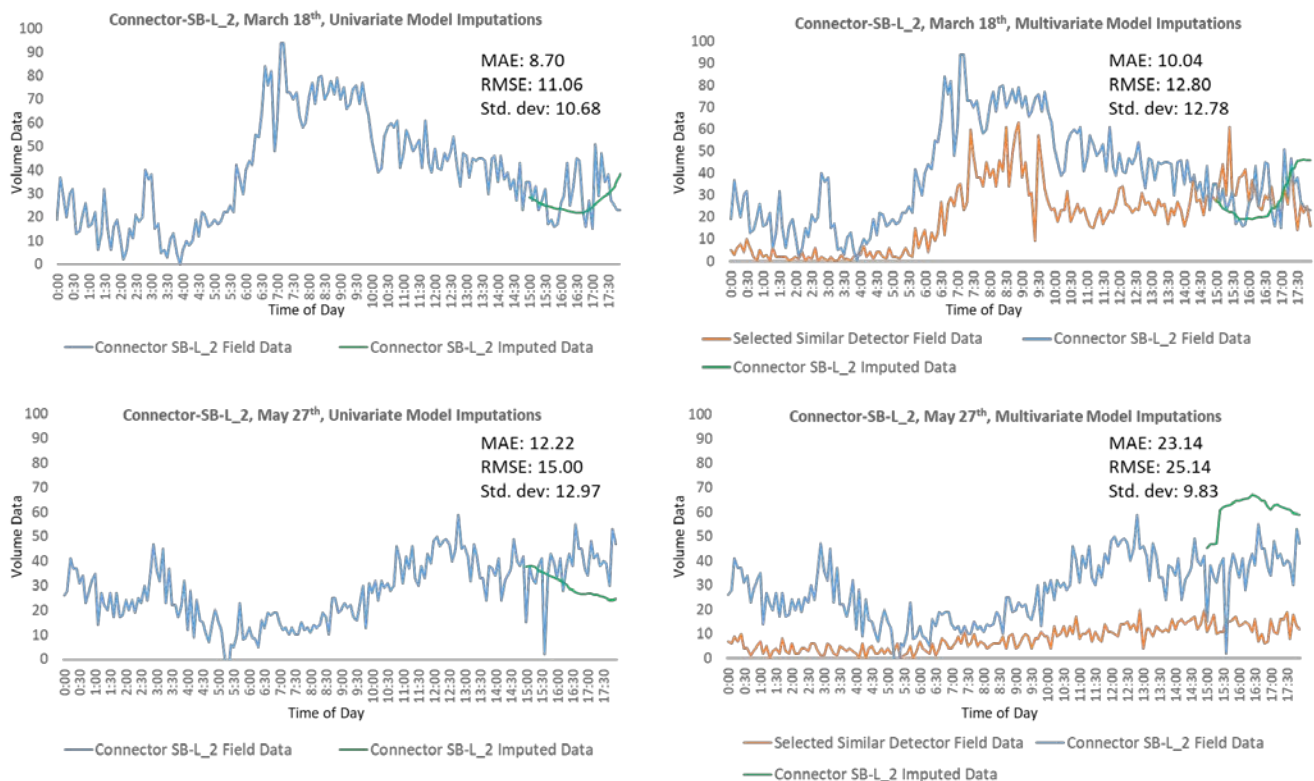


Figure 2. Model predictions for Connector-SB-L_2 starting at 3 PM.

As seen in Figure 2, one reason for the weaker performance of the multivariate model at Connector-SB-L_2 is a poor correlation between the traffic pattern on the matched detector and the given detector. The identification of detectors with similar patterns was undertaken utilizing data streams from multiple “typical” days [10]. This raises a possibility that detectors that are reasonably correlated under typical conditions may not be well correlated under atypical conditions. Thus, future improvements to the method may be achieved by identifying different matching detectors for different conditions.

In further exploring the simulation performance given the imputed data it is noted that the field volume data was available in six-minute bins. Therefore, the imputation was also set to generate six-minute binned data. Thus, in the simulation implementation the volume data is entered into the model in six-minute intervals, randomly distributed (shifted poisson distribution interarrival times) over the interval length. However, if the entry link is oversaturated (i.e., a vehicle queue extends to the link entrance point) the new vehicles will not be

able to enter the network during their set interval, and will instead be held until space becomes available. Figure 3 shows the volumes that entered the simulation model at Peachtree St. SB for the ten univariate and multivariate replicate trials, for the typical and atypical day scenarios, as well as the imputed volume that sought to enter. A low difference between the imputed and processed entry volume per six-minute interval suggests under-saturated conditions. However, the variation in the March 18th volume entry counts for the multivariate imputations suggests that this approach operates near saturation state during the typical day PM peak period, with the slightly higher multivariate imputed volumes sufficient to create over-saturated conditions. The variation in entry volumes given the univariate imputed values is significantly less, as the imputed volumes are lower than those of the multivariate model (Figure 4). Another clear observation from Figure 4 is that both imputation approaches have a tendency to smooth the volumes relative to the field conditions.

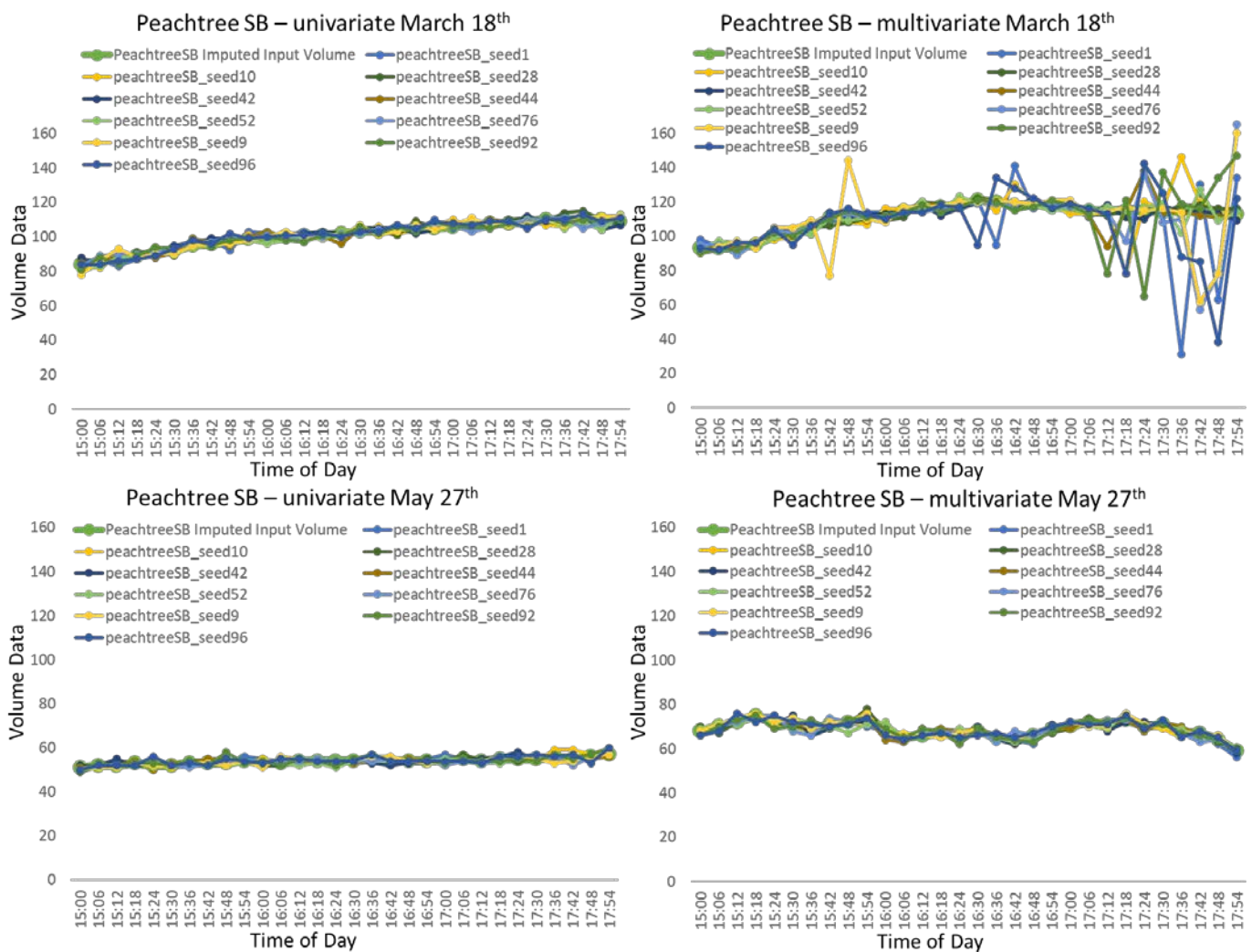


Figure 3. Imputed vs Entry volume (10 replicate seeds) for Peachtree St. SB approach volumes in the six-minute bins.

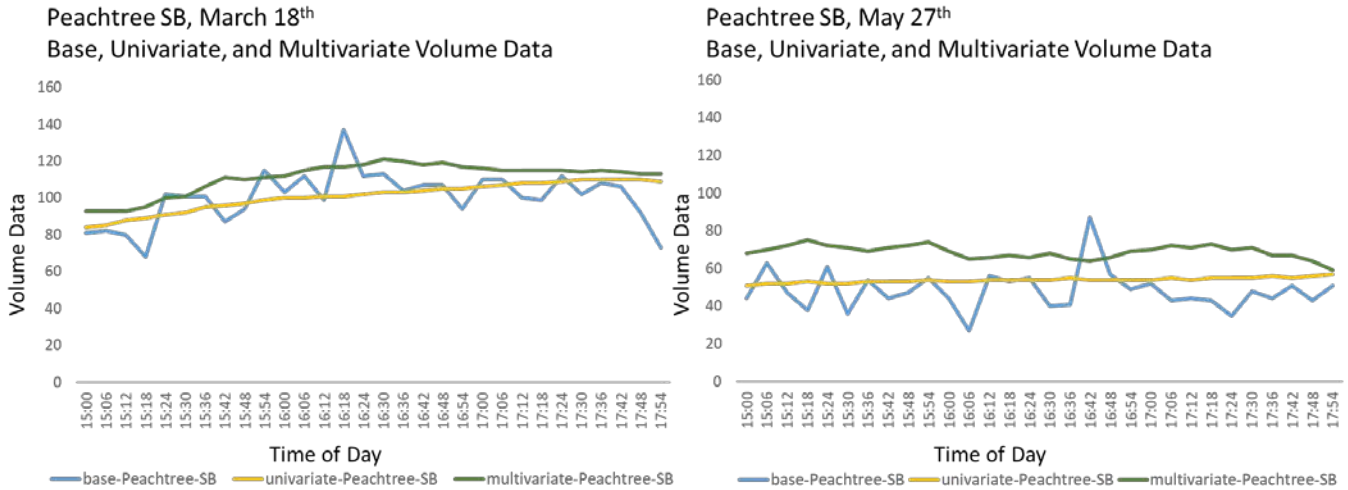


Figure 4. Approach volume for three cases for 3-6 PM at Peachtree St. SB.

4.1 Impact on Digital Twin Generated Travel Time

Figure 5 presents box plots of the 85th percentile travel times obtained from the replicate trials at the nine routes, for the three data input scenarios, for the typical and atypical days. It is observed that for the typical and atypical scenarios, travel times simulated using the multivariate imputations are generally closer to that of the base day than those simulated using the univariate imputations. For 8 of 9 routes under typical conditions multivariate provide closer results than univariate, reducing errors on average by 4%. Under atypical conditions, also for 8 of 9 routes, multivariate provides closer results, reducing errors on average by 3%.

The impact of underlying corridor demand, i.e., saturated vs under-saturated, can be seen on the simulated travel times. For example, lower travel time variation is seen on the atypical traffic day across cases, likely due to the lower holiday traffic. In addition, on Route 19, there are observable travel time differences for the three cases under typical conditions likely due to a saturated traffic state for the PM peak on the typical day versus the under-saturated holiday traffic. The low travel time differences for typical and atypical holiday traffic on Routes 75 and 84 are a result of under-saturated conditions on both days, even though for these routes the univariate imputation provides higher volume estimates than both the multivariate model and base data on the atypical day.

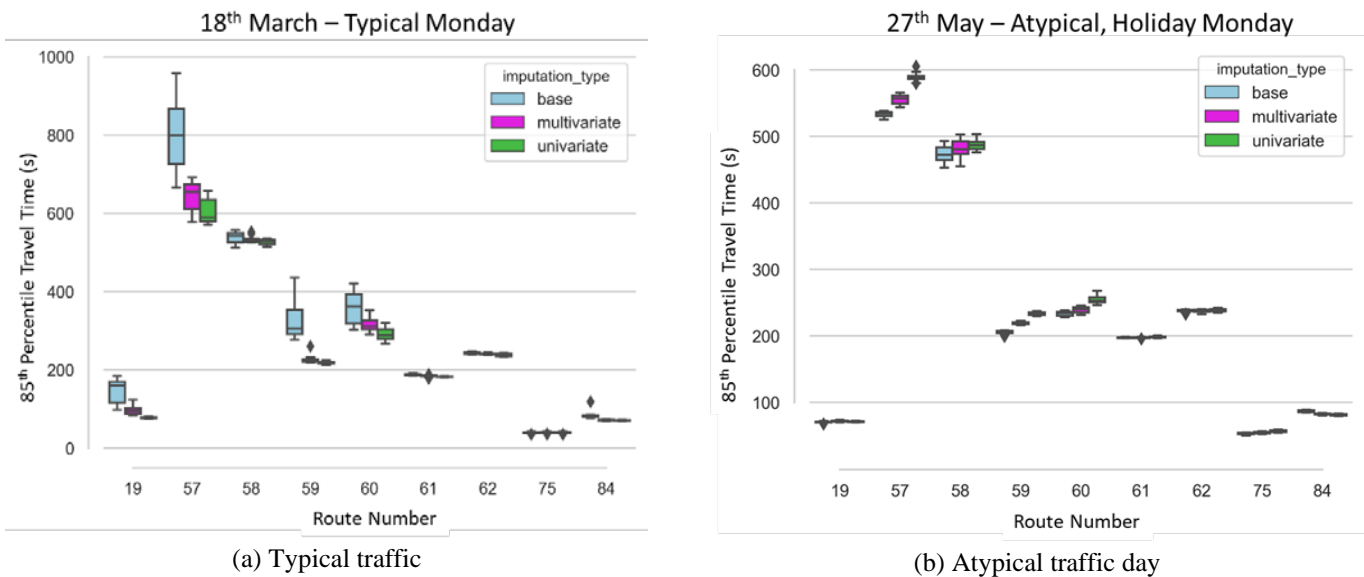


Figure 5. Boxplots of 85th percentile travel time at the nine study routes for (a) March 18th and (b) Monday, March 27th

Given the under-saturated conditions the overestimated volumes were not sufficiently erroneous to impact travel times. However, over-estimation of volumes is likely a factor that contributed to the increased travel time for the univariate model scenario compared to the base case on Routes of 57, 58, 59, and 60 (Figure 5b). Here, the volume estimation error combined with the underlying near-saturation conditions were sufficient to negatively influence the predicted travel times.

These observations clearly indicate that when developing smart applications, it is critical to identify those locations with the most potential to influence results. Key attributes of the applications, such as identification of a matching detector in the given example, should be assured as well as increased data control and data quality efforts at these locations.

5. CONCLUSIONS AND FUTURE WORK

This effort investigated the impact of the previously developed LSTM RNN multivariate and univariate model imputations on Digital Twin generated travel times. The results indicate that for the studied typical and atypical traffic, the multivariate imputations lead to simulated travel times that are closer to that of the base day. However, additional improvements in the multivariate method may be achieved by improved matching detector selection. Next, the importance of the underlying corridor conditions, i.e., saturation level, is observed. It is demonstrated that when developing smart applications both the imputation methodology and the local conditions must be considered. Specific to this effort, to improve the performance of the LSTM RNN models, future investigations may consider additional atypical training and test data as well as hyperparameter tuning.

6. ACKNOWLEDGMENTS

The information, data, or work presented herein was funded in part by the City of Atlanta (Research Project FC-9930- Smart Cities Traffic Congestion Mitigation Program) and in part by National Center of Sustainable Transportation (NCST). The authors thank the sponsors for their support of this research. Any opinions, findings, conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors. This paper does not constitute a standard, specification, or regulation.

REFERENCES

- [1] USDOT. *Connected Vehicle Pilot Deployment Program*. n.d. [cited 2020 February 26]; Available from: <https://www.its.dot.gov/pilots/>.
- [2] Frost, A. *Singapore to develop first 5G C-V2X research testbed on NTU campus*. 2019 October 22, 2019 [cited 2020 February 26]; Available from: <https://www.traffictoday.com/news/connected-vehicles-infrastructure/singapore-to-develop-first-5g-c-v2x-research-testbed-on-ntu-campus.html>.
- [3] CenterForTransportationResearch. *City of Austin Connected Corridors*. 2018 [cited 2020 February 26]; Available from: <https://ctr.utexas.edu/nmc/research-2/projects/current-and-ongoing-projects/city-of-austin-connected-corridors/>.
- [4] HighwaysEngland. *Signs of the future: new technology testbed on the A2 and M2 in Kent*. 2018, Highways England: <https://www.gov.uk/government/news>.
- [5] California CV Testbed. n.d. [cited 2020 February 26]; Available from: <http://www.caconnectedvehicletestbed.org/index.php/>.
- [6] INRIX. *INRIX IQ*. n.d. [cited 2020 27 December 2020]; Available from: <http://docs.inrix.com/traffic/speed/>.
- [7] HERE. *HERE Technologies*. n.d. [cited 2020 27 December 2020]; Available from: <https://www.here.com/platform/traffic-solutions>.
- [8] Saroj, A., S. Roy, A. Guin, M. Hunter, and R.M. Fujimoto, *Smart city real-time data-driven transportation simulation*, in *Proceedings of the 2018 Winter Simulation Conference*. 2018, IEEE Press: Gothenburg, Sweden. p. 857–868.
- [9] Hunter, M., R. Guensler, A. Guin, A. Saroj, and S. Roy, *Smart Cities Atlanta - North Avenue*, in *City of Atlanta Research Project*. 2019: <http://realtime.ce.gatech.edu/RenewAtlanta-GeorgiaTech-Final-Report.pdf>. p. 82.
- [10] Saroj, A., *Development of a real-time connected corridor data-driven digital twin and data imputation methods*, in *School of Civil and Environmental Engineering*. 2020, Georgia Institute of Technology: <https://smartech.gatech.edu/handle/1853/63642>.
- [11] Kostadinov, S. *How Recurrent Neural Networks work*. 2017 December 2017 [cited 2020 February 22]; Available from: <https://towardsdatascience.com/learn-how-recurrent-neural-networks-work-84e975feaf7>.
- [12] See, A. *Vanishing gradients and fancy RNNs*. Natural Language Processing with Deep Learning 2019 2019 [cited 2020 July 14]; Available from: <https://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture07-fancy-rnn.pdf>.
- [13] Hochreiter, S. and J. Schmidhuber, *Long Short-term Memory*. *Neural computation*, 1997. **9**: p. 1735-80. <https://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735>
- [14] Weber, N. *Why LSTMs stop your gradients from vanishing: A view from the backwards pass*. 2017 November 15 2017 [cited 2020 February 23]; Available from: <https://weberna.github.io/blog/2017/11/15/LSTM-Vanishing-Gradients.html>.
- [15] Yao, Y. and Z. Huang. *Bi-directional LSTM Recurrent Neural Network for Chinese word segmentation*. in *Neural Information Processing*. 2016. Cham: Springer International Publishing. https://link.springer.com/chapter/10.1007/978-3-319-46681-1_42