

1 **THE INTEGRATED ANALYSIS OF ACCIDENT REPORTS AND TRAFFIC FLOW**  
2 **DATA SETS WITH EARLY TRAFFIC DISRUPTION DETECTION AND**  
3 **SEGMENTATION**

4

5

6

7 **Artur Grigorev**

8 University of Technology Sydney, Ultimo, NSW 2007, Australia

9 Corresponding Author Artur.Grigorev@student.uts.edu.au

10

11 **Adriana-Simona Mihăiță**

12 University of Technology Sydney, Ultimo, NSW 2007, Australia

13 adriana-simona.mihaita@uts.edu.au

14

15 **Khaled Saleh**

16 University of Technology Sydney, Ultimo, NSW 2007, Australia

17 khaled.aboufarw@uts.edu.au

18

19 **Fang Chen**

20 University of Technology Sydney, Ultimo, NSW 2007, Australia

21 fang.chen@uts.edu.au

22

23

24 Word Count: 6581 words + 1 table(s) × 250 = 6831 words

25

26

27

28

29

30

31 Submission Date: 01 August 2022

**1 ABSTRACT**

2 Traffic accidents are often miss-reported with regards to either the exact location or the start and  
3 end time of the disruption due to several external factors - communication delays or misreporting,  
4 delays in accident clearance and non reporting of exact number of lanes, etc. Misleading informa-  
5 tion can lead to wrong decisions being made and can also affect the accuracy of any data-driven  
6 model that is responsible to predict either the severity or the disruption length. Several studies so  
7 far are using the reported incident data logs as truth ground which may wrongly affect any model  
8 build on top of this. To address these issues, our paper presents a novel machine learning frame-  
9 work that can be used for the early detection and prediction of incident durations. Firstly, we start  
10 by mapping and fusing several data sets related to reported traffic incidents (flows, speed, incident  
11 locations and incident details). Secondly, we propose several mathematical metrics (among which  
12 the Wasserstein and the Chebyshev metrics) are fed into an early detection and disruption segmen-  
13 tation algorithm which allows us to choose the best performing metric for future model training.  
14 Thirdly, we proposed a modelling interpretation of the traffic speed distributions based on the win-  
15 ning metric to correctly identify between single versus cascade incidents. Last we train and predict  
16 using various machine learning models and show that by using our enhanced modelling approach  
17 we can reduce the RMSE by almost 41% - 45.6% to as compared to the traditional case of using  
18 only historical incident logs.

19

20 *Keywords:* traffic accident, incident duration prediction, machine learning, traffic management

## 1 INTRODUCTION

2 The number of vehicles has been substantially increasing during the past decades, which currently  
3 leads to an increase in the number of traffic accidents (1). The National Highway Traffic Safety  
4 Administration (NHTSA) reported more than 5 million traffic accidents happening in the United  
5 States during year 2013 (2). Traffic Managements Agencies usually rely on Traffic Incident Man-  
6 agement Systems (TIMS) to collect data on traffic accidents, including information on various  
7 accident, traffic state and environmental conditions. Accurately predicting the total duration of  
8 an incident shortly after it is being finished, will help in improving the effectiveness of accident  
9 response by providing important information to decide the required resources to be allocated (re-  
10 sponse team size, equipment, traffic control measures) (3). Traffic accident is a rare event with  
11 stochastic nature. The effect of the accident can be observed as an anomalous state in the time  
12 series of traffic flow (4).

13 **Challenges:** The traffic accident analysis may be a challenging task due to incorrect or  
14 incomplete accident reports, including the set and the quality of the accident characteristics that  
15 have been reported. Accident reports can contain user-input errors related to the accident duration  
16 such as: 1) an approximate reporting of accident's start and end time 2) reporting of the accident  
17 start time could have been done after the incident finished in reality 3) a 'placeholder' accident  
18 duration reporting. In our previous research (5, 6) we found that timeline-related errors are present  
19 in accident reports across three different data sets from both Australia and the United States of  
20 America, which could be also the case with other data sets from around the world, due to multiple  
21 human and technical factors that can arise. To forecast the accident impact it is crucial to have  
22 an accurate and correct data regarding the observed disruption timeline. We emphasise that dis-  
23 ruptions observed in a recorded traffic state can be automatically segmented and associated with  
24 a reported accident at the same time and location as when the accident occurred, which allows  
25 to eliminate user-input errors from reports and improve the accident duration prediction perfor-  
26 mance in many traffic management centres around the world. To help address this issue, in our  
27 paper we propose various methods for a correct traffic disruption segmentation, the method for an  
28 association between vehicle detector stations and accident reports.

29 Another important challenge is that many incident data sets around the world are private and  
30 not shared for public investigation; for those open data sets, there are several missing information  
31 fields, or even worse, incomplete information regarding the traffic conditions in the vicinity of  
32 the accidents. Even often publish crash data sets are limited in size as well and contain a very  
33 small number of records. This represents a tight constraint when testing one framework over  
34 multiple countries with different traffic rules and regulations. For our studies we have oriented our  
35 attention towards two big open data sets - CTADS (Countrywise traffic accident data set) which  
36 contains 1.5 million accident reports and the Caltrans Performance Measurement System (PeMS)  
37 which provides data on traffic flow, traffic occupancy and traffic speed across California. Despite  
38 both being extensive data sets, vehicle detector station readings from PeMS are not associated  
39 with traffic accident reports from CTADS either by time, location or coverage area. The lack of  
40 such association makes it impossible to analyse the relation between accidents and their effects on  
41 traffic flow and speed. To address this challenge, in our paper we introduce the following mapping  
42 algorithm which will secure several steps such as :

- 43 • an association of Vehicle Detection Stations (VDS) with reported accidents in their prox-  
44 imity,
- 45 • a segmentation of traffic speed disruptions from detector readings,

1           • an association of detector stations with reported accidents (we will further show that this  
2           step is necessary due to many detected user-input errors in accident reports).

3 As a result, we obtain traffic disruptions segmented by the traffic flow associated with reported  
4 accidents. This association makes it possible to perform various important tasks of the accident  
5 analysis: 1) prediction of the traffic accident impact on the traffic speed based on accident reports,  
6 2) prediction of the traffic accident duration derived directly from the effect of disruption on the  
7 traffic speed (impact-based duration), 3) analysis of disruption propagation (each detected disrup-  
8 tion can be studied for spatial-temporal impact within the traffic network). Through this work,  
9 we will focus on the prediction of the impact-based accident duration and lay the foundation for  
10 a further research. Overall, the main contributions (summarised in Figure 1) of our paper is as  
11 follow:

12           1. We conduct a fusion methodology of two large data sets (CTADS and PeMS) for a  
13 detailed traffic accident analysis. To the best of our knowledge, this is the first research study  
14 proposing the methodology for merging of two large data sets of such nature, which allows an  
15 association between observed disruptions in traffic flow and the reported accidents.

16           2. We propose a novel methodology for the disruption mining using a combination of  
17 different metrics: a) the Wesserstein metric, which allows us to measure the disruption severity  
18 and b) the Chebyshev metric, which provides a higher selectivity for the disruption mining and  
19 a rectangular shape of the disrupted segments, allowing an automated disruption segmentation.  
20 We detail all unique properties of both metrics utilized together to allow an accurate disruption  
21 segmentation.

22           3. We perform the estimation of traffic accident disruption duration from traffic speed via  
23 the above metrics which allows us to alleviate user-input errors in accident reports.

24           4. We evaluate multiple machine learning models by comparing both the reported and the  
25 estimated accident duration predictions extracted from traffic speed disruptions.

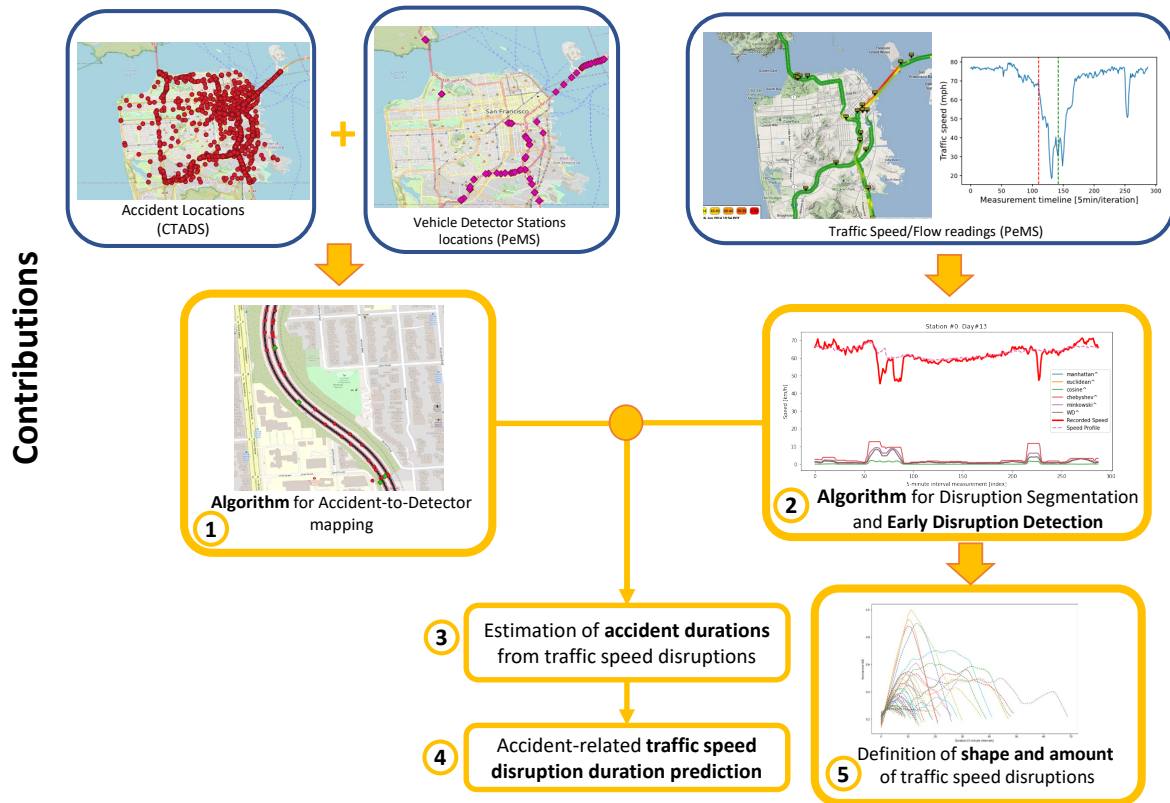
26           5. We introduce a new modelling approach which focuses on predicting the amount of  
27 the the disruption associated with an accident, rather than predicting via only historical logs. Our  
28 advanced disruption mining associated with traffic accident reports allows us to extract actual  
29 shapes of the disruptions, which allows a further analysis and modelling of the traffic flow impact  
30 properties. We also detect multiple accidents which produce secondary disruptions by using this  
31 approach.

32           Overall, this research forms the foundation for a new early traffic accident disruption de-  
33 tection, traffic flow disruption shape analysis and the use of observed traffic accident durations for  
34 correcting errors in user reports. Moreover, this work contributes to our ongoing objective to build  
35 a real-time platform for predicting traffic congestion and to evaluate the incident impact (see our  
36 previous works published in Mihaita et al. (6)-Shafiei et al. (7)-Mao et al. (8)).

37           The paper is further organised as follows: Section 3 discusses related works, Section 4  
38 presents the data sources available for this study, Section 5 showcases the methodology, Section 6  
39 presents the disruption segmentation results, showcases the result of data set fusion and Section 7  
40 provides conclusions and future perspectives.

## 41 **RELATED WORKS**

42 Multiple studies rely on user-input-based incident reports from Traffic Management Centers (TMC)  
43 with different machine learning models to predict the traffic incident duration (9). The use of traffic  
44 flow features is found to be rare and mostly specific - incident detection and incident impact pre-



**FIGURE 1 Contributions and data-flow schema for association of traffic speed readings with accident reports**

1 diction by using traffic flow (10). In other words, traffic flow data is rarely combined with actual  
 2 incident reports since it requires a higher system complexity and extensive data collection.

3 There were numerous studies related to accident detection from traffic flow using anomaly  
 4 detection techniques (11). Various methods used for anomaly detection in time series are applicable  
 5 for the task of traffic disruption detection. The ability to perform the detection of actual disruption,  
 6 which should give us actual shapes of disruptions and time intervals allows in-depth analysis of  
 7 usual accident statistics including the effect of the type of accident on the pattern of disruption  
 8 in traffic flow. By integrating data on traffic state with accident reports we are able to further  
 9 connect traffic flow disruption patterns to various accident characteristics (hour of the day, weather  
 10 conditions, crash type, type of vehicle involved - truck/car (12), the effect of road pavement types  
 11 (13), road design and road operation (14), etc).

12 Analysis of the effect of traffic incidents has been performed previously using Caltrans  
 13 PeMS data, where the measure of incident impact was represented as a cumulative travel time  
 14 delay (15), which is an aggregated value. However, traffic state recovery from disruptions is not  
 15 necessarily following a single pattern - it may be slowly dissipating, we may observe secondary  
 16 crashes, it may have a high or low impact, etc. Traffic accident duration prediction methodology  
 17 relies on reported traffic accidents, but actual reports may contain user-input errors and be mis-  
 18 aligned with the actual shape of disruption produced by the accident. Therefore, the approach  
 19 for disruption segmentation may provide the accident duration estimated from the actual shape of

1 disruption in traffic flow.

2 Various machine learning models are used to solve the task of traffic accident duration  
3 prediction(9) including k-nearest neighbours and Bayesian networks (16), Recursive Boltzman  
4 Machines and Support Vector Machines (17), Random Forests and XGBoost (6).

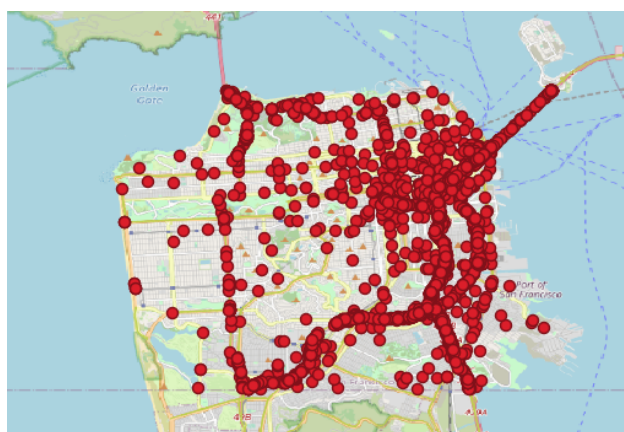
5 The definition of traffic incident duration phases is provided in the Highway Capacity Man-  
6 ual (18) and includes the following time-intervals: 1) incident detection - the time interval between  
7 the incident occurrence and its reporting, 2) incident response - time between the incident reporting  
8 and the arrival of the response team, 3) incident clearance time between the arrival of the response  
9 team and the clearance of the incident, 4) incident recovery - the time between the clearance of  
10 the incident and the return of traffic state to normal conditions. In this research, we rely on total  
11 incident duration - the time between incident occurrence and return of the state to normal condi-  
12 tions. Also, we analyse the subset of traffic incidents - traffic accidents. As we found during the  
13 data investigation, traffic accident duration is reported at the time when the incident is cleared by  
14 the response team, which doesn't include the duration of the effect that the accident produces on  
15 traffic flow.

## 16 CASE STUDY

17 Before diving into the methodology, we provide a brief introduction into the data sets in use for  
18 showcasing our approach, which helps establishing the modelling base and understanding of the  
19 steps taken. We make the observation that the current methodology can be applied on any incident  
20 and traffic state data set which can contain a time component, and is not bounded to the chosen  
21 data sets for exemplification.

### 22 CTADS: Accident reports data set

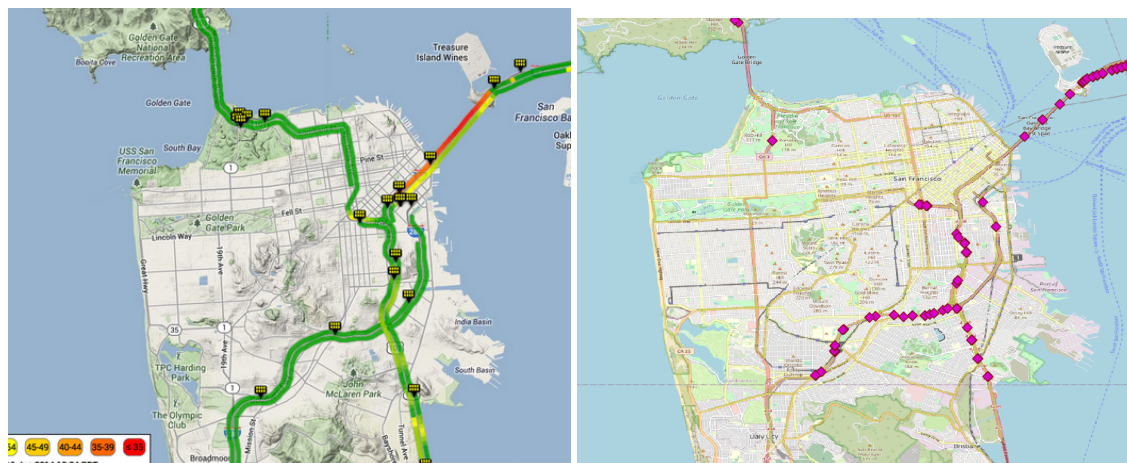
23 We rely on accident reports from the "Countrywide Traffic Accident Dataset" (CTADS), recently  
24 released in 2021 (19, 20), which contains 1.5 million accident reports collected for almost 4.5 years  
25 since March 2016, each report containing 49 features obtained from MapQuest and Bing services.  
26 We select the area of San-Francisco, U.S.A and extract data for 9,275 accidents (see Figure 2).



**FIGURE 2 CTADS reported accidents for San-Francisco**

## 1 PeMS: Traffic speed and flow data set

2 We rely on Caltrans Performance Measurement System (PeMS) (21) to collect data on traffic flow  
 3 and speed. This data set provides aggregated 5-minute measurements of traffic flow, speed and  
 4 occupancy across California. We decided to extract the data for the area of San-Francisco (see 3a),  
 5 which contains 83 Vehicle Detection Stations (VDS) placed in that area (see 3b), and we try to  
 6 associate each traffic accident occurred with each of San-Francisco VDS in their 500m proximity  
 7 using the algorithm detailed in the following section. In total, from 9,275 accidents in the area  
 8 (extracted from CTADS) we have obtained 1,932 traffic incident reports which we were able to  
 9 associate with the correct and complete traffic flow and speed readings from a VDS.



**FIGURE 3 1) PeMS data set area coverage for San-Francisco 2) Mapping of the Vehicle Detection Stations from PeMS data set.**

## 10 METHODOLOGY

11 The new framework we propose in this paper is represented in Figure 1 which we support across  
 12 some initial definitions for our modelling approach (see next sub-section). First, we associate the  
 13 road segments with their corresponding Vehicle Detector Stations (VDS) from the Caltrans PeMS  
 14 data set, as well with the locations of reported accidents (see Algorithms 1 and 2 proposed in sub-  
 15 section 5.2). The main outcome of this algorithm is that traffic accidents will get associated with  
 16 the traffic flow, speed and occupancy readings from the VDS stations.

17 Second, we propose an new algorithm for an early disruption detection and segmentation,  
 18 detailed in sub-section 5.3. By segmenting the disruptions occurred in time-space proximity of  
 19 reported traffic accidents, we obtain the estimated traffic accident duration. This gives us much  
 20 more information to include in the model training than just the simple accident duration: 1) the  
 21 disruption shape in terms of modifications of speed data profiles from the standard patterns 2) the  
 22 accident duration estimated from the impact on the traffic flow 3) the cumulative accident impact  
 23 estimation.

## 24 Definitions

### 25 *Speed difference estimation definitions*

26 In the current study we compare the performance of multiple difference metrics that will help us  
 27 to correctly estimate the impact of an accident and the deviation from the historical speed patterns.

1 These metrics are defined as follows:

2 a) the Chebyshev difference is expressed as:

$$3 \quad D_{\text{Cheb}}(x, y) := \int \max_i (|x_i - y_i|) \quad (1)$$

4 b) the Wasserstein difference:

$$6 \quad D_{\text{WD}}(u, v) = \inf_{\pi \in \Gamma(u, v)} \int_{\mathbb{R} \times \mathbb{R}} |x - y| d\pi(x, y) \quad (2)$$

7 where  $\Gamma(u, v)$  is the set of (probability) distributions on  $\mathbb{R} \times \mathbb{R}$  whose marginals are  $u$  and  $v$   
 8 on the first and second factors respectively.

9 c) the Cosine difference:

$$11 \quad D_{\text{C}}(u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2}. \quad (3)$$

12 where  $u \cdot v$  is the dot product of  $u$  and  $v$ .

13 d) the Euclidean difference between 1-D arrays  $u$  and  $v$ , is defined as

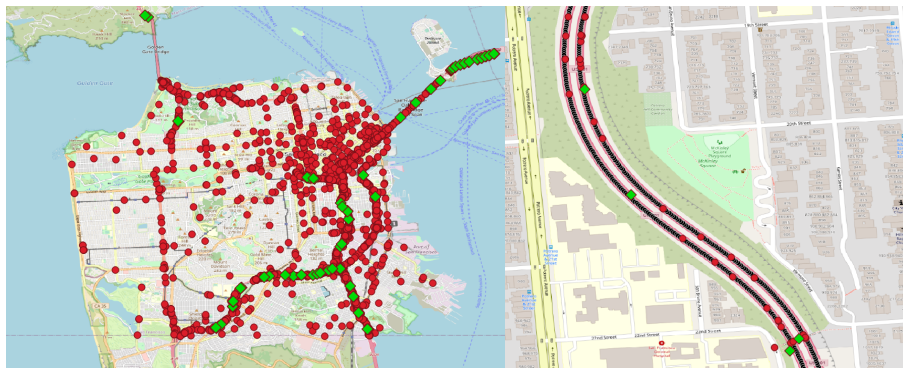
$$15 \quad D_{\text{E}}(u, v) = \left( \sum (w_i |u_i - v_i|^2) \right)^{1/2} \quad (4)$$

16 e) the Minkowski difference between 1-D arrays  $u$  and  $v$ , is defined as

$$18 \quad D_{\text{M}}(u, v) = \left( \sum |u_i - v_i|^p \right)^{1/p} \cdot \left( \sum w_i (|u_i - v_i|^p) \right)^{1/p}. \quad (5)$$

## 20 5.1 Algorithm for vehicle detector station to accident association

21 In order to match correctly what traffic conditions reflect best the effects of each incident, we  
 22 further define the association procedure between traffic accidents and VDS (A-to-VDS), for the  
 23 San Francisco area, as shown in Figure 4. We observe that only a few traffic accidents have VDS  
 24 stations in their proximity to allow a good traffic speed and flow extraction.



**FIGURE 4 a) Traffic incidents [marked in red dots] and VDS stations [marked in green diamonds] across the San Francisco city area. b) A closer look into VDS and traffic incidents locations where black dots are road points.**

25 In order to find the traffic incidents for which we can have traffic flow and speed data, we  
 26 develop a mapping algorithm (A-to-VDS) which consists of two parts (see Algorithm 1-2), defined  
 27 by the following steps:

- 28 1. We first extract the primary and secondary road lines from Open Street Maps.
- 29 2. Road segments are then transformed (segmented) into points at 2-meters equal distance.
- 30 3. Each VDS station and accident are mapped to the closest road point (up to 10m distance).
- 31 Accidents that are not in the proximity of road points removed.



- 1           4. From this step we use Algorithm 1 to process the point-based representation of VDS,  
 2           accidents and road segments. The *vdsPoints* array contains a tuple of the form (VDS ID,  
 3           x and y coordinates); each point in *accidentPoints* contains an array *visitedBy* (initialized  
 4           as empty) to maintain a list of stations in the proximity of the accident; *assignedVDS* is  
 5           a final nearest VDS station close to the accident along the road.  
 6           The algorithm relies on a recursive function to implement the process of visiting road points  
 7 (see Algorithm 2). The association part of the algorithm works as follows:  
 8           1. We select the current VDS station.  
 9           2. We move (jump by points) in all possible directions available from the starting and forth-  
 10           coming points in a 3m radius. This radius allows us to move along the road jumping  
 11           between road points. Movement in all possible directions allows to grasp the propa-  
 12           gation of the traffic congestion associated with the accident. The maximum available  
 13           distance is set to 500m (250 jumps) and allows to limit the observable impact distance.  
 14           3. By moving across points we collect traffic incidents in the 5m proximity of each point  
 15           and associate them with the current VDS station.

---

**Algorithm 1:** A-to-VDS: Accident to VDS mapping algorithm

---

**Input:** *point*  
**Output:** *None*  
**Access global arrays:** *roadPoints, accidentPoints, vdsPoints*  
**Function** *visitNearestPoints(VDSID, point, currentHops)*  
*accidents := findNearestAccidents(point, accidentPoints, 10m)*  
**for** *i = 0 to length(accidents)* **do**  
  | *a := accidents[i]*  
  | *a.visitedBy.append([VDSID, currenthops])*  
**end**  
**if** *currenthops < 500/2* **then**  
  | *roadpoints:=findNearestRoadPoints(point, roadPoints, 3m)* **for** *i = 0 to*  
  | *length(roadpoints)* **do**  
  | *rp := roadpoints[i]*  
  | **if** *VDSID not in rp.visitedBy* **then**  
  | | *rp.visitedBy.append(VDSID)* *visitNearestPoints(point, currentHops + 1)*  
  | **end**  
  | **end**  
**else**  
  | **Return**  
**end**

---

16           The algorithm is recursive and relies on the list of visited points for each VDS. At the end  
 17 of the algorithm, we have a subset of traffic accidents with their associated VDS which allows us  
 18 to extract the traffic flow and speed in the vicinity of the accident. Ideally, all traffic accidents  
 19 should have associated traffic flow but given their unavailability (due to detector coverage), we  
 20 select accident reports which have associated traffic flow information currently available from the  
 21 PeMS data set.

**Algorithm 2:** The recursive function for traveling across road points

---

```

Input: roadPoints, accidentPoints, vdsPoints
Output: assignedAccidents
for  $i := 0$  to  $length(vdsPoints)$  do
  |  $vds := vdsPoints[i]$ 
  |  $visitNearestPoints(vds, 0)$ 
end
assignedAccidents = []
for  $i := 0$  to  $length(accidentPoints)$  do
  |  $accident := accidentPoints[i]$ 
  | if  $length(accident.visitedBy) > 0$  then
    |  $accident.assignedVDS = sort(accident.visitedBy, sortvalue = hops)[0]$ 
    |  $assignedAccidents.append(accident)$ 
  | end
end
return assignedAccidents

```

---

1 **5.2 Algorithm for automated disruption segmentation (ADS)**

2 Once the accidents have been mapped and associated to their VDS stations which allows us to  
 3 select the flow/speed that match the day of the incident, etc, we are using the extracted traffic state  
 4 parameters to propose a new automated disruption segmentation (ADS) method. The algorithm  
 5 for the segmentation of disruptions via traffic speed works as follows:

- 6 1. A time series pre-processing step prepares all the data for segmentation (see Alg. 3):  
 7 (a) Calculate the average monthly profile for daily traffic speed measurements;  
 8 (b) Iterate over the traffic speed time series using a moving window of 1-hour time  
 9 interval (in total there are twelve measurements of 5-minute each)  
 10 (c) On each iteration perform a comparison of a 12-unit window between the monthly  
 11 profile and the current day of measurements. The resulting single value is added to  
 12 the resulting time series sequence.  
 13 (d) Calculated the time series differences (TS) choosing the above defined metrics will  
 14 be then adjusted by selectivity (using the power function, which will keep values  
 15 closer to one for the least affected by the function and minor values the most sup-  
 16 pressed) and normalized to produce nTS and pTS arrays respectively.
- 17 2. The time-series segmentation step (see Alg. 4):  
 18 (a) A first order derivative (dTTS) is calculated for the resulting time series of the pre-  
 19 vious stage (nTS), which returns positive peaks when entering the disruption and  
 20 negative peaks when exiting the disruption state.  
 21 (b) We iteration over resulting derivative time series to record the opening and closing  
 22 of each disruption in each time series. If two consecutive positive peaks (opening  
 23 times) are observed then we choose the largest one between the two (we will further  
 24 debate on this aspect in our future work plans). We repeat the same for consecutive  
 25 negative peaks.  
 26 (c) We then associate the detected disruptions with the accident reports: for each acci-  
 27 dent report, we extract the traffic speed time series on the day of the accident and

1 if both opening and closing times are recorded, we perform an association of the  
 2 accident with these times and extract the actual time series sequence for further  
 3 analysis.

4 **Enhancing selectivity:** We use the convolution with the kernel (1,1,1), which attributes to  
 5 the morphological dilation operation, to facilitate the work of the segmentation algorithm. By ap-  
 6 plying this convolution we make multiple consequent differences to be accumulated ; for example,  
 7 assuming we have a sequence of 0.3, 0.1, 0.1, 0.2 and 0.2 as differences for each 5-minute step,  
 8 therefore a total of 0.9 change over 4 iterations. The convolution (1, 1, 1) will produce the values  
 9 of 0.5, 0.4, and 0.5 by making a sequence of high values from the sequence of small changes (see  
 10 Figure 5). The dilation operation is primarily used in computer vision tasks to make connected  
 11 groups from closely placed scattered points to facilitate a further image analysis.



**FIGURE 5 The application of dilation operation to an image and time series**

12 To obtain the monthly profile, the traffic speed measurement sequence was obtained for  
 13 a duration of 1 month from the VDS before the accident occurred, and was done separately for  
 14 each accident. This sequence then gets reshaped into a matrix of the form  $[number\_of\_days; 288]$ ,  
 15 where columns contain the total number of measurements across an entire day ( $24 \cdot 12 = 288$ ). The  
 16 monthly average was then calculated across axis 1 to obtain a vector with 288 values of measure-  
 17 ments. This vector gets recalculated for a number of days of observations from each detector to be  
 18 comparable with the VDS daily measurements.

19 As an observation, the constants  $pThreshold$  and  $nThreshold$  represent thresholds for change  
 20 that observed in the time series of the metric derivative; they allow us to define a positive and nega-  
 21 tive change of the difference metric, the selectivity controls power function coefficient to suppress  
 22 the non-significant and filter the most significant disruptions.

### 23 **5.3 Modification of the algorithm for automated real-time early disruption detection**

24 Since our proposed algorithm doesn't look into the future and calculates different metrics based  
 25 on the currently observed traffic speed and a few measurements in the past (11 units in the current  
 26 study), we can perform an early accident detection which will consist in calculating and comparing  
 27 the FOD of Chebyshev metric based on the monthly profile. The detection of significant positive  
 28 peaks (e.g. 0.3-0.5 of normalized difference metric) can identify the amount of disruption in real-  
 29 time. The end of the disruption can be detected using the same approach in real-time as well by  
 30 observing a significant negative peak.

**Algorithm 3:** Algorithm for automated disruption segmentation. Part 1

---

```

Input: monthlyProfile, speedReadings
Output: cTS
; //Accident array contains a day number, starting and ending index
  for segmented traffic disruptions
step := 1
windowSize := 12
i := windowSize
lastDiff = 0
DS = []
while i < length(speedReadings) do
  | A := speedReadings[i - windowSize + 1 : i]
  | B := monthlyProfile[i - windowSize + 1 : i]
  | diff := metric(A, B)
  | DS.append(diff)
  | lastDiff = diff
end
for i = 0 to windowSize do
  | DS.append(lastDiff)
end
pTS = power(TS, selectivity)
nTS = normalize(pTS)
dTS = derivative(nTS)
cTS = convolution(dTS, [1, 1, 1])
return cTS

```

---

**1 5.4 Accident duration prediction definitions**

2 Using all available data sets and the incident information, we first denote the matrix of traffic  
3 incident features as:

$$4 \quad X = [x_{ij}]_{i=1..N_i}^{j=1..N_f} \quad (6)$$

6 where  $N_i$  is the total number of traffic incident records used in our modelling and  $N_f$  is the total  
7 number of features characterising the incident (accident severity, vehicles involved, number of  
8 lanes, etc) according to the accident report data set.

9 Traffic speed is represented as a vector with 5-minute averaged readings from Vehicle De-  
10 tector Stations:

$$11 \quad S = [s_i]_{i=1..N} \quad (7)$$

13 where  $N$  is the total amount of traffic speed readings.

14 As detailed in our methodological framework, the last step of our approach for improving  
15 the incident duration prediction is to assess the performance of various Machine Learning models  
16 on the tasks of predicting the reported and the estimated accident duration. This is mainly to  
17 understand if the approach we have proposed for an early accident detection is better or worse than  
18 by simply using reported incident logs from the traffic management centres. We define the task of  
19 the accident duration prediction as a regression problem.

**Algorithm 4:** Algorithm for automated disruption segmentation. Part 2

---

```

Input:  $cTS, pThreshold, nThreshold, selectivity$ 
Output: Accidents
; //Accidents array contains a day number, starting and ending index
  for segmented traffic disruptions
state := 0
Accidents = []
for  $i := 0$  to  $length(cTS)$  do
  | if  $cTS[i] > pThreshold$  then
  | | if  $state \langle \rangle +1$  then
  | | |  $state = +1$ 
  | | |  $enteridx = i$ 
  | | else
  | | | if  $cTS[i] > cTS[enteridx]$  then  $enteridx = i$ ;
  | | end
  | end
  | if  $cTS[i] < nThreshold$  then
  | | if  $state \langle \rangle -1$  then
  | | |  $state = -1$   $exitidx = i$ 
  | | else
  | | | if  $cTS[i] < cTS[enteridx]$  then  $exitidx = i$ ;
  | | end
  | end
  | if  $i \bmod 288 == 0$  and  $i > 0$  then
  | |  $state = 0$ 
  | |  $Accident.append([i \div 288, enteridx, exitidx])$ 
  | end
end
return Accidents

```

---

1 The incident duration regression vector ( $Y_r$ ) is represented as:

$$2 \quad Y_r = [y_i^r]_{i \in 1..N}, y_i^r \in \mathbb{N} \quad (8)$$

3 and the regression task is to predict the traffic accident duration  $y_i^r$  based on the traffic incident  
4 features  $x_{i,j}$ . The regression models go via an 10-fold cross-validation procedure with intense  
5 hyper-parameter tuning. Several models have been trained and compared and their performance  
6 assessed in our results section.

7 The estimate the accident duration prediction performance we use the root mean squared  
8 error:

$$9 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (A_i - F_i)^2} \quad (9)$$

10 where  $A_i$  represents the actual value of the incident duration logs (which can be either used raw  
11 from the incident logs or from our automated segmentation technique) and  $F_i$  the predicted value  
12 of the regression models.  
13  
14

## 1 RESULTS

### 2 Data exploration and setup

3 CTADS data set contains traffic accident reports, which after an initial data mining investigation,  
4 we found to contain several user-input errors; for example, a lot of traffic accident durations have  
5 been rounded to 30 or 360 minutes (see Fig. 6d)); or the incident start time which was reported  
6 is unrelated to any disruptions observed by the vehicle detector stations in the proximity - see  
7 Figure 6 in which we have provided two different examples of speed recorded during two different  
8 accidents A-1390 and A-7102; the red lines indicate the official reported start and end time of the  
9 accidents, while in reality the accidents have had a long lag in spreading across the network - see  
10 Fig. 6a) or were reported much later that the official speed drop was recorded - see Fig. 6b).

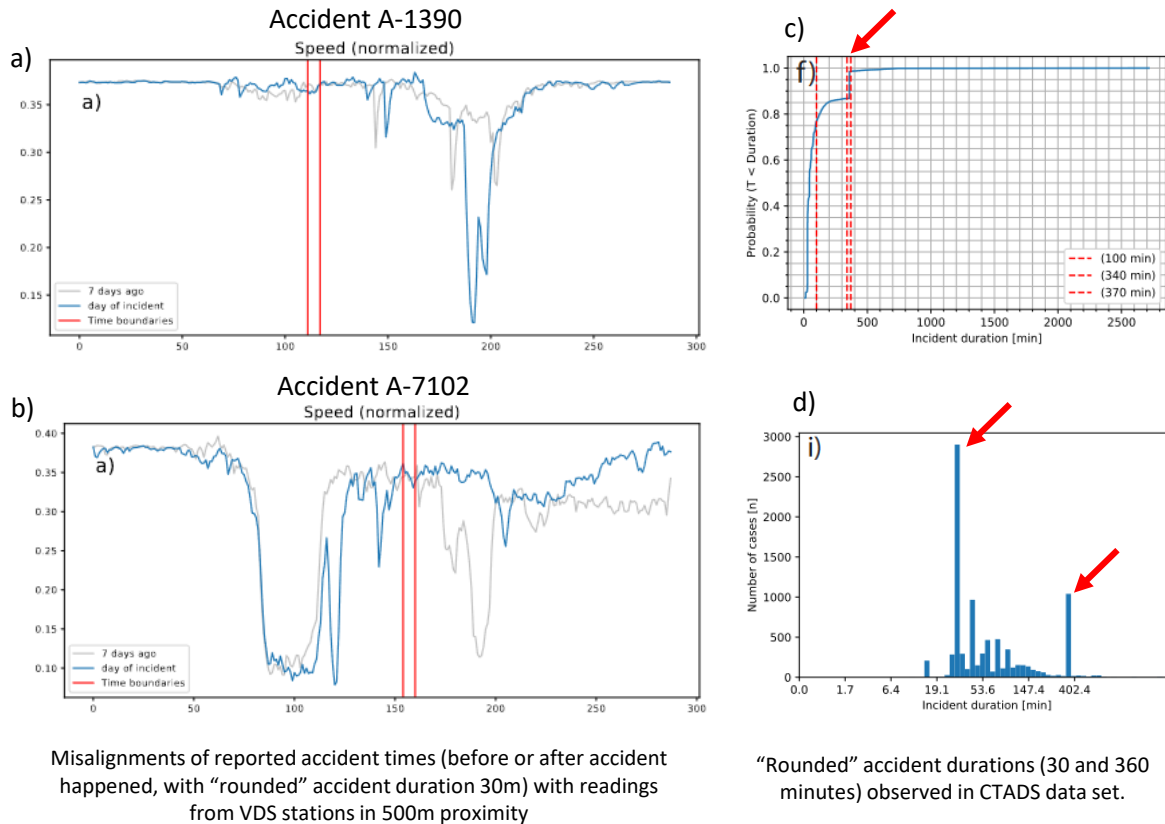
11 At this step we observed a significant amount of user-input errors in accident reports, which  
12 affect the accident duration/impact analysis: 1) accidents can be reported earlier or later than its oc-  
13 currence (observable disruption misalignment in time) 2) a report can be filled with "placeholder"  
14 duration values not representing the actual accident duration 3) there may be no observable disrup-  
15 tion in traffic speed despite the accident report (due to placement and management of the accident)  
16 (false positive) 4) there may be accident-related traffic disruptions not grasped by accident reports  
17 (false negative). Therefore, incorrect accident start time, duration and end time, unreported pres-  
18 ence or absence of disruption make it necessary to estimate accident duration characteristics from  
19 traffic state data instead of relying on user reports. In this paper our proposed methodology is really  
20 meant to solve the user-reporting issues related to traffic accidents and to be applied automatically  
21 on any data set, regardless of its nature or geo-location.

22 The use of PeMS data set allows to estimate the impact of accidents on the traffic states  
23 (flow, speed). For our scenarios, we choose the area of San-Francisco with accidents recorded  
24 from 2016 to 2020 in the CTADS data set. We then obtain Vehicle Detector Station locations from  
25 PeMS, the road network shape from OpenStreetMap and we perform an association of CTADS  
26 accident reports with VDS stations along the road within 500m proximity. We then try to segment  
27 the disruption time interval occurred on the day of an accident. Further, we associate observed  
28 disruptions in the traffic speed series with actual accident reports. The purpose of this step is to  
29 reduce user-input errors in accident reports and to enhance the modelling of traffic disruptions with  
30 an analysis of traffic speed.

### 31 Metric performance comparison

32 We apply the difference metrics detailed earlier in Section 5.1 to a monthly traffic speed/flow pro-  
33 file (monthly readings averaged to one day) and reading on the day of traffic accident. There are  
34 two approaches to applying the difference calculation: 1) a global difference - when we try to  
35 find the difference between monthly profile and traffic flow/speed readings on the day of accident;  
36 the global approach is too broad and will not allow the actual comparison between disruptions  
37 localized in time (very long but subtle disruptions can be measured as the same amount of dis-  
38 ruption as abrupt but impactful one). We measure the amount of difference that occurred within a  
39 moving time window (we choose twelve 5-minute time intervals equivalent to one hour). Traffic  
40 speed/flow readings from the moving window are taken right before the currently observed value  
41 to ensure that the difference estimation algorithm is not looking into the future. It will be further  
42 necessary to perform the prediction of difference values.

43 To compare the metric performances we provide an example of speed readings from one of  
44 the detector stations. Each difference metric demonstrates its specifics as represented in Figure 7:

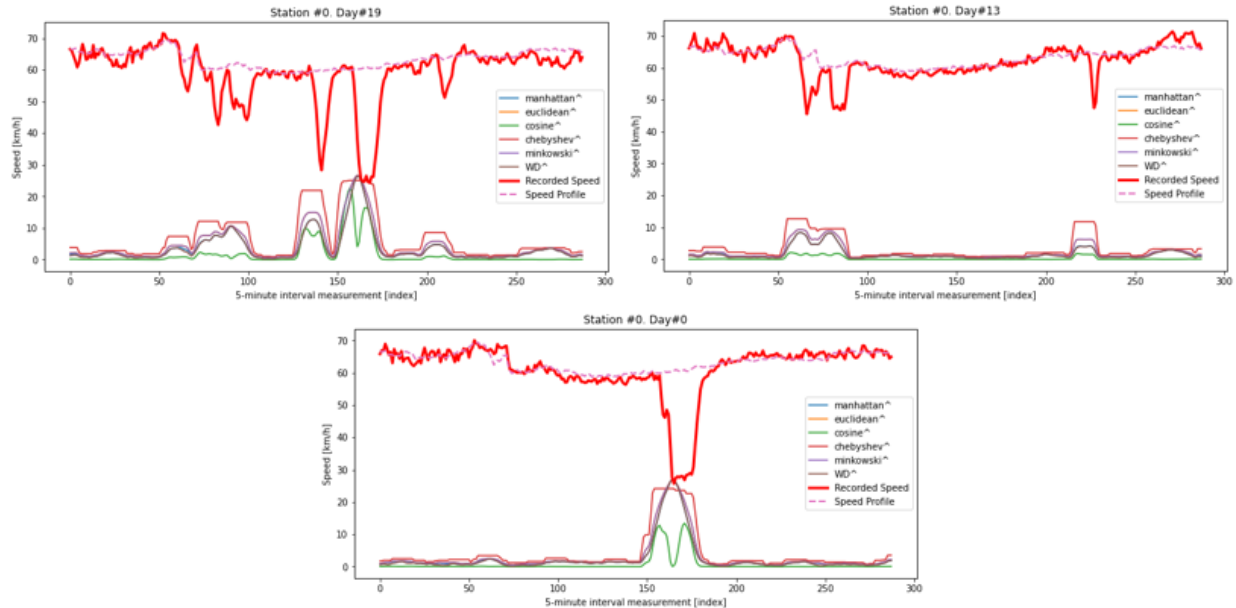


**FIGURE 6 User input errors located within the CTADS data set**

1) the Chebyshev metric, which we define as the maximum difference between the monthly profile and the observed readings, produces a noticeably rectangular shape and demonstrates a higher selectivity towards major disruptions than other metrics; the Chebyshev metric will be further used for the automated accident segmentation; 2) the Cosine distance detects a change in the traffic flow state - speed fall and increase both represented as positive peak values, 3) the  $D_{WD}$  allows for smooth representation of the amount of disruption (conceptually, it measures the amount of work necessary to change one shape into another, which we can rephrase as the amount of work produced by an accident to deviate the traffic state from the normal operation), 4) the Minkowsky, Euclidean and Manhattan difference metrics show little to no difference to the Wasserstein distance; we choose to use the  $D_{WD}$  since its connection to physical interpretation.

### Automated disruption segmentation results

Figure 8 presents the results obtained from our algorithm for the automated disruption segmentation. The segmentation line (dotted blue) represents the estimated disruption intervals represented as 0 and 1 to perform our visualisation investigation better. Figure 8a) shows that there may be multiple observed disruptions in a  $300 \times 5 = 1500$  time interval. Due to errors in accident reports regarding the starting time and the duration of the accident, it is non-trivial to determine which disruption is associated with the accident. The situation may be easier in the case when only one disruption is observed during the day. According to our algorithm we select the largest disruption on the day the accident reported. Figures 8b) and 8c) highlight additional specific situations which



**FIGURE 7 Various metrics applied to difference between recorded speed and speed profile**

1 need to be considered: 1) higher traffic speed at the end of the day than observed from the monthly  
 2 profile, 2) unstable traffic speed approaching normal traffic conditions with high frequency, 3)  
 3 slight misalignment of disruption intervals with the visually observed disruption intervals. All  
 4 these problems can be addressed by instead of using manual segmentation we can deploy several  
 5 Deep Learning models since it has advanced computer vision methods proposed in recent years.

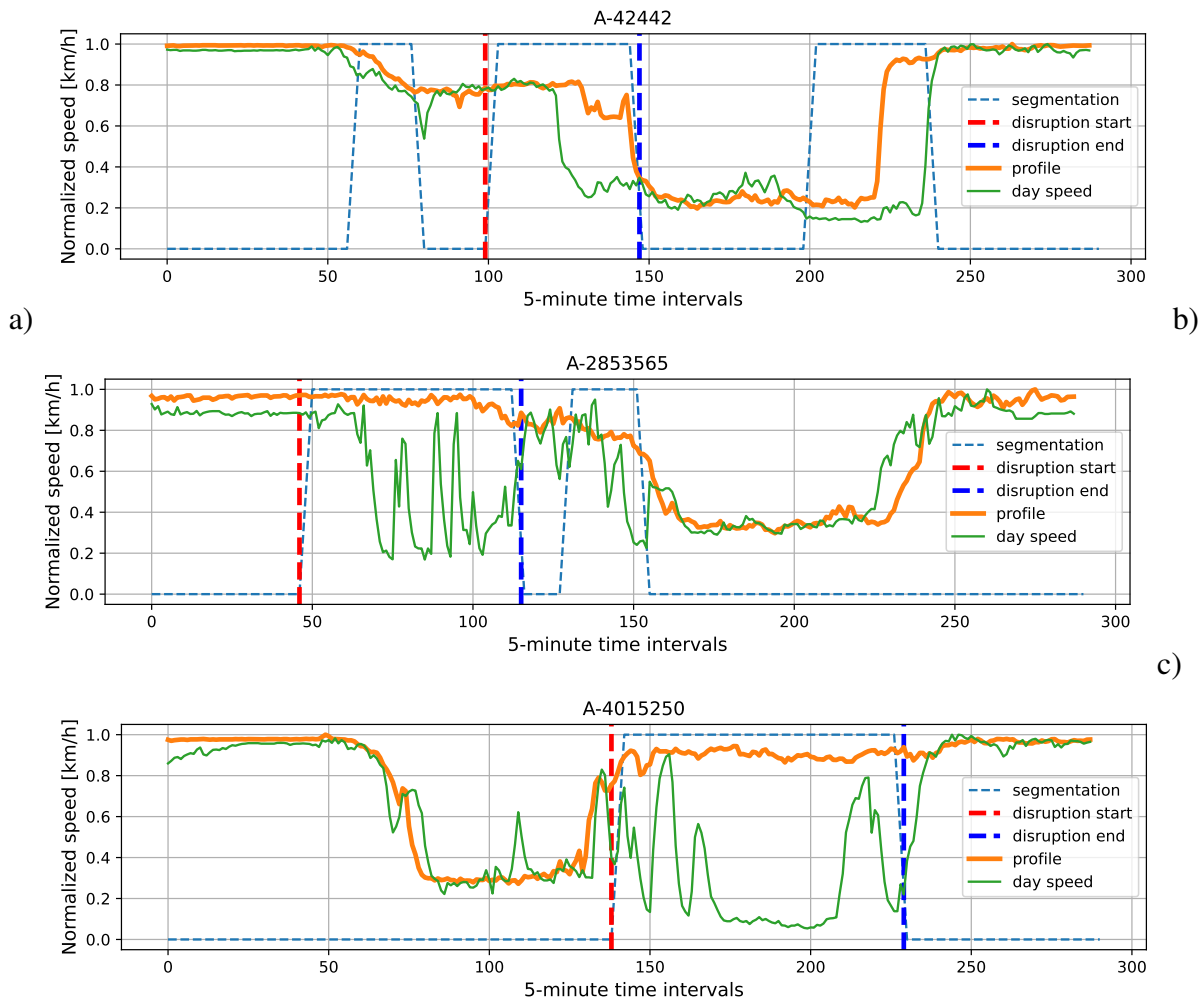
#### 6 **Comparison of our estimated durations versus the official reported accident durations**

7 There is a significant difference between the estimated and the reported accident durations that we  
 8 would like to highlight: 1) the reported accident durations contain a large amount of 30 and 360  
 9 minutes duration values (nearly 40% of data - see Figure 6d)), 2) the estimated accident durations  
 10 using our approach have an average duration of 58 minutes, while the reported is 108 minutes  
 11 (which is by assumption skewed due to 360 placeholder values), 3) the estimated accident durations  
 12 are distributed between 29 and 69 minutes (0.10 and 0.90 quantiles correspondingly) (see Figure  
 13 9a)), while the reported durations are distributed between 29 and 360 minutes (see Figure 9b)),  
 14 which highlights that disruptions observed from traffic speed are much shorter in reality than the  
 15 ones reported in the original data set, 4) there is no noticeable correlation between the observed and  
 16 reported durations - see Figure 10 in which we observe that the majority of correlation points are  
 17 on the lower side of the Ox axis, indicating that in reality the accident reports are much longer than  
 18 what we have succeeded to estimate via our proposed segmentation and early detection approach.

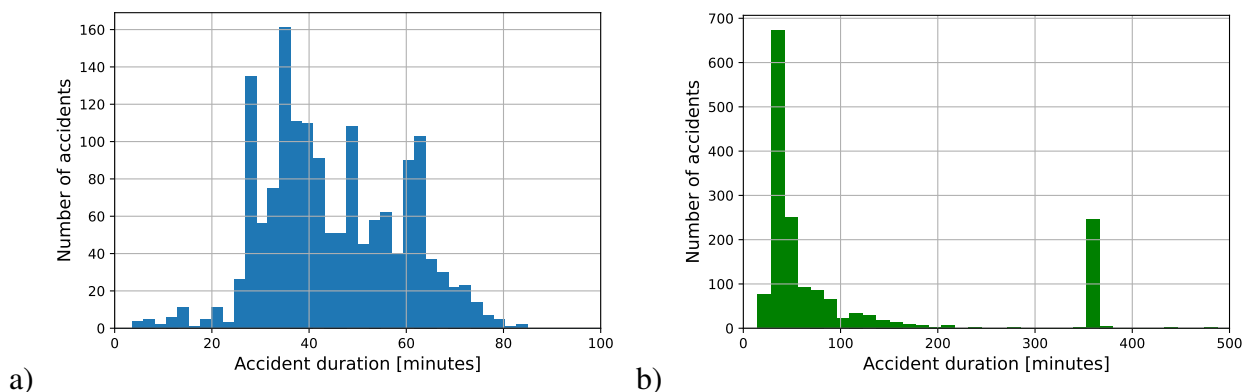
#### 19 **Extraction of disruption shapes**

20 Previously we applied a Chebyshev metric to perform the segmentation of disruptions. By applying  
 21 the Wasserstein difference between monthly speed profile and daily traffic speeds and extracting  
 22 the corresponding disruption interval we are able to observe the shape of disruption impact. The  
 23 Wasserstein distance (or metric), originally named an Earth Mover distance, has intuitive physical  
 24 interpretation - the minimum "cost" of altering one pile of earth into the other, which is assumed

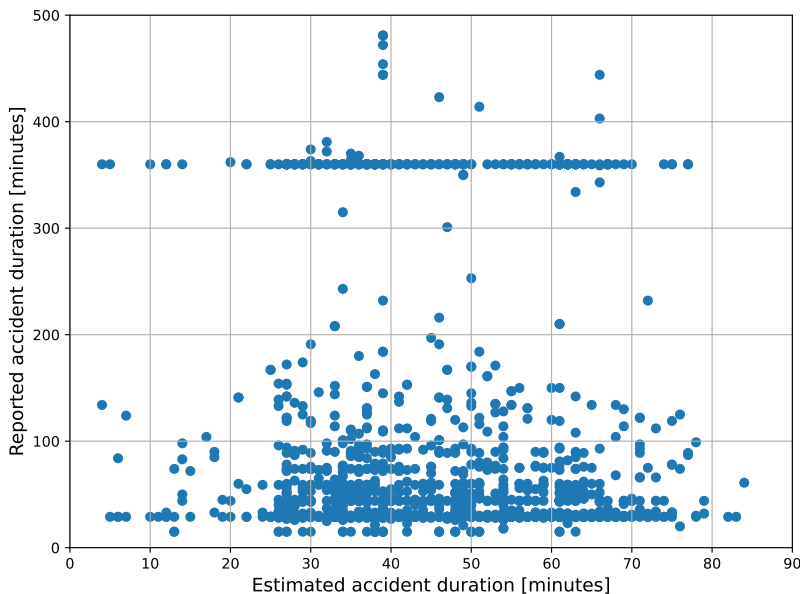




**FIGURE 8 Results for automated disruption segmentation algorithm**



**FIGURE 9 Distribution of accident durations for a) estimated and b) reported accident durations for the area of San Francisco**

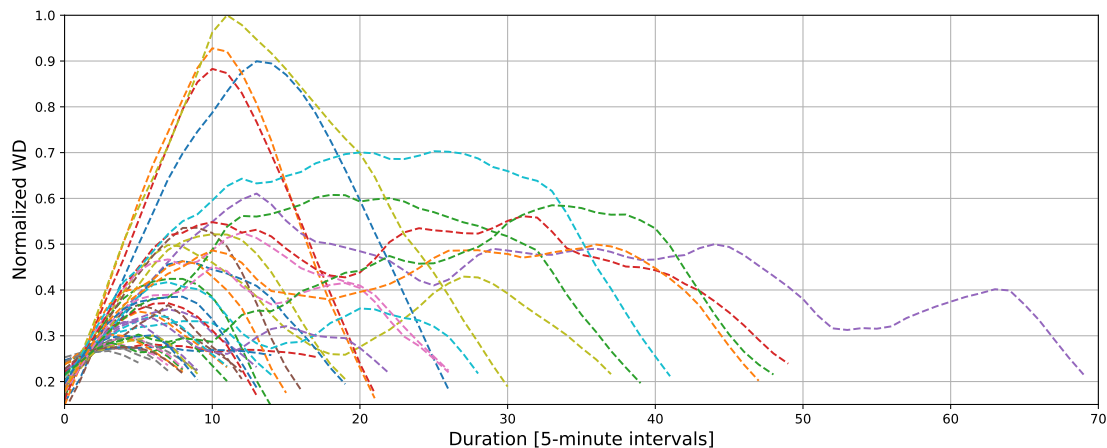


**FIGURE 10** Scatter plot for a) estimated and b) reported accident durations for the area of San Francisco

1 to be the amount of earth that needs to be moved times the mean distance it has to be moved. In  
 2 relation to the traffic state, it is the minimum amount of work necessary to alter the traffic state  
 3 to a disrupted condition, or in other words - the amount of traffic disruption. We compare the  
 4 normalized metric values since at every vehicle detector station there is a different average traffic  
 5 speed. As in our proposed algorithm, we use a 12-units moving window (one hour) to estimate the  
 6 Wasserstein difference between the traffic speed measurements and provide the plot for the first 40  
 7 segmented disruptions. It allows to conduct a shape analysis investigation in relation to the traffic  
 8 disruption amount (see Figure 11). Firstly, we observe the similarity between multiple disruptions  
 9 - they have a 'hill' shape with a high slope that almost doubles or triples the initial Wasserstein dis-  
 10 tance (e.g. from 0.2 to 0.5-0.6). Secondly, there are secondary relapses of the Wasserstein distance  
 11 (a double 'hill') which indicate long-lasting disruptions or cascading disruptions (a harder situation  
 12 to predict in our future works). The observed shapes can be defined through the parametric equa-  
 13 tion to perform the classification of disruption effects and facilitate the prediction of the disruption  
 14 impact timeline, since we observe that high-peak fast-ascending disruptions have a probability to  
 15 end sooner than slowly ascending ones (which indicate a long term effect or possibly large scale  
 16 accidents happening in the network).

### 17 Accident duration prediction results

18 We further compare a regression model prediction performance on the CTADS data set by using  
 19 on the training data set both our estimated versus the reported accident durations (see Table 1).  
 20 We report results of a 10-fold cross-validation over 1,792 accident reports for which we performed  
 21 a Vehicle Detector Station association. Firstly, we need to consider that the performance using  
 22 reported durations from CTADS can be affected because of the presence of user-input errors in  
 23 the form of placeholder values. Secondly, the nature of estimated accident durations is different  
 24 since accident response teams usually report the end of the accident at the moment they finished



**FIGURE 11 Normalized Wasserstein distance plot for disruption shapes extracted for segmented intervals**

- 1 the accident clearance, without estimating the time for the traffic to return to a normal condition,
- 2 which would require additional presence, calculations and access to measurements.

**TABLE 1 Accident duration prediction results**

Regression Model	$RMSE_{est}$	$RMSE_{rep}$	$MAPE_{est}$	$MAPE_{rep}$
KNN	<b>60.1</b>	85.7	46.0	<b>32.9</b>
RF	<b>57.8</b>	81.1	40.4	<b>28.6</b>
LR	<b>60.5</b>	129.6	<b>40.8</b>	72.7
SVM	<b>61.9</b>	113.8	<b>36.1</b>	50.3
GBDT	<b>59.3</b>	69.2	41.4	<b>31.1</b>

3 When we are using accident reports to predict the estimated accident duration, we obtain  
 4 a better performance using the RMSE metric across all the regression models, which may be con-  
 5 nected to the lower amount of long accident durations than reported (see Table 1 first column of  
 6  $RMSE_{est}$  which is much lower than the  $RMSE_{rep}$  by almost 41% in the case of RF for example,  
 7 and 45.6% in the case of SVM). When we compare the prediction results using the MAPE met-  
 8 ric, we observe a significantly better performance only for the Linear Regression (LR) and the  
 9 Support Vector Machines (SVM) - 48.87% and 28.2% improvement respectively, but for k-nearest  
 10 neighbours (KNN), Random Forest (RF) and Gradient-Boosted Decision Trees (GBDT), results  
 11 are lower, which can be explained by the presence of the high amount of 30-minute reported dura-  
 12 tions (30% of all data set), which may skew the model performance towards this value and gain an  
 13 misleading performance improvement. Therefore the better model performance may be deceiving  
 14 due to uncorrected user-input errors in accident reports. This makes the estimation of accident  
 15 duration from traffic state important and necessary allowing an in-depth analysis to be undertaken.

## 1 CONCLUSION

2 The proposed methodology in this paper aims at detecting, segmenting and extracting the observed  
3 disruptions in the traffic speed which was modelled together with reported traffic accidents by  
4 traffic management centers. The approach is innovative in its distance metric approach for an  
5 automatic incident detection coupled with an incident segmentation which has shown to improve  
6 the incident prediction by almost 41% in RMSE across multiple machine learning models. By  
7 obtaining shapes of disruptions we lay the foundation for accident impact modelling. Many studies  
8 still rely on the modelling of reported accident durations and pre-defined parameters, while they  
9 can be estimated from traffic state measurements, which gives us more data than just aggregated  
10 variables (duration, start time, etc). By having information on how each accident affect the traffic  
11 flow, we can study the accident impact with precision.

12 **Limitations of this work:** The current modelling approach has been applied to a San  
13 Francisco data set due to its public availability and easiness to access. However, we would like  
14 to test the approach on multiple other countries and incident databases across the globe; the main  
15 challenge is the lack of both traffic states and traffic accidents logs to be released with synchronised  
16 timelines. **Future works:** We are currently modelling the cascading effect on traffic disruptions  
17 and how these can be automatically identified based on multiple incoming traffic state streams; the  
18 main challenge of detecting subsequent incidents lie in the time-span duration of the first incident  
19 which is normally stochastic in nature.

## 20 AUTHOR CONTRIBUTION STATEMENT

21 The authors confirm contribution to the paper as follows: study conception and design: A. Grig-  
22 orev, A-S. Mihaita, F. Chen; data collection: A. Grigorev; analysis and interpretation of results: A.  
23 Grigorev, K. Saleh; draft manuscript preparation: A. Grigorev, A-S. Mihaita. All authors reviewed  
24 the results and approved the final version of the manuscript.

## 25 ACKNOWLEDGMENT

26 This work has been done as part of the ARC Linkage Project LP180100114. The authors are highly  
27 grateful for the support of Transport for NSW, Australia. This research is funded by iMOVE  
28 CRC and supported by the Cooperative Research Centres program, an Australian Government  
29 initiative."

1 **REFERENCES**

- 2 1. Organization, W. H., *Global status report on road safety 2015*. World Health Organization,  
3 2015.
- 4 2. Administration, N. H. T. S., *Traffic safety facts 2013*. U.S. department of transportation,  
5 2013.
- 6 3. Kim, W. and G.-L. Chang, Development of a Hybrid Prediction Model for Freeway Inci-  
7 dent Duration: A Case Study in Maryland. *International Journal of Intelligent Transporta-*  
8 *tion Systems Research*, Vol. 10, 2011.
- 9 4. Theofilatos, A., G. Yannis, P. Kopelias, and F. Papadimitriou, Predicting Road Accidents:  
10 A Rare-events Modeling Approach. *Transportation Research Procedia*, Vol. 14, 2016, pp.  
11 3399–3405, transport Research Arena TRA2016.
- 12 5. Grigorev, A., A.-S. Mihaita, S. Lee, and F. Chen, Incident duration prediction using a bi-  
13 level machine learning framework with outlier removal and intra–extra joint optimisation.  
14 *Transportation Research Part C: Emerging Technologies*, Vol. 141, 2022, p. 103721.
- 15 6. Mihaita, A. S., Z. Liu, C. Cai, and M. Rizoio, Arterial incident duration prediction using a  
16 bi-level framework of extreme gradient-tree boosting. *CoRR*, Vol. abs/1905.12254, 2019.
- 17 7. Shafiei, S., A. Mihaita, H. Nguyen, C. Bentley, and C. Cai, Short-Term Traffic Predic-  
18 tion Under Non-Recurrent Incident Conditions Integrating Data-Driven Models and Traf-  
19 fic Simulation. In *Transportation Research Board 99th Annual Meeting*, 2020, pp. –.
- 20 8. Mao, T., A.-S. Mihăiță, F. Chen, and H. L. Vu, Boosted Genetic Algorithm using Machine  
21 Learning for traffic control optimization. *IEEE Transactions on Intelligent Transportation*  
22 *Systems*, 2021.
- 23 9. Li, R., F. C. Pereira, and M. E. Ben-Akiva, Overview of traffic incident duration analysis  
24 and prediction. *European transport research review*, Vol. 10, No. 2, 2018, p. 22.
- 25 10. Fukuda, S., H. Uchida, H. Fujii, and T. Yamada, Short-term prediction of traffic flow  
26 under incident conditions using graph convolutional recurrent neural network and traffic  
27 simulation. *IET Intelligent Transport Systems*, Vol. 14, No. 8, 2020, pp. 936–946.
- 28 11. Parsa, A. B., H. Taghipour, S. Derrible, and A. K. Mohammadian, Real-time accident  
29 detection: coping with imbalanced data. *Accident Analysis & Prevention*, Vol. 129, 2019,  
30 pp. 202–210.
- 31 12. Eboli, L., C. Forciniti, and G. Mazzulla, Factors influencing accident severity: an anal-  
32 ysis by road accident type. *Transportation Research Procedia*, Vol. 47, 2020, pp. 449–  
33 456, 22nd EURO Working Group on Transportation Meeting, EWGT 2019, 18th – 20th  
34 September 2019, Barcelona, Spain.
- 35 13. Tsubota, T., C. Fernando, T. Yoshii, and H. Shirayanagi, Effect of Road Pavement  
36 Types and Ages on Traffic Accident Risks. *Transportation Research Procedia*, Vol. 34,  
37 2018, pp. 211–218, international Symposium of Transport Simulation (ISTS’18) and  
38 the International Workshop on Traffic Data Collection and its Standardization (IWT-  
39 DCS’18)Emerging Transport Technologies for Next Generation Mobility.
- 40 14. Yannis, G., A. Dragomanovits, A. Laiou, T. Richter, S. Ruhl, F. La Torre, L. Domeni-  
41 chini, D. Graham, N. Karathodorou, and H. Li, Use of Accident Prediction Models in  
42 Road Safety Management – An International Inquiry. *Transportation Research Procedia*,  
43 Vol. 14, 2016, pp. 4257–4266, transport Research Arena TRA2016.
- 44 15. Miller, M. and C. Gupta, Mining traffic incidents to forecast impact. In *Proceedings of the*  
45 *ACM SIGKDD international workshop on urban computing*, 2012, pp. 33–40.

- 1 16. Kuang, L., H. Yan, Y. Zhu, S. Tu, and X. Fan, Predicting duration of traffic accidents  
2 based on cost-sensitive Bayesian network and weighted K-nearest neighbor. *Journal of*  
3 *Intelligent Transportation Systems*, Vol. 23, No. 2, 2019, pp. 161–174.
- 4 17. Xiao, S., Traffic accident duration prediction based on natural language processing and a  
5 hybrid neural network architecture. In *2021 International Conference on Neural Networks,*  
6 *Information and Communication Engineering*, SPIE, 2021, Vol. 11933, pp. 194–202.
- 7 18. Alkaabi, A. M. S., D. Dissanayake, and R. Bird, Analyzing clearance time of urban traf-  
8 fic accidents in Abu Dhabi, United Arab Emirates, with hazard-based duration modeling  
9 method. *Transportation Research Record*, Vol. 2229, No. 1, 2011, pp. 46–54.
- 10 19. Moosavi, S., M. H. Samavatian, S. Parthasarathy, R. Teodorescu, and R. Ramnath, Ac-  
11 cident risk prediction based on heterogeneous sparse data: New dataset and insights. In  
12 *Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Ge-*  
13 *ographic Information Systems*, 2019, pp. 33–42.
- 14 20. Moosavi, S., M. H. Samavatian, S. Parthasarathy, and R. Ramnath, A countrywide traffic  
15 accident dataset. *arXiv preprint arXiv:1906.05409*, 2019.
- 16 21. Chen, C., K. Petty, A. Skabardonis, P. Varaiya, and Z. Jia, Freeway performance mea-  
17 surement system: mining loop detector data. *Transportation Research Record*, Vol. 1748,  
18 No. 1, 2001, pp. 96–102.