# Traffic incident duration prediction via a deep learning framework for text description encoding

Artur Grigorev [1,*], Adriana-Simona Mihăiţă [1], Khaled Saleh[1], Massimo Piccardi [1]

*Abstract*—**Predicting the traffic incident duration is a hard problem to solve due to the stochastic nature of incident occurrence in space and time, a lack of information at the beginning of a reported traffic disruption, and lack of advanced methods in transport engineering to derive insights from past accidents. This paper proposes a new fusion framework for predicting the incident duration from limited information by using an integration of machine learning with traffic flow/speed and incident description as features, encoded via several Deep Learning methods (ANN autoencoder and character-level LSTM-ANN sentiment classifier). The paper constructs a cross-disciplinary modelling approach in transport and data science. The approach improves the incident duration prediction accuracy over the top-performing ML models applied to baseline incident reports. Results show that our proposed method can improve the accuracy by** $60\%$ **when compared to standard linear or support vector regression models, and a further** $7\%$ **improvement with respect to the hybrid deep learning auto-encoded GBDT model which seems to outperform all other models. The application area is the city of San Francisco, rich in both traffic incident logs (Countrywide Traffic Accident Data set) and past historical traffic congestion information (5-minute precision measurements from Caltrans Performance Measurement System).**

*Index Terms*—**traffic incident prediction, deep learning, LSTM-ANN, sentiment classification**

## I. INTRODUCTION

When traffic accidents occur, the majority of traffic management centres (TMCs) store a brief textual description and the GPS coordinates of the incident. There is a lot of uncertainty at the beginning of disruptions with regards to how long the traffic incident will last, and most of the time, centres do not have an overview of the length or severity of the disruptions. Therefore, it is extremely insightful for TMCs to be able to utilise the data on historical traffic flows or readily available accident description to predict or improve predictions of the incident duration. In order to improve predictions, we need more information on the factors (both readily-available and historical) which can have an effect on the incident duration prediction accuracy. This paper presents an advanced incident duration prediction framework which makes use of additional incident report variables and past incidents records, merged into a hybrid machine learning (ML) modelling approach with deep learning encoding of additional features (e.g. textual incident description and historical traffic flow in the vicinity of the section). Feature encoding is justified since the traffic incident description and traffic flow/speed measurements have a high dimensionality,

which can lead to overfitting when using ML models and it may be worsened by the small size of a typical incident report data set.

This paper is organised as follows: Section I presents the challenges and reviews the related works; Section II introduces the data sources we have used as well as our traffic flow mapping algorithm for feature construction; Section III proposes our modelling framework and explains the ML models we have used, the LSTM sentiment encoder for textual incident descriptions, and the ANN encoder for traffic flow speed; Section IV introduces the results before summarising all findings in the Conclusions section.

### A. RELATED WORK

There are multiple research papers which use baseline incident reports from TMC with different machine learning models to predict the traffic incident duration [1]. The use of traffic flow and incident description features is found to be rare and mostly specific - topical text modelling [2] for the task of the incident duration detection, modelling or incident impact prediction by using traffic flows [3]. And its scarcity is highlighted since it requires the involvement of additional specific models with a feature fusion approach. In other words, traffic flow data is rarely combined with textual incident description and an actual incident reports since it requires a higher system complexity.

But feature combination can be observed in some specific research studies related to the traffic incident impact prediction, which rely heavily on the historical traffic flow data with and without consideration of features that are describing the incident [3]; other works have addressed a similar approach [4], [5], [6]. Also, these works don't focus on the incident duration prediction.

Sometimes, researchers try to apply uniform ML approaches or specific models for all the sub-tasks. Separate RBM models were applied to different kinds of features and feature fusion representing a uniform application of ML method to different data sets [7]. Also, kNN and Bayesian cost-sensitive networks were combined for the task of the incident duration prediction [8]. But neither of these research studies investigated a deep dive into their model selection.

Since we have the incident description and incident severity values in our incident reports, we can utilise specific models for the task of sentiment classification. Previously, the LSTM architecture has been compared with Support Vector Machines, Artificial Neural Networks, Deep Belief Networks and Latent Dirichlet Association on the task of detection of incidents from social media data [9]. LSTM

[1] University of Technology, Sydney
[*] Corresponding author: Artur.Grigorev@student.uts.edu.au (A. Grigorev)

was also successfully used for stock price prediction [10], making it applicable for modelling of traffic flow/speed time-series data. Despite its superior performance, we need to uplift and bring significant modifications to this architecture. Since we are planning to use encoded time series with machine learning methods, we need a controllable size of the feature vector to simultaneously avoid overfitting and provide enough information for ML methods. This is why we propose to use LSTM coupled with ANN, where the ANN feature vector size and the activation function are varied.

## II. CASE STUDY

In this study we assume that textual incident reports as well as historical traffic flows and speed data (including the ones from the moment when an incident happened) are readily available at the moment the incident was reported and sufficient to make the prediction of its duration.

### A. Incident description data set and baseline feature set

A Countrywide Traffic Accident Data set (CTADS) has been recently published [11]-[12], which contains about 1.5 million traffic accident records across 49 states of United States of America from February 2016 to December 2020 (version 4). Each incident report contains 47 features describing the traffic accident. The majority of these traffic accidents were recorded in the state of California. The most notable features include: a) Incident Severity (valued from 1 to 4), b) Start and End Time of the incident (from which the traffic incident duration is derivable), c) The road extent affected by the accident, d) textual Incident Description, d) weather and lighting conditions. For the extended description of features please refer to the original paper describing the data set [12]. This data set allows us to use the textual incident description and, hence, apply a sentiment analysis methodology (based on the incident severity) [13]. We further refer to these features as a baseline feature set, excluding the textual incident description.

### B. Traffic flow and speed data

To collect the data on traffic flows and speed we rely on the Caltrans Performance Measurement System (PeMS) [14], which provides aggregated 5-minute precision measurements of traffic movements across California. Although there is a lot more data for the Los Angeles area (which may be considered in our future research), we decided to concentrate on the area of the city of San Francisco. We focus on 83 Vehicle Detection Stations (VDS) placed in that area, and we try to manually associate each incident occurred in that area with a VDS in their 500m proximity. VDS in PeMS may have detector failures and incomplete readings, which is common across the data set and should be taken into account. Even though the PeMS data set contains data on reported incidents, we decided to use the descriptions from the Countrywise Traffic Accident Data Set since it provides a high-quality description of each incident (47 features in each incident report) extracted from Bing and MapQuest services.
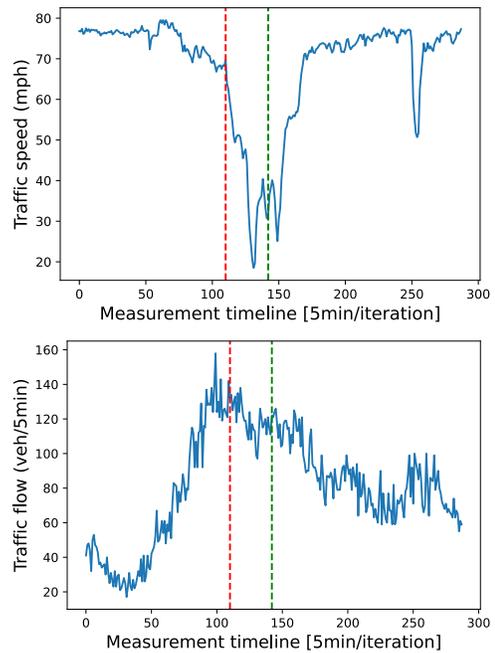


Fig. 1: a) Traffic speed and b) Traffic flow plots for the VDS associated to incident A-4798 (accident on US-101 Southbound with duration of 31 5-minute iterations - actual reported incident clearance time, without considering the incident recovery time). The red line denotes the start of the accident, and the green line the end of the accident. The blue line denotes the speed evolution in the vicinity of the incident location (drops almost to 20km/h) while the flow is still running at high values due to large numbers of vehicles blocked in traffic.

In total, from 9,275 incidents in the area (extracted from CTADS) we have obtained 1,932 traffic incident reports in a 500m proximity next to VDS stations, which we were able to associate with the correct (no detector faults) and complete traffic flows and speed readings. Incident to VDS association is necessary since both are represented as points and it is not clear which incident is related to which detector since incidents on different separate roads can be in proximity of one detector (also, since we have a different representation of street names in VDS and the incident data sets). The task of VDS-to-incident assignment can be a topic for additional research, but in this paper we summarize our extracted mapping strategy as follows. We extract the following speed and flow readings from each VDS station:

1. Speed – Traffic Speed from the 24h leading to the incident occurrence.

2. Flow – Traffic Flow from the 24h leading to the incident occurrence.

3. Speed7 – Traffic Speed on the same weekday, the week before the incident.

4. Flow7 – Traffic Flow on the same weekday, the week before the incident.

5. SD – the vector difference between the traffic speed on the day of the incident and on the same weekday, the week before the incident.

6. FD – the vector difference between the traffic flow on the day of the incident and on the same weekday, the week
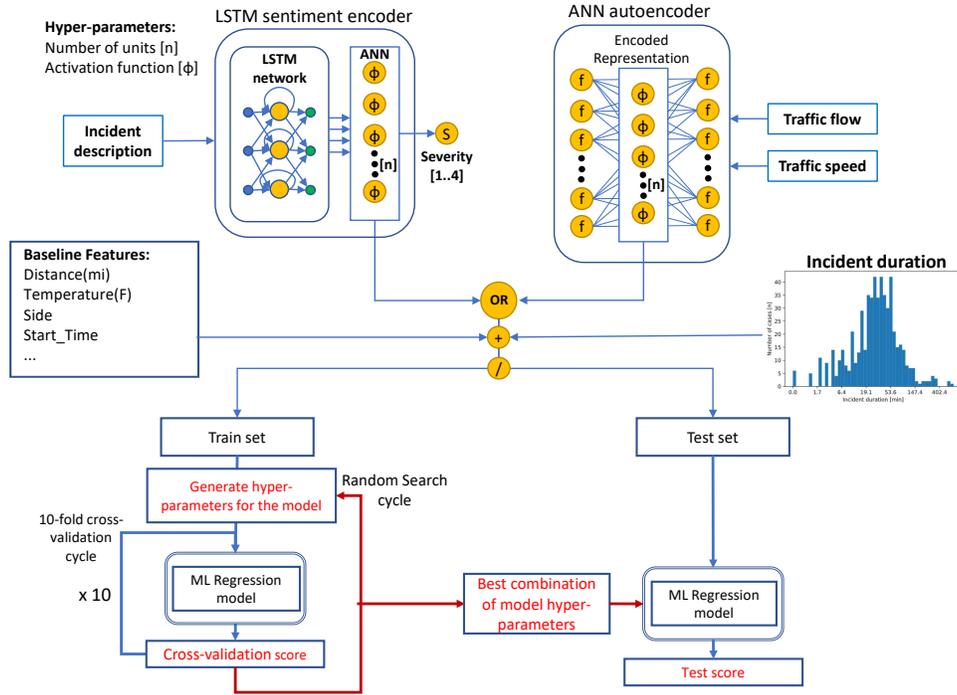
Fig. 2: The structure of the proposed framework

before the incident.

Each of these feature vectors contains 288 values, which correspond to 5-minute readings throughout the day. Since each of these vectors have a high dimensionality, we decide to perform dimensionality reduction via an ANN autoencoder.

The use of dimensionality reduction is justified since a large number of explanatory variables can cause model overfitting [15], [16].

**Regarding the 288 input values on the day of incident:** the traffic data is taken from the time between the incident start and minus 24h before of its occurrence and not during the entire day after the incident has been lodged.

Figure 1 shows an example of a traffic speed drop during the incident A-4798. After we have analysed different traffic flows and speed plots we expect that the traffic speed will be the most useful single feature for the task of incident duration prediction as the traffic flow measurement seems to be not affected by the accident (as the majority of vehicles will be waiting for the congestion to clear off the road, and will still be counted as part of the traffic flow). We will also use speed measurements from the weekday, 7 days before the incident in order to obtain the complete picture between what is a regular traffic flow condition versus disrupted traffic condition on the same time and same day of the week. We make the observation that we have also conducted a detailed feature ranking and selection (via SHAP values, forward feature selection, etc.) to several incident data sets which are not presented here due to space limitations.

The point A-4798 point was selected just as an example

for a traffic speed drop and its usefulness to the prediction problem; in reality, we have analysed about 100 traffic flow and speed plots before drawing the conclusions (we provide several shapshots of flow and speed reading in the supplementary material [17]). As an observation, by adding severity classification probabilities (from the LSTM-ANN model) to the feature vector for the task of incident duration prediction doesn't seem to be useful since we already included Severity, which is a strong feature.

## III. METHODOLOGY

Figure 2 shows how we use the data to perform the incident duration prediction. We combine the baseline feature set with either the encoded textual description or the encoded traffic flow/speed values. The encoder parts of both LSTM-ANN network and the ANN autoencoder have hyper-parameters in the form of number of units and used activation functions to ensure an optimal encoding for the specific ML method. After obtaining encoded representations associated with the incident, we search for the optimal hyper-parameters for each ML regression model at each case of the encoded representation. It allows us to adapt the model parameters to work with encoded data and provide the best cross-validation results.

### A. LSTM-ANN for the textual incident description encoding

Textual Incident Description in the CTADS data set describes type of disruption caused by the incident and/or location (Table I).

To perform the encoding of the textual description of the incident we use a combination of character-level LSTM and

Accident on I-280 Northbound at Exit 57 King St.
Right hand shoulder blocked due to accident on I-280 Northbound after Exits 54 54A 54B US-101.
Lane blocked due to accident on US-101 Presidio Pkwy Southbound at Exit 438 CA-1.
Accident on I-80 Westbound at Exits 1 1C / Bryant St / 8th St.
Second lane blocked due to accident on I-80 Eastbound at Exits 2B 2C Harrison St.
Lane blocked due to accident on US-101 Golden Gate Brg Southbound at Exit 439 Transit Transfer Facility.
Right hand shoulder blocked due to accident on I-280 Northbound at Exit 52 San Jose Ave.
Right hand shoulder blocked due to accident on US-101 Southbound at Exits 429B 429C Bay Shore Blvd.
Lane blocked on exit ramp due to accident on I-280 Northbound at Exit 55 Cesar Chavez.
Right hand shoulder blocked due to accident on I-280 Northbound at Ocean Ave.

TABLE I: Example of the Incident Description values

ANN for the sentiment analysis (Figure 3). We use the textual incident description from all the available traffic incident reports for the San Francisco area (9,275 incident records). Firstly, we set the target variable for the LSTM classification model as the incident severity (values 1 to 4).
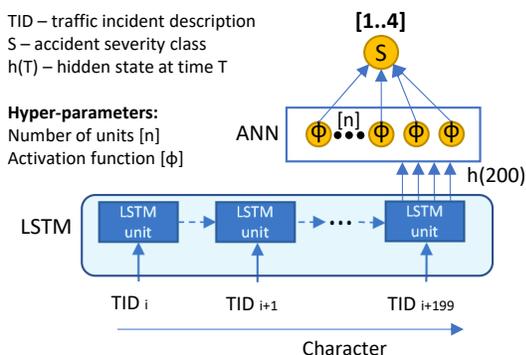


Fig. 3: LSTM sentiment encoder structure.

Secondly, we use the encoded representation of the textual description extracted from the LSTM sentiment classification model to use it as additional features for the task of incident duration prediction.

The incident description text is only provided at the beginning of the incident reporting timeline, and no temporal evolution is found across multiple countries for which we analysed the incident logs in our previous work [18].

Each textual description is formed into repeated strings up to 200 characters in length and each character in that string is then encoded by using binary encoding.

In order to showcase the importance of the textual incident description for the tasks of incident duration prediction and incident severity classification, we perform a word importance analysis using the LIME method (provided in the supplementary material [17]). We further train a an LSTM model with 80-units hidden state vector. We use the encoding of the incident description by using different numbers of neurons and different activation functions. An example of training results for one of the variants is shown on Figure 4. Traffic incidents descriptions were used to predict the incident severity. The data set was split into train, validation and test sets by proportion 70:20:10. Training results show that the LSTM sentiment encoder needs at least 15 epochs to converge, so we decided to train each variant of the LSTM sentiment encoder for 15 epochs. We use Mean Squared Error (MSE) as the loss function.
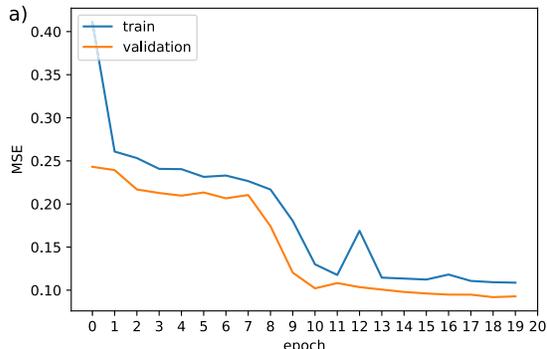


Fig. 4: Example of LSTM network training results using 12 units, a ReLU activation function, 10 epochs, 80 hidden units. a) Train-validation score over 20 epochs

*1) The use of MSE versus cross-entropy:* MSE is a legitimate metric for the classification when the target feature is represented as an ordered variable [19] in which MSE is preferred instead of the Cross-Entropy (CE) loss in order to reduce the model complexity and the probability of over-fitting. In our research we determined that CE required Nx5 sized matrix for the intermediate feature vector to the target value classification, while the MSE solution requires only Nx1 matrix, where N is size of the intermediate feature vector). MSE loss is also superior to CE loss for class-imbalanced datasets [20] and our incident severity feature distribution poses an imbalanced classification problem.

### B. Artificial Neural Network Encoder for the traffic flow/speed encoding

As additional data sources apart from the incident baseline features, we use the general structure of Artificial Neural Network (ANNs) Autoencoder [21] with varying number of neurons and different activation functions in the bottleneck layer to produce the encoded speed/flow data sets. Flow and speed values are normalised to the corresponding maximum observed traffic speed and flow in the data set. To improve the performance of the encoding model we use all the time series data available for each incident which could be matched to a VDS station. We combine normalised flow and speed data sets to perform the ANN model training which allows the model to grasp actual time series without focusing on speed and flow on an individual level. We do make the observation

that while speed and flow could be used as raw features in any ML prediction framework, the benefit of using ANN for auto-encoding is mainly a dimensionality reduction and improved accuracy in case of extreme outliers. Last, we extract the outputs of the ANN autoencoder bottleneck layer and use them as features in the ML models shown in Fig. 5.
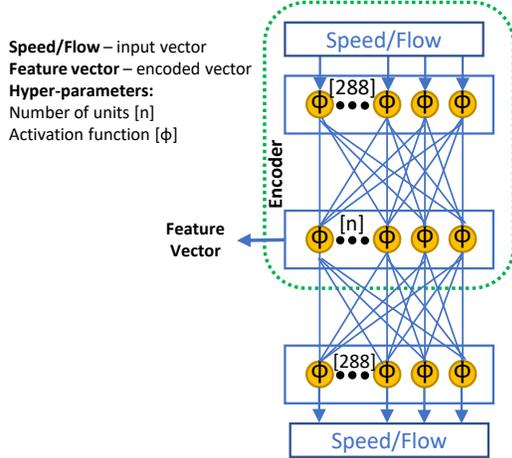


Fig. 5: The structure of the ANN autoencoder

The following activation units were used in the bottleneck layers of the ANN autoencoder and the LSTM sentiment encoder: a) the Rectified Linear Unit (ReLU) [22] which is a piecewise linear function (output values are $[0; +\infty]$ b) the Exponential Linear Unit (ELU) [23], which was developed to reduce bias shift (which leads to weight oscillations) c) the Tanh - a hyperbolic tan function which has the property of equalizing training over layers [24]; its output can take values in the interval $(-1; +1)$ d) the Sigmoid activation function which output can take values in the interval $(0; 1)$.

### C. Baseline Machine Learning model selection

When all encoding has been finalised, we first use the following ML regression models as a baseline to perform the incident duration prediction:

a) gradient boosting decision trees - GBDT [25] which rely on training a sequence of models, where each model is added consequently to reduce the residuals of prior models; b) extreme gradient decision trees - XGBoost [26] which rely on an exhaustive search of split values by enumerating over all the possible splits on all the features and contains a regularisation parameter in the objective function; c) random forests - RF [27] which applies a bootstrap-aggregation (bagging, which consists of training models on randomly selected subsets of data) and uses the average (or majority of votes) of multiple decision trees in order to reduce the sensitivity of a single tree model to noise in the data d) Support Vector Regression (SVR) machines [28] which are characterized by the use of kernels and symmetrical loss function (equal penalization of high and low errors), e) Decision Trees (DT) regression models [29] which rely on the repetitive process of splitting and generates a set of rules which can be used for the value prediction, f)

Linear Regression (for which we use standard Ordinary Least Squares optimisation) which represents the relation between features and the target variable as a linear equation targeting to minimize the residual sum of squares between the actual and the predicted values of the target variable.

*1) Model performance evaluation:* To evaluate the regression models on the task of the incident duration prediction we use the mean absolute percentage error and the root mean squared error defined as:

$$MAPE = \frac{1}{n} \sum_{t=1}^{n} \left| \frac{A_t - F_t}{A_t} \right| \qquad (1)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (A_t - F_t)^2} \qquad (2)$$

where $A_t$ are the actual values and $F_t$ - the predicted values, $n$ - the number of samples. We do make the observation that other performance metrics have been obtained (MAE, SMAPE), but given the current page limitations, we focus on MAPE, RMSE results only.

*2) Hyper-parameter tuning for the proposed regression model:* We use 10-fold cross-validation to overcome the over-fitting problem [30] and to assess the generalization performance of the ML models. In each scenario, the data set is partitioned into 10 folds. The ML regression model is trained on 9 folds to make prediction on the remaining fold. The procedure is then repeated 10 times and the accuracy results are averaged across several repetitions.

### D. MAPE versus RMSE comparison and their non-linear relationship

There is a non-linear relationship between MAPE and RMSE when performing regression, which can be verified by using different regression data sets. We tested this hypothesis on the Concrete Compressive Strength (CCS) Data Set from UCI Machine Learning Repository by using 1000 evaluations of random 9:1 train-test splits using Random Forest evaluated against MAPE and RMSE. Fig. 6a) presents the MAPE versus RMSE plot in which we observe that, the same MAPE result (e.g. 12%) may be attributed to multiple RMSE results (e.g. from 3.5 to 6.5). A similar situation observed for 45% of MAPE on CTADS using Random Forest (see Fig. 6)b). Therefore, the occurrence of a higher RMSE error when MAPE becomes lower (as in our paper) and vice-versa is a correct result. MAPE vs RMSE compared between 100-units random vectors with 1-10 value interval using 10,000 evaluations (see Fig. 6)c). As can be seen from all three sub-plots, the decrease in MAPE doesn't necessarily mean a decrease in RMSE. For our study we focused on discussing the MAPE metric, which is widely used in the literature on the topic of incident duration prediction since its intuitive meaning (e.g. a 30% MAPE means a 30% deviation of prediction from the actual incident duration) and a less inclination to high errors from outliers such as the case of RMSE. The results are part of the optimal Pareto Front [marked in orange] which showcases that our proposed method can obtain the set of optimal feature combination scenarios rather than only

one winning scenario. To conclude, despite an assumption on linear dependence between the RMSE and the MAPE metrics (assumption that both metrics should be reduced in an efficient solution), both in our incident duration case and the CCS data set, we observe a Pareto front of efficient solutions (no solution is sufficient in both metrics, making our results stand strong).

### E. Comparison to other baselines

It is hard to perform a comparison between different studies on the traffic incident duration prediction since different data sets are used for research purposes [1]. Majority of these data sets are also private and rely on different sets of features. CTADS data set appeared only recently (2019) and there is still no uniform convention on which data subset to use as a baseline, since the data set is big (1.5 million records) and heterogeneous (it includes reports from all kinds of traffic networks around United States). Indeed, in our previous work we have compared various ML-DL approaches against logs from Australia and USA, which can be used as extended results.

## IV. RESULTS

### A. Best model selection

First, we try to find the three best models which show high performance of the baseline feature set consisting of traffic accident reports for which we have available traffic flow counter data. We do so by performing a cross-validation as described in III-C.1 and a performance evaluation as detailed in III-C.2. Figure 7 shows the average MAPE score for the 10-fold cross-validation obtained across several ML models such as Random Forests (RF), GBDT, XGBoost, kNN, Decision Trees (DTs), Linear Regression (LR) and Support Vector Regression (SVR). Given that the majority of traffic incident duration prediction methods published previously have reported a MAPE score below 50% [1], we select RandomForest, GBDT and XGBoost as the best performing models as their MAPE score falls below 46%. Next, we evaluate these three models against the baseline feature set when we apply our novel modelling approach as previously explained in sections III-A-III-B: traffic flow, speed via ANN autoencoding and textual incident description via LSTM sentiment encoding.

There are in total 140 scenarios describing combinations of additional features [7 speed/flow/text features x 5 unit count x 4 activation functions] for each of the top three ML models. Given the restricted space allocation for this article, in Tables II-IV we present only the top 8 best scenario results ranked against MAPE for each ML model.

Findings reveal that the encoded textual description is among the top 3 configurations for every regression model as seen from Tables II-IV. Models also demonstrate a preference for the way of encoding: 1) the Tanh activation function forms a majority in the top results for GBDT both for encoding the incident description and flow/speed features (Table II), 2) the ReLU activation function forms a majority in the case of XGBoost (Table IV). This observation can

| AdditionData | units | activation | MAPE | RMSE |
|---|---|---|---|---|
| baseline | | | 44.99 | 58.4 |
| LSTM-sent | 12 | relu | 41.89 | 65.03 |
| Flow7 | 8 | tanh | 41.92 | 64.61 |
| LSTM-sent | 16 | tanh | 42.05 | 65.04 |
| Speed7 | 16 | tanh | 42.28 | 63.97 |
| LSTM-sent | 8 | tanh | 42.43 | 64.13 |
| LSTM-sent | 16 | relu | 42.56 | 65.82 |
| Flow | 16 | sigmoid | 42.57 | 64.53 |
| Speed7 | 2 | sigmoid | 42.59 | 64.76 |

TABLE II: Top 8 best scenario results for GBDT-enabled framework

| AdditionData | units | activation | MAPE | RMSE |
|---|---|---|---|---|
| baseline | | | 44.58 | 57.6 |
| Flow | 4 | tanh | 43.02 | 63.88 |
| LSTM-sent | 4 | elu | 43.02 | 65.04 |
| LSTM-sent | 12 | relu | 43.06 | 63.32 |
| Flow7 | 16 | sigmoid | 43.19 | 64.04 |
| Flow | 16 | elu | 43.30 | 63.92 |
| FD | 4 | elu | 43.32 | 64.12 |
| Flow7 | 16 | tanh | 43.33 | 63.47 |
| FD | 4 | sigmoid | 43.39 | 64.53 |

TABLE III: Top 8 best results for RF

point on a preference in the way of encoding features when using specific regression models. The best performing model among the top three finalists, when using all additional features seems to be GBDT: the best results are obtained when encoding the traffic incident description and when using the traffic flow 7 days before the incident with 12 units and the ReLU activation function [$MAPE = 41.89\%$, Table II] (therefore including the information on the regular traffic flow profile on the same weekday, together with the incident report proves important for the task of incident duration prediction).

Other models show a higher MAPE or RMSE results for the incident duration prediction (see RF enabled results in Table III with lowest $MAPE = 43.2\%$ for a combination of baseline, regular traffic flow, 4 layer units and a tanh activation function); similar findings appear for XGBoost-enabled results in Table IV with the lowest $MAPE = 43.44\%$, when using again the regular flow features, 8 layer units and ReLU activation function. This experiment shows that an accurate incident duration prediction immediately after the event has occurred is possible, leveraging the incident description and the measured traffic flow on the day of accident, which may prove very useful for TMCs to incorporate directly in their incident management platforms. Lower MAPE does not necessarily mean lower RMSE as seen from the baseline and additional data scenarios, but the LSTM sentiment encoding seems to be the approach that obtains the best RMSE score (64.13) when combined indeed with other variations of the activation function and number of hidden units (as shown in Table II).

### B. Parallel coordinates for scenario setup

To supplement the findings, we also provide a parallel categories representation of all the 140 scenarios for the GBDT
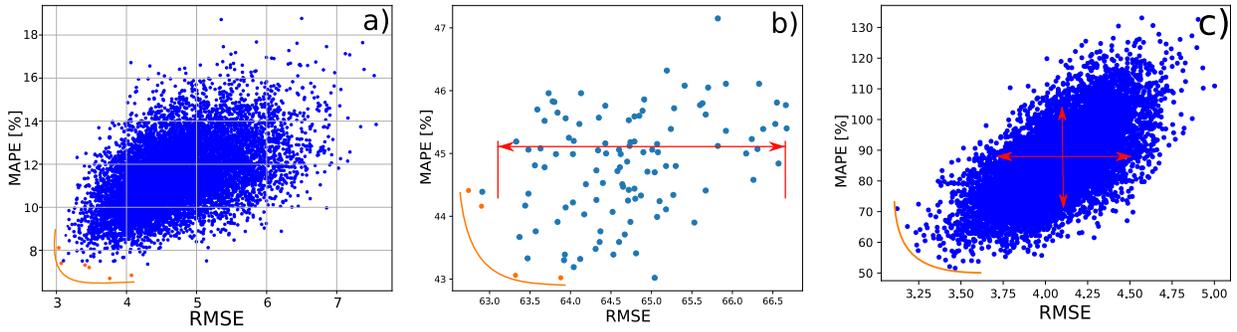
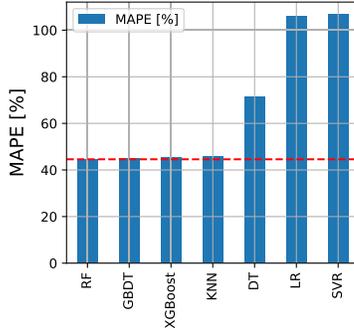Fig. 6: RMSE vs MAPE results for a) CCS data set, b) CTADS - incident duration c) Random vectors



Fig. 7: Regression results for baseline feature set across different ML models.

| AdditionData | units | activation | MAPE | RMSE |
|---|---|---|---|---|
| baseline | | | 45.44 | 63.41 |
| Flow | 8 | relu | 43.44 | 69.93 |
| LSTM-sent | 4 | tanh | 43.58 | 71.03 |
| Speed7 | 16 | tanh | 43.63 | 71.62 |
| SD | 4 | relu | 43.73 | 70.58 |
| Speed7 | 16 | relu | 43.80 | 71.92 |
| LSTM-sent | 16 | elu | 43.81 | 70.45 |
| LSTM-sent | 8 | relu | 43.82 | 72.19 |
| Flow7 | 2 | relu | 43.85 | 72.94 |

TABLE IV: Top 8 best results for XGBoost

model in Figure 8, which highlights the best combination of activation functions that seem to be working best alongside the character-level LSTM sentiment encoder of traffic flow incident textual description and speed information - mostly from previous daily speed profiling using historical data. The worst results seem to be the ones obtained when using only the speed or flow difference vector alongside the baseline incident features.

Encoding using Sigmoid and Tanh activation units on average performs best, probably because of the limited value range: Tanh an Sigmoid allow encoded representations to take values in ranges $[-1; +1]$ and $(0; 1)$ correspondingly, ReLU and ELU can take unlimited positive values. These results indicate which value ranges work best for encoded representation.

Comparison of MSE and CE implementations of lstm severity classification metric for the purpose of obtaining feature vector representation of Incident Description (see Fig. 8) shows that a sentiment classifier with Cross-entropy (lstmsentCE) as a target metric with one-hot encoded severity values is more efficient (left column attributed to lstmsentCE shows more blue rows associated with low metric values than lstmsentMSE - sentiment encoder which predicts severity as a single value). Comparison between the number of units shows preference for 4 units since the presence of the lowest error and absence of the highest error rows. Among the activation units, the Sigmoid is the best performer showing more low error results than other units. This scenario is to show how feature vector representing incident description may may be efficiently encoded to be used with conventional GBDT machine learning method: using cross-entropy for the severity classification, using 4 units and the Sigmoid as activation function.

## V. CONCLUSION

In this paper, we have proposed a novel framework to predict the incident duration using an integration of machine learning with traffic flow and description features encoded via several Deep Learning methods. This approach demonstrates the stable and noticeable improvement across all the performing models. The results give evidence to the importance of using specific deep-learning encoding approaches for all regression models which provide a further boost-up in the model performance from past historical traffic information and the textual incident description. Efficiently encoding incident-related features for the task of incident duration prediction is the first step to model the traffic incident impact on the traffic flow. Further work is currently being focused on exploring the spatial and the temporal dynamic prediction of the incident impact via graph-based modelling approaches. The research has the following limitations: a) we used as study area only San Francisco, but there is a data availability on traffic accidents and traffic flow for the area of California, b) traffic speed and flow were taken into account only before the incident; by collecting traffic count data for longer periods it possible to build traffic speed/flow profiles which may provide more accurate predictions. The societal impact of the research is as follows: the data availability of the predicted incident duration can improve for TMC incident and traffic management (e.g. TMC can announce when an incident is expected to dissipate, how many resources to allocate, etc), which in turn will reduce the time spent by people in the traffic congestion
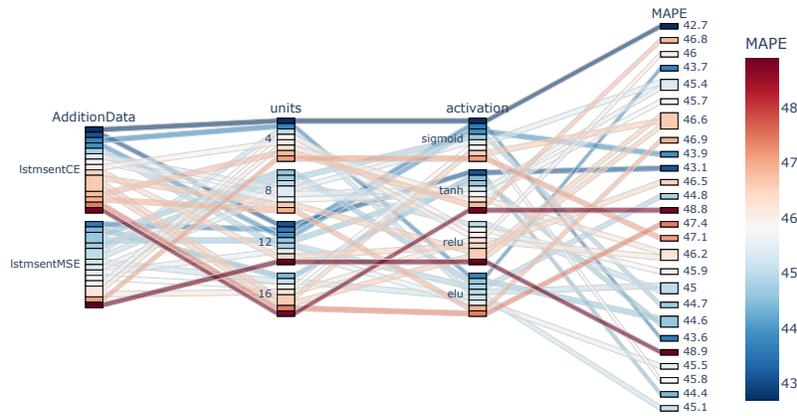
Fig. 8: Parallel categories representation for all regression scenarios with GDBT.

caused by the incident. The code for the paper can be found: https://github.com/Future-Mobility-Lab/TIDP_2022.

REFERENCES

[1] R. Li, F. C. Pereira, and M. E. Ben-Akiva, "Overview of traffic incident duration analysis and prediction," *European transport research review*, vol. 10, no. 2, pp. 1–13, 2018.

[2] Y. Gu, Z. S. Qian, and F. Chen, "From twitter to detector: Real-time traffic incident detection using social media data," *Transportation research part C: emerging technologies*, vol. 67, pp. 321–342, 2016.

[3] S. Fukuda, H. Uchida, H. Fujii, and T. Yamada, "Short-term prediction of traffic flow under incident conditions using graph convolutional recurrent neural network and traffic simulation," *IET Intelligent Transport Systems*, vol. 14, no. 8, pp. 936–946, 2020.

[4] T. Wen, A. S. Mihăiță, H. Nguyen, and C. Cai, "Integrated incident decision support using traffic simulation and data-driven models," *Transportation Research Board - 97th Annual Meeting, Washington, D.C.*, Oct. 2018.

[5] A.-S. Mihaita, H. Li, Z. He, and M.-A. Rizoiu, "Motorway traffic flow prediction using advanced deep learning," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, 2019, pp. 1683–1690.

[6] A.-S. Mihaita, Z. Papachatgis, and M.-A. Rizoiu, "Graph modelling approaches for motorway traffic flow prediction," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, 2020, pp. 1–8.

[7] L. Li, X. Sheng, B. Du, Y. Wang, and B. Ran, "A deep fusion model based on restricted boltzmann machines for traffic accident duration prediction," *Engineering Applications of Artificial Intelligence*, vol. 93, p. 103686, 2020.

[8] L. Kuang, H. Yan, Y. Zhu, S. Tu, and X. Fan, "Predicting duration of traffic accidents based on cost-sensitive bayesian network and weighted k-nearest neighbor," *Journal of Intelligent Transportation Systems*, vol. 23, no. 2, pp. 161–174, 2019.

[9] Z. Zhang, Q. He, J. Gao, and M. Ni, "A deep learning approach for detecting traffic accidents from social media data," *Transportation research part C: emerging technologies*, vol. 86, pp. 580–596, 2018.

[10] J. Sen and T. D. Chaudhuri, "Stock price prediction using machine learning and deep learning frameworks," in *Proceedings of the 6th International Conference on Business Analytics and Intelligence, Bangalore, India*, 2018, pp. 20–22.

[11] S. Moosavi, M. H. Samavatian, S. Parthasarathy, R. Teodorescu, and R. Ramnath, "Accident risk prediction based on heterogeneous sparse data: New dataset and insights," in *Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2019, pp. 33–42.

[12] S. Moosavi, M. H. Samavatian, S. Parthasarathy, and R. Ramnath, "A countrywide traffic accident dataset," *arXiv preprint arXiv:1906.05409*, 2019.

[13] S. Alkheder, M. Taamneh, and S. Taamneh, "Severity prediction of traffic accident using an artificial neural network," *Journal of Forecasting*, vol. 36, no. 1, pp. 100–108, 2017.

[14] C. Chen, K. Petty, A. Skabardonis, P. Varaiya, and Z. Jia, "Freeway performance measurement system: mining loop detector data," *Transportation Research Record*, vol. 1748, no. 1, pp. 96–102, 2001.

[15] Z. Sawalha and T. Sayed, "Transferability of accident prediction models," *Safety science*, vol. 44, no. 3, pp. 209–219, 2006.

[16] M. A. Sahraei, E. Kuşkapan, M. Çodur, and A. Tortum, *An Overview of Traffic Accident Prediction Models*, 02 2021.

[17] O. supplement, "Appendix: Traffic incident duration prediction via a deep learning framework for text description encoding," 2022, http://www.simonamihaita.com/papers/ITSC2022_TID_NLP_supplement.pdf.

[18] A. Grigorev, S. Lee, F. Chen, and A.-S. Mihaita, "Incident duration prediction using a bi-level machine learning framework with outlier removal and intra-extra joint optimisation," *Transportation Research Part C (major revisions as of March 2022)*, 2022.

[19] L. Gaudette and N. Japkowicz, "Evaluation methods for ordinal classification," in *Canadian conference on artificial intelligence*. Springer, 2009, pp. 207–210.

[20] S. Kato and K. Hotta, "Mse loss with outlying label for imbalanced classification," *arXiv preprint arXiv:2107.02393*, 2021.

[21] M. A. Kramer, "Nonlinear principal component analysis using autoassociative neural networks," *AIChE journal*, vol. 37, no. 2, pp. 233–243, 1991.

[22] A. F. Agarap, "Deep learning using rectified linear units (relu)," *arXiv preprint arXiv:1803.08375*, 2018.

[23] L. Trottier, P. Giguere, and B. Chaib-Draa, "Parametric exponential linear unit for deep convolutional neural networks," in *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2017, pp. 207–214.

[24] B. L. Kalman and S. C. Kwasny, "Why tanh: choosing a sigmoidal function," in *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*, vol. 4. IEEE, 1992, pp. 578–581.

[25] Y. Xia and J.-C. Chen, "Traffic flow forecasting method based on gradient boosting decision tree," in *Advances in Engineering Research*, 2017.

[26] T. Chen, T. He, M. Benesty, V. Khotilovich, Y. Tang, H. Cho *et al.*, "Xgboost: extreme gradient boosting," *R package version 0.4-2*, vol. 1, no. 4, pp. 1–4, 2015.

[27] Y. Liu and H. Wu, "Prediction of road traffic congestion based on random forest," in *2017 10th International Symposium on Computational Intelligence and Design (ISCID)*, vol. 2, 2017, pp. 361–364.

[28] H. Drucker, C. J. Burges, L. Kaufman, A. Smola, V. Vapnik *et al.*, "Support vector regression machines," *Advances in neural information processing systems*, vol. 9, pp. 155–161, 1997.

[29] L. Breiman, J. Friedman, R. Olshen, and C. Stone, "Cart," *Classification and Regression Trees; Wadsworth and Brooks/Cole: Monterey, CA, USA*, 1984.

[30] S. Geisser, "An introduction to predictive inference," 1993.

[31] B. R.-N. R. Baeza-Yates, "R. baeza-yates and b. ribeiro-neto: Modern information retrieval, addison wesley (1999)," vol. 17, no. 1, pp. 110–110, 2002.

[32] M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.

## APPENDIX

### A. Word importance for severity classification

To estimate word importance in the Incident Description feature, word count matrix has been transformed to a normalized TF-IDF representation (term frequency–inverse document frequency) [31]. N-gram value range is (1,2). Then linear dimensionality reduction has been performed using truncated singular value decomposition to 50 componenets for 7 iterations. Then we used GBDT classification model to fit incident severity and three quantiled groups (ratio 33%:33%:33% to represent equaly sized groups with duration intervals 0-29min, 30-71min and 72-2750min) of the incident duration. Classifer predictions were then analyzed for feature importance using LIME method [32], where every feature represents 1 word or 2 word combination presence in the incident description. One or more combinations of word in the description can contribute to the incident being classifed into one of severity groups (Fig. 9) - presence of "lanes blocked" and "two lanes blocked" has the highest contribution to the incident being classifed into highest (3) or lowest (0) severity group. Severity 1 or 2 is more related to the actual location, which represented as word describing Cesar Chavez St and I-280 Interstate Highway. High positive and opposite high negative contribution of words towards severity group observed for severity groups 1 and 2, where "280" and "chavez" have high opposite contributions, making this groups easily separable. When we perform classification towards equaly sized incident duration groups, "lanes blocked" has the highest positive contribution of the incident to be classified into low duration group. If accident happens on Cesar Chavez St, it can be easily classified into low duration group signifying importance of location for the task of incident duration prediction. High negative contribution of "lanes blocked" observed for duration group 1 with the highest contribution of "280" word meaning that incident appears on I-280 Interstate Highway.

### B. Traffic flow and traffic speed on the day of the incident

The following plots represent recorded traffic speed and flow on the day of the incident and week before in 500m proximity of the incident along the road (see Fig. 11 and 12). Reports in CTADS data set indicate that the highest impact of traffic incident is attributed to significant decrease in traffic speed, while traffic flow stays the least affected by disruption.
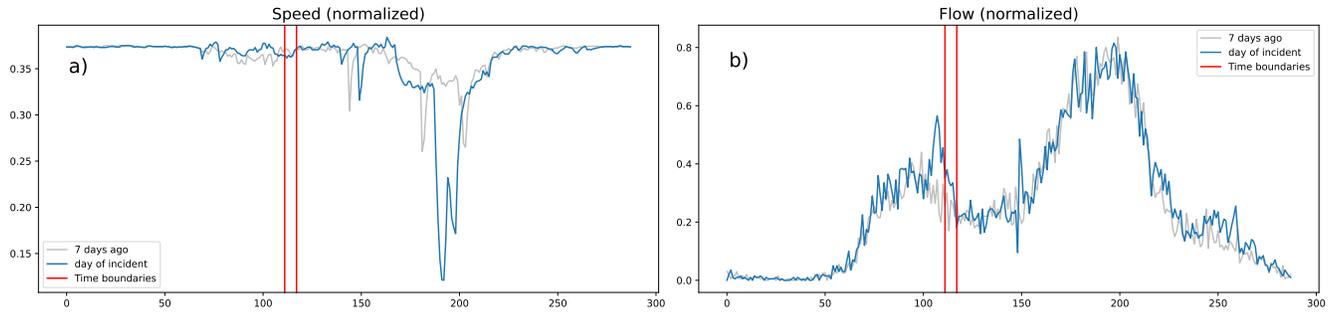
| Severity group=0 | | Severity group=1 | | Severity group=2 | | Severity group=3 | |
|---|---|---|---|---|---|---|---|
| Weight? | Feature | Weight? | Feature | Weight? | Feature | Weight? | Feature |
| +0.644 | lanes blocked | +1.559 | chavez | +2.805 | 280 | +0.982 | lanes blocked |
| +0.345 | two lanes | +1.190 | cesar | +1.909 | 280 northbound | +0.427 | two lanes |
| +0.231 | blocked due | +0.973 | <BIAS> | +0.828 | blocked | +0.365 | blocked due |
| +0.034 | due to | +0.894 | st | +0.740 | to accident | +0.174 | on i |
| +0.007 | to accident | +0.475 | i | +0.736 | accident | +0.008 | to accident |
| -0.407 | lanes | +0.467 | to | +0.721 | i 280 | -0.076 | cesar chavez |
| -0.620 | blocked | +0.465 | on | +0.697 | chavez st | -0.127 | lanes |
| -0.689 | i | +0.351 | at | +0.677 | accident on | -0.546 | blocked |
| -0.704 | <BIAS> | +0.309 | northbound | +0.448 | two lanes | -0.621 | i |
| -0.748 | to | +0.307 | cesar chavez | +0.375 | lanes blocked | -0.666 | on |
| -0.760 | st | +0.289 | two | +0.336 | at cesar | -0.672 | <BIAS> |
| -0.769 | accident | +0.153 | due | +0.194 | due to | -0.678 | to |
| -0.793 | two | -0.031 | northbound at | +0.187 | blocked due | -0.710 | at |
| -0.797 | due | -0.101 | blocked due | +0.138 | lanes | -0.762 | two |
| -0.800 | on | -0.125 | due to | +0.070 | northbound at | -0.773 | due |
| -0.818 | at | -0.310 | at cesar | +0.022 | on i | -0.918 | accident |
| -0.869 | northbound | -0.330 | two lanes | -0.160 | northbound | -0.953 | st |
| -0.924 | 280 | -0.372 | lanes | -0.208 | due | -0.961 | cesar |
| -0.979 | chavez | -0.466 | lanes blocked | -0.354 | cesar chavez | -0.997 | chavez |
| -1.149 | cesar | -0.647 | accident on | -0.358 | at | -1.116 | 280 |
| | | -0.684 | i 280 | -0.369 | two | -1.140 | northbound |
| | | -0.692 | chavez st | -0.498 | on | | |
| | | -0.711 | accident | -0.509 | to | | |
| | | -0.728 | to accident | -0.534 | i | | |
| | | -0.913 | blocked | -0.891 | st | | |
| | | -1.993 | 280 northbound | -0.994 | <BIAS> | | |
| | | -2.728 | 280 | -1.110 | cesar | | |
| | | | | -1.479 | chavez | | |

Fig. 9: Word importance estimation using LIME method for incident severity groups

| duration group=0 | | duration group=1 | | duration group=2 | |
|---|---|---|---|---|---|
| Weight? | Feature | Weight? | Feature | Weight? | Feature |
| +1.307 | lanes blocked | +0.548 | 280 | +0.389 | chavez st |
| +0.653 | two lanes | +0.444 | northbound | +0.256 | 280 northbound |
| +0.461 | blocked due | +0.357 | blocked | +0.149 | blocked due |
| +0.422 | lanes | +0.218 | chavez | +0.132 | northbound at |
| +0.326 | to accident | +0.214 | st | +0.092 | at cesar |
| +0.324 | on i | +0.213 | accident | +0.075 | cesar chavez |
| +0.255 | at cesar | +0.182 | cesar chavez | +0.068 | to accident |
| +0.230 | due to | +0.095 | two lanes | +0.062 | cesar |
| +0.216 | northbound at | +0.091 | cesar | +0.017 | to |
| +0.211 | chavez st | +0.050 | due to | -0.036 | <BIAS> |
| +0.177 | accident on | +0.039 | i 280 | -0.057 | lanes blocked |
| +0.026 | i 280 | +0.034 | lanes | -0.080 | due |
| -0.123 | st | +0.029 | 280 northbound | -0.088 | at |
| -0.153 | cesar chavez | -0.013 | on | -0.133 | two lanes |
| -0.232 | blocked | -0.030 | <BIAS> | -0.232 | accident |
| -0.232 | 280 northbound | -0.037 | two | -0.264 | chavez |
| -0.275 | i | -0.069 | to accident | -0.383 | st |
| -0.290 | at | -0.072 | northbound at | -0.502 | northbound |
| -0.348 | on | -0.077 | i | -0.580 | lanes |
| -0.405 | chavez | -0.129 | blocked due | -0.594 | 280 |
| -0.437 | northbound | -0.204 | chavez st | -0.633 | blocked |
| -0.439 | 280 | -0.655 | lanes blocked | | |
| -0.440 | to | | | | |
| -0.449 | due | | | | |
| -0.485 | accident | | | | |
| -0.544 | two | | | | |
| -0.724 | cesar | | | | |
| -0.918 | <BIAS> | | | | |

Fig. 10: Word importance estimation using LIME method for incident duration groups
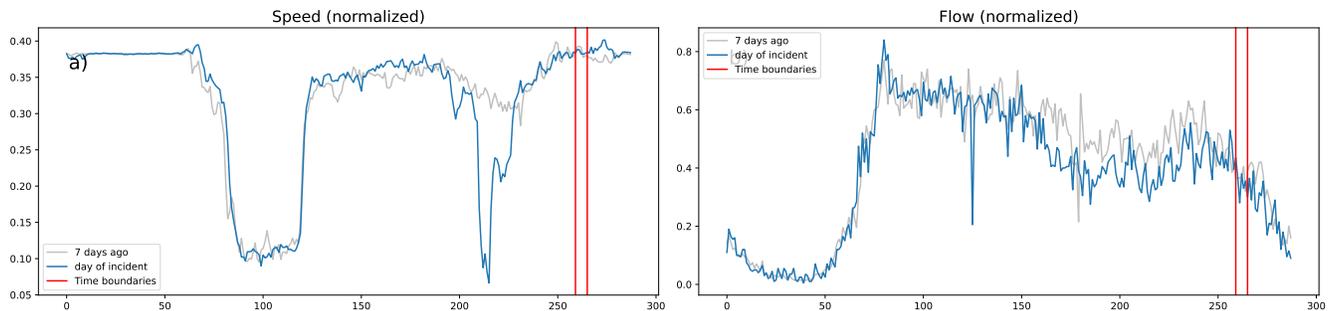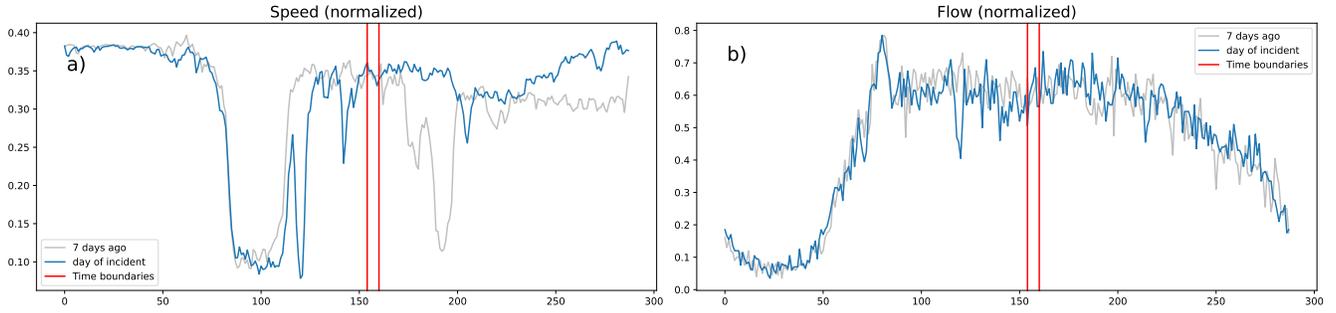
Fig. 11: Traffic speed and flow during the day of the incident. Part #1

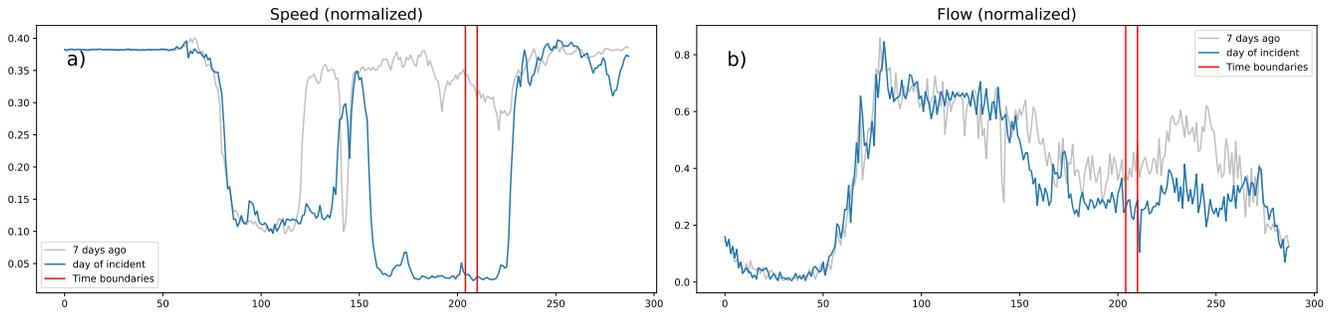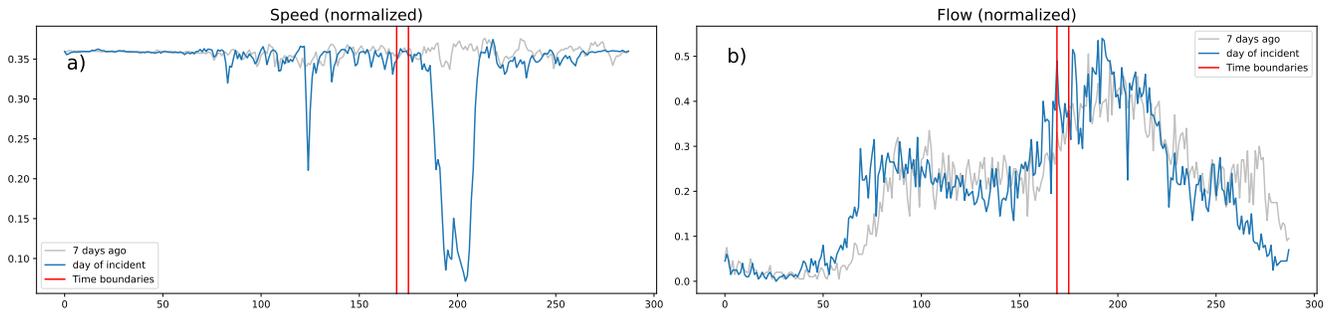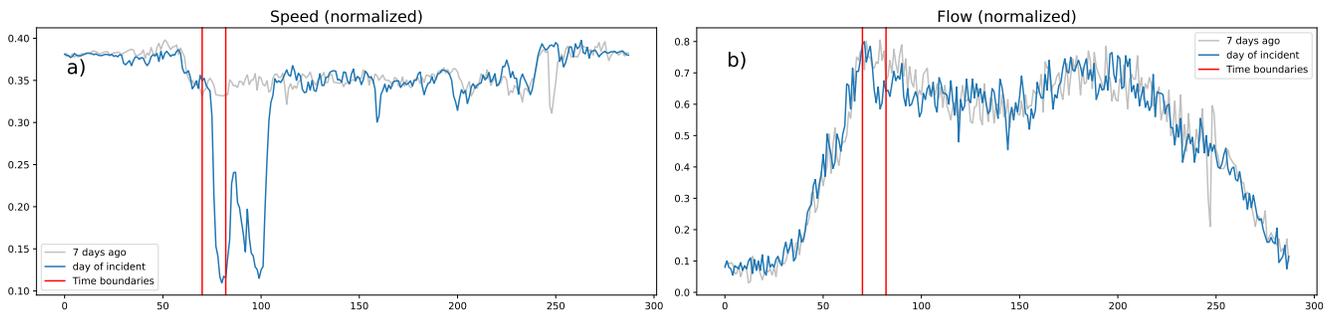Fig. 12: Traffic speed and flow during the day of the incident. Part #2