

1 **THE WHO AND WHERE OF ROAD SAFETY: EXTRACTING SURROGATE INDICATORS**  
2 **FROM SMARTPHONE-COLLECTED GPS DATA IN URBAN ENVIRONMENTS**

3  
4  
5  
6 **Joshua Stipanovic, Corresponding Author**, PhD Candidate  
7 Department of Civil Engineering and Applied Mechanics, McGill University  
8 Room 391, Macdonald Engineering Building, 817 Sherbrooke Street West  
9 Montréal, Québec, Canada H3A 0C3  
10 Email: joshua.stipanovic@mail.mcgill.ca

11  
12 **Luis Miranda-Moreno**, Associate Professor  
13 Department of Civil Engineering and Applied Mechanics, McGill University  
14 Room 268, Macdonald Engineering Building, 817 Sherbrooke Street West  
15 Montréal, Québec, Canada H3A 0C3  
16 Phone: (514) 398-6589  
17 Fax: (514) 398-7361  
18 Email: luis.miranda-moreno@mcgill.ca

19  
20 **Nicolas Saunier**, Associate Professor  
21 Department of Civil, Geological and Mining Engineering  
22 Polytechnique Montréal, C.P. 6079, succ. Centre-Ville  
23 Montréal, Québec, Canada H3C 3A7  
24 Phone: (514) 340-4711 x. 4962  
25 Email: nicolas.saunier@polymtl.ca

26  
27  
28 Word count: 5490 words + 8 tables/figures x 250 words (each) = 7490 words  
29  
30  
31  
32  
33  
34  
35

August 1<sup>st</sup>, 2015

**1 ABSTRACT**

2 Environment and driver behaviour are significant contributory factors in traffic collisions. Surrogate safety  
3 measures, non-crash measures that are physically and predictably related to crashes, provide opportunities  
4 for user-centric approaches to road safety and reduce dependency on crash data in environment-centric  
5 approaches. The purpose of this study is to extract surrogate safety measures from the smartphone-collected  
6 GPS data of regular drivers and to analyze those measures from an environment-centric and user-centric  
7 perspective. GPS travel data was collected using the Mon Trajet smartphone application in Quebec City,  
8 Canada over 21 days. Crash data was obtained from the Ministry of Transportation Quebec for a five year  
9 period from 2006 to 2010. The selected surrogate indicator, hard braking events (HBEs), demonstrated a  
10 spatial correlation of 0.67 with collision occurrence. Despite strong correlation, HBEs tend to overestimate  
11 risk on highway facilities and underestimate risk on local and arterial streets as the sample data collected  
12 from regular drivers likely over-represents travel on highways and under-represents travel on urban streets.  
13 The user-centric analysis showed that more HBEs occur during the AM and PM peak periods, and that  
14 braking in the PM peak period tends to be more severe, demonstrating that HBEs are not only spatially  
15 correlated with actual collision occurrence, but also make sense intuitively with respect to the behaviours  
16 related to collision occurrence. Future work will determine if other surrogate indicators that are more  
17 closely correlated with collision occurrence can be extracted, and disaggregating the analyses by facility  
18 type should improve the results.

19

20

21

22

23

*Keywords:* surrogate safety, smartphone, GPS, urban, collision prediction, behaviour, probe vehicles

## 1 INTRODUCTION

2 Environment and driver behaviour are significant contributory factors in traffic collisions and important  
3 influencers of road safety (1). Naturally, effective safety improvements should be environment-centric  
4 (addressing the ‘where’ of collisions) or user-centric (addressing the ‘who’ of collisions) or both.  
5 Environment-centric approaches involve identifying and remediating high-risk sites that “create an  
6 increased risk of unforeseeable accidents” due to their design or operation (2). Screening methods based on  
7 crash frequency or severity ranking criteria have traditionally been used to identify hazardous locations.  
8 Unfortunately, crash-based methods are reactive (2), require long collection periods to accumulate the  
9 necessary volume of crash data for analysis (3), are subject to errors and omissions in collision databases,  
10 and are sensitive to crash underreporting (4). These issues are particularly important in developing countries  
11 where the lack of reliable crash data inhibits implementation of crash-based techniques. User-centric  
12 approaches attempt to understand the relationship between driver behaviour and crash occurrence (1), often  
13 using naturalistic driving data collected unobtrusively in crashes, near crashes, and normal conditions.  
14 Naturalistic methods provide information difficult to observe by other techniques (5, 6) and allow for the  
15 use of surrogate safety measures based on behaviour rather than indicators based on collision statistics.  
16 Surrogate safety measures are non-crash measures that are physically and predictably related to crashes (7),  
17 and provide opportunities for user-centric approaches to road safety while reducing dependency on crash  
18 data in environment-centric approaches (8).

19 Naturalistic approaches typically yield large volumes of data from which surrogate indicators must  
20 be identified (5). Various methods for analyzing naturalistic data have been proposed, including the use of  
21 human observers. Though human observation practically limits the amount of data that can be analyzed and  
22 measurements may be subjective (8), human judgement provides a level-of-detail beyond what is currently  
23 possible through objective techniques (8). Compared to human observation, roadside-based sensors  
24 increase the sampling rate of road users and improve objectivity. Among methods for surrogate safety  
25 analysis, the traffic conflict technique using video-based sensors and computer vision techniques has been  
26 popular for before and after studies (2). Though video-based sensors provide high temporal resolution (2)  
27 and rich positional data beyond counts and speed (9), the analysis of video data is potentially time and  
28 resource intensive (2, 8), and interpretation of video data in behaviour terms requires additional  
29 consideration (8). Indicators based on traffic parameters collected by traditional point sensors including  
30 loops, radar, or other sensors (10, 11, 12) have yet to be proven as reliable surrogate safety measures, and  
31 the costs of these technologies make it impractical to implement them across an urban network (13).

32 In-vehicle sensors provide the best opportunity for collecting spatio-temporal naturalistic driving  
33 data within a road network. Instrumented vehicles (probe vehicles or floating car data) act “as moving  
34 sensors, continuously feeding information about traffic conditions” (14). GPS devices are reliable sources  
35 of naturalistic driving data (15) and may be complemented by additional vehicle kinematics from  
36 accelerometers or gyroscopes and environmental factors collected by external sensors such as radars. These  
37 sensors provide long periods of continuous data for a small sample of road users (2). Though the method is  
38 limited in terms of the studied population of drivers, the spatial coverage of GPS data makes it ideal for  
39 studying environmental factors, and the naturalistic nature of GPS data makes it ideal for addressing  
40 behavioural factors (1). Furthermore, new technologies, including GPS-enabled smartphones, have made  
41 and will continue to make obtaining GPS data from vehicles easier over time. This leads to opportunities  
42 for real time data collection and safety analysis which is potentially interesting for emergency services.

43 The purpose of this study is to examine surrogate safety measures derived from probe vehicle data  
44 collected by the GPS-enabled smartphones of regular drivers. The objectives of this research are to correlate  
45 GPS-based surrogate measures to actual collision occurrence, to analyze those surrogate measures from  
46 both an environment-centric and user-centric perspective, and to discuss the strengths and limitations of  
47 GPS data in surrogate safety analysis.

## 48 LITERATURE REVIEW

49 Though probe vehicles have been widely used in spatio-temporal applications such as traffic monitoring  
50 and origin-destination studies (13), applications in road safety have been less common. Automated incident  
51

1 detection (AID) involves the identification of “non-recurring events such as accidents” through pattern  
2 classification of traffic flow (16), and improves safety by reducing secondary collisions (17). Existing  
3 techniques using dedicated probe vehicles have low penetration rates and are insufficient for providing “an  
4 exhaustive coverage of the transportation network” (13). Therefore, probe vehicles are often used in  
5 conjunction with traditional roadside sensors (14). In research applications, traffic simulation has been used  
6 to achieve proportions of dedicated probe vehicles beyond that which is possible in the field. Sethi et al.  
7 (17) found that probe vehicles improve successful incident detection rates and decreased false alarm rates  
8 over roadside sensors alone, though only when using a proportion of dedicated probe vehicles beyond what  
9 could be expected in practice (17). Dia and Thomas (16) similarly achieved the best results when probe  
10 vehicles comprised 20% of the traffic stream (16).

11 User-centric approaches using probe vehicles have been somewhat infrequent. Fazeen et al. (19)  
12 used smartphone accelerometer data to classify ‘safe’ accelerations and decelerations from ‘unsafe’ ones  
13 (approximately 3 m/s<sup>2</sup> or greater), though failed to demonstrate whether ‘unsafe’ behaviour led to increased  
14 collision risk. Jun, Ogle, and Guensler (15) analyzed the relationship between spatio-temporal driving  
15 activity and likelihood of crash involvement. Using dedicated GPS devices and self-reported safety data,  
16 the study found that drivers involved in crashes tended to travel longer distances and at higher speeds, and  
17 also “engaged in hard deceleration events” (greater than 2.7 m/s<sup>2</sup>) more frequently (15). Though failing to  
18 show a causal link between decelerations and collision risk, the authors suggest that decelerations “may be  
19 employed as roadway safety surrogate measures” (15). Although behavioural studies typically consider  
20 differences in demographics, they frequently fail to consider temporal and spatial factors (1). Ellison,  
21 Greaves, and Bliemer (1) studied 106 drivers using GPS devices to collect speed, speed limit, location, and  
22 timestamp for every second of vehicle operation, along with demographic surveys for each driver. By  
23 controlling for temporal and spatial factors including geometry, weather, time of day, trip purpose, and  
24 vehicle occupancy, the authors found that the road environment was a significant influencer of driver  
25 behaviour (1). However, as 90% of all traffic collisions involve behavioural factors (20), this points to the  
26 strong indirect effect of road environment on safety through behavioural influence.

27 Probe-based surrogate safety measures aim to identify drivers avoiding collisions through evasive  
28 manoeuvres including steering, braking, or accelerating (21). Although speed is often regarded as an  
29 important surrogate measure, changes in speed (acceleration or jerk) may be more important (8). Algerholm  
30 and Larhmann (2) used data collected from 6 drivers over a 3 month period using GPS devices and  
31 accelerometers. The authors stated that “braking was the evasive action [...] in 88% of the accidents in  
32 built-up areas” (2), making decelerations a logical indicator to extract. Jerk was found to be correlated with  
33 accident occurrence both across drivers (user-centric) and across sites (environment-centric) (2). Bagdadi  
34 (5) noted that the most common crashes are rear-end collisions, and used GPS, accelerometer, and radar  
35 data from 109 participants. The proposed surrogate measure based on jerk was used to correctly identify  
36 self-reported near misses at an 86% success rate (5). One shortcoming of this study is that the ground truth  
37 data used was itself a surrogate measure (near misses) and not actual collision data.

38 Several shortcomings are apparent in the existing literature, which this study attempts to address.  
39 First, there has been no attempt to derive surrogate safety measures from smartphone-collected GPS data  
40 of regular drivers alone. Existing studies have used dedicated probe vehicles (resulting in sample sizes of  
41 100 drivers or less) or dedicated GPS devices with supplemental accelerometer data. Second, there has been  
42 no comprehensive comparison of GPS-based surrogate indicators to large quantities of crash data. Instead,  
43 studies have compared indicators to sample safety data, which is often self-reported. Thirdly, there has been  
44 little effort to consider user-centric and environment-centric approaches simultaneously.

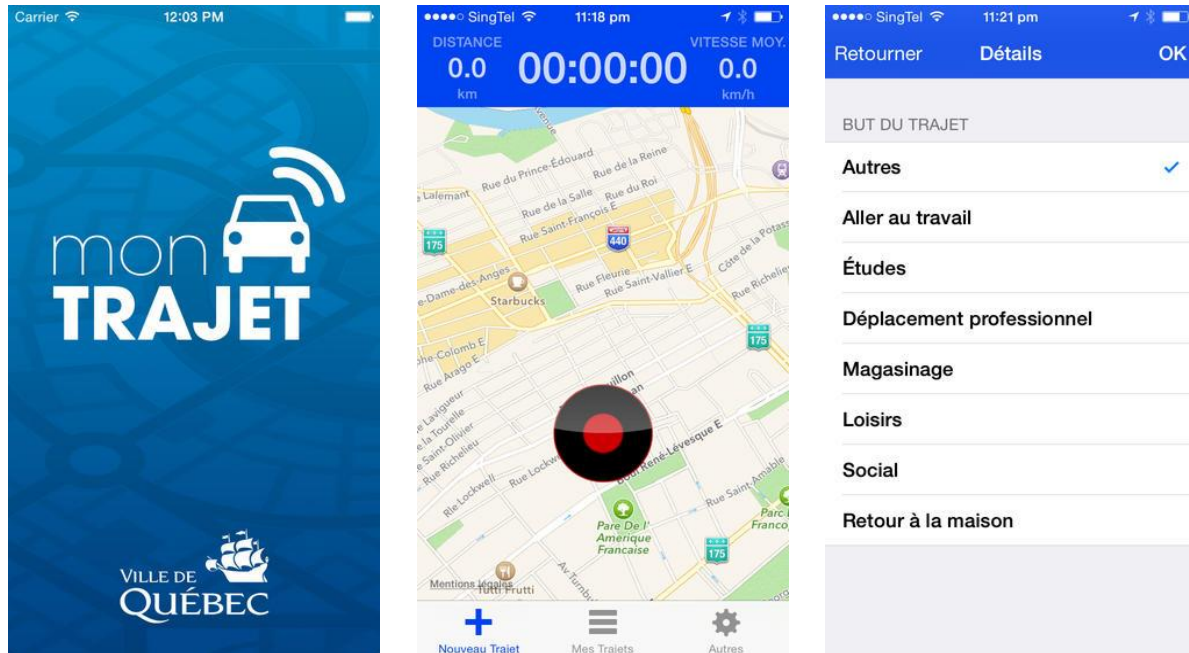
## 45 **METHODOLOGY**

### 46 **Data Collection**

47 Naturalistic driving data should be collected as unobtrusively as possible, to ensure data accurately  
48 represents normal driving conditions. Collecting GPS data from smartphones allows for the study of regular  
49 drivers using a system that minimally impacts their behaviour, and the implementation of a smartphone  
50  
51

1 application makes use of devices already widely available to the driving population reducing cost and  
 2 increasing potential sample size. Smartphone applications, such as Mon Trajet (22, 23) by Brisk Synergies  
 3 (24), shown in FIGURE 1, are installed voluntarily by drivers and collect GPS data anonymously. General  
 4 trip information, including route, origin and destination, and start and end time, are captured for every user-  
 5 reported trip logged in the application. Travel is described by observations including user speed, latitude,  
 6 longitude, and altitude captured for every 1-2 seconds of vehicle operation. Other socio-demographic  
 7 information may also be available depending on the configuration of the application. Once a trip has been  
 8 collected and reported by the user, initial pre-processing of the data is completed using Kalman filtering to  
 9 reduce data variability. The GPS data is then stored in remote databases, from which the raw GPS  
 10 observations are exported for further analysis.

11



12

13 **FIGURE 1 Smartphone application interfaces**

#### 14 **Data Cleaning**

15 Although GPS data from a smartphone application is rich in spatio-temporal data, raw GPS traces contain  
 16 variability in both position and speed. Even with pre-processing of the user data, additional data cleaning  
 17 methods are required. Although Kalman filtering is popular, the method only smooths vehicle positions in  
 18 terms of a latitude and longitude and does not explicitly link trips to the road network. This study used a  
 19 map-matching process to ensure that trips are correctly matched to the links in the road network where they  
 20 occurred. Vehicle speed measurements were cleaned using exponential smoothing.

21

#### 22 *Map Matching*

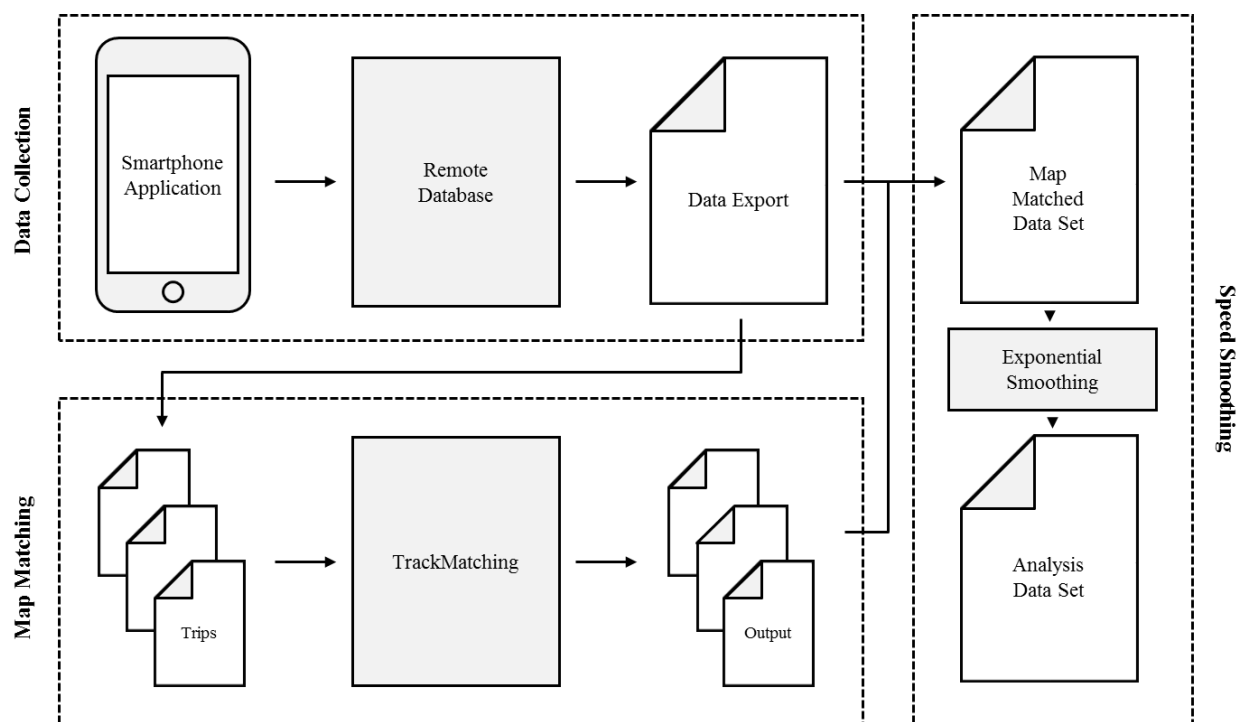
23 TrackMatching is a commercially available, cloud-based web map-matching software service (25) that  
 24 matches GPS data trip data to the OpenStreetMap (OSM) road network (26). Before GPS data is sent to  
 25 TrackMatching, the data must be split into individual trips and formatted according to the input  
 26 requirements of the software, including only the coordinate id, timestamp, latitude, and longitude for each  
 27 observation. The software returns a map-matched OSM ID (link ID), map-matched latitude and longitude,  
 28 and source and destination nodes along the OSM link for each GPS observation. Importantly, because much  
 29 important information (user ID, speed, timestamp, etc.) is lost through the map matching process, the results  
 30 must be merged back with the original data to preserve the complete data set.

### 1 *Speed Smoothing*

2 Simple exponential smoothing is used to eliminate noise and outliers within the GPS-collected speed data  
3 by generating new speed estimates for each observation (27). The smoothing equation is given by

$$4 \hat{Y}_t = \alpha Y_t + (1 - \alpha) \hat{Y}_{t-1} \quad (1)$$

7 where  $\hat{Y}_t$  is the estimated speed at the current time step,  $t$ ;  $Y_t$  is the speed at the current time step as measured  
8 by the GPS,  $t$ ;  $\hat{Y}_{t-1}$  is the estimated speed at the previous time step,  $t - 1$ ; and  $\alpha$  is the smoothing parameter.  
9 Increasing alpha increases the effect of the observed speed (less smoothing) while decreasing alpha  
10 increases the effect of the last estimated speed (more smoothing). If alpha is too large, then smoothing is  
11 minimal, and noise in the GPS data may be wrongly interpreted as braking (type II error). If alpha is too  
12 small, then smoothing can potentially eliminate actual braking events (type I error). Smoothing parameters  
13 of 0.4, 0.6, and 0.8 were tested. The process of collecting and cleaning the GPS data is illustrated in  
14 FIGURE 2.



16

17 **FIGURE 2 Collection and cleaning of smartphone-collected GPS data**

### 18 **Extraction of Surrogate Indicators**

19 After the data has been collected and cleaned, the surrogate safety measures are extracted from the analysis  
20 data set. Recognizing that deceleration is perhaps the most common evasive manoeuvre in urban areas (2),  
21 selecting hard braking events (HBEs) as the surrogate indicator of interest is logical. Studies focussed on  
22 deceleration have used jerk, observed using accelerometers, to define the surrogate indicator (2, 5). When  
23 using GPS data alone, calculating jerk-based surrogate safety measures is not possible, as GPS observations  
24 are too infrequent to capture the required detail. However, Fazeen et al. (19) suggested that decelerations  
25 exceeding  $3 \text{ m/s}^2$  were an indicator of ‘unsafe’ behaviour. Therefore, using a deceleration threshold may be  
26 sufficient to define HBEs. Although the  $3 \text{ m/s}^2$  threshold is a starting point to develop GPS-based surrogate  
27 indicators, thresholds of  $4 \text{ m/s}^2$  and  $5 \text{ m/s}^2$  were also tested. An algorithm was developed to automatically  
28 identify all instances where a vehicle exceeded the threshold. HBEs were then analyzed from both  
29 environment-centric and user-centric perspectives.

## 1 Environment-Centric Analysis

### 3 *Spearman Rank Correlation*

4 As surrogate safety measures must be predictably related to crashes (7), any proposed measure must  
 5 demonstrate correlation with actual safety or risk. Spearman's Rank Correlation Coefficient, or Spearman's  
 6 rho, indicates how strongly the dependency between two variables is described by a monotonic function  
 7 and is a popular choice for correlating surrogate indicators with crash data. Locations with the most  
 8 collisions should also have the most HBEs, and sites with fewer collisions should have fewer HBEs. A rho  
 9 of 1.0 indicates positive correlation, 0.0 indicates no correlation, and -1.0 indicates negative correlation.  
 10 Spearman's rho,  $\rho$ , is calculated using

$$12 \quad \rho = 1 - \frac{6 \sum (x_i - y_i)^2}{n(n^2 - 1)} \quad (2)$$

13 where  $x_i$  and  $y_i$  are the ranks of site  $i$  in the two data sets and  $n$  is the total number of sites. The ranks,  $x_i$   
 14 and  $y_i$ , were created by generating buffers around each intersection (any point where two OSM links meet  
 15 or intersect) in the road network using GIS. The total numbers of collisions and HBEs within the buffers  
 16 were then counted, and intersections were ranked based on these counts. The effect of buffer size on  
 17 correlation was determined by comparing results generated with 100 m, 200 m, and 500 m buffers.  
 18

### 20 *Hot Spot Analysis*

21 Although Spearman's rho generally quantifies the correlation between a surrogate safety measure and  
 22 collision occurrence, additional analysis is necessary to observe discrepancies between the data sets. Heat  
 23 maps generated using GIS can be compared visually to determine where the surrogate measure performs  
 24 well (has strong agreement with the crash data) and where performance is poor. Heat maps were generated  
 25 in GIS for both collisions and HBEs using a 500 m radius and 50 m pixel width.

## 27 User-Centric Analysis

28 Rather than just considering the locations in the road network where hard braking events occur, a user-  
 29 centric approach to surrogate safety can show which characteristics or behaviours of drivers contribute to  
 30 the occurrence or severity of HBEs and therefore collisions (if the link between HBEs and collisions is  
 31 demonstrated). Though driver socio-demographics were not collected, consideration for spatio-temporal  
 32 driving behaviour (15) is possible based on GPS data. The user-centric analysis was completed using two  
 33 ordered logit models. The first model was used to analyze the occurrence of hard braking events. In this  
 34 model, trips were divided into trips with at least one HBE above 3 m/s<sup>2</sup> (Alternative 1), and trips with none  
 35 (Alternative 0). The dependent variables included trip characteristics of length and average speed, and time-  
 36 of-day characteristics indicating whether the trip occurred during the AM peak period (6:00 AM to 9:00  
 37 AM), PM peak period (4:00 PM to 7:00 PM) or at night (10:00 PM to 4:00 AM), and whether the trip was  
 38 made on a weekday or on the weekend.

39 A second ordered logit model was used to analyze the severity of braking events. In this model,  
 40 trips without HBEs above 3 m/s<sup>2</sup> were ignored. The remaining trips were grouped according to the hardest  
 41 braking event experienced during the trip; Alternative 0, 3-4 m/s<sup>2</sup>; Alternative 1, 4-5 m/s<sup>2</sup>; and Alternative  
 42 2, 5 m/s<sup>2</sup> or greater. The same dependent variables were considered with the addition of instantaneous  
 43 vehicle speed immediately before the HBE occurred.

## 45 DATA DESCRIPTION

46 This study made use of three primary data sources. GPS travel data was collected in Quebec City, Canada  
 47 using the Mon Trajet application by Brisk Synergies (24). In total, approximately 5000 driver participants  
 48 have logged nearly 50,000 trips using the application. However, the sample for this study contained 2413  
 49 drivers and 12,724 individual trips during the period between April 28 and May 18, 2014. Over the 21 days

1 sampled, 19.7 million individual data points were logged, with observations available every 1-2 seconds  
 2 during a trip. Crash data was obtained from the Ministry of Transportation Quebec for a five year period  
 3 from 2006 to 2010. 9248 collisions identified across the 5-year period involved at least one vehicle. Map  
 4 data used for the environment-centric analysis was obtained from OpenStreetMaps in order to maintain  
 5 consistency with the map matching results.

## 7 RESULTS

### 9 Extraction of Surrogate Indicators

10 TABLE 1 provides the number of HBEs identified for each combination of deceleration threshold and  
 11 smoothing parameter. Both of these variables were observed to greatly influence the total number of HBEs  
 12 that were identified. In the most restrictive case, only 1444 events were extracted, while the least restrictive  
 13 case found nearly 80,000 events.

14 **TABLE 1 Number of Hard Braking Events Identified**

		Alpha		
		0.8	0.6	0.4
Threshold	3 m/s <sup>2</sup>	78457	43119	13870
	4 m/s <sup>2</sup>	21744	9356	2719
	5 m/s <sup>2</sup>	6958	3021	1444

### 16 Environment-Centric Analysis

#### 18 *Spearman Rank Correlation*

19 Spearman's rho was calculated for three different deceleration thresholds, with three different smoothing  
 20 parameters, and three different buffer sizes, for a total of 27 tests. The results are presented in TABLE 2.  
 21 Smoothing parameters of 0.6 and 0.8 were found to provide roughly equivalent results, only differing by a  
 22 few percentage points. An alpha value of 0.4 was inferior in all test cases (as were alpha values less than  
 23 0.4). Although a higher alpha value results in improved correlation, reducing the alpha value significantly  
 24 reduces the number of events identified (as shown in TABLE 1), with only a minor reduction in correlation  
 25 strength. In this case, an alpha value of 0.6 provides good correlation with far fewer observations than an  
 26 alpha of 0.8. All cases with 100 m buffers failed to provide correlation above 0.50. A 200 m buffer  
 27 performed better, with correlations between 0.50 and 0.60, and the 500 m provided the best results, with  
 28 correlations between 0.60 and 0.70. A deceleration of threshold of 3 m/s<sup>2</sup> consistently provided the highest  
 29 correlation, up to 0.669. For these reasons, the remainder of this paper focusses on a threshold of 3 m/s<sup>2</sup>,  
 30 alpha of 0.6, and a buffer size of 500 m (rho = 0.644).

#### 32 *Hot Spot Analysis*

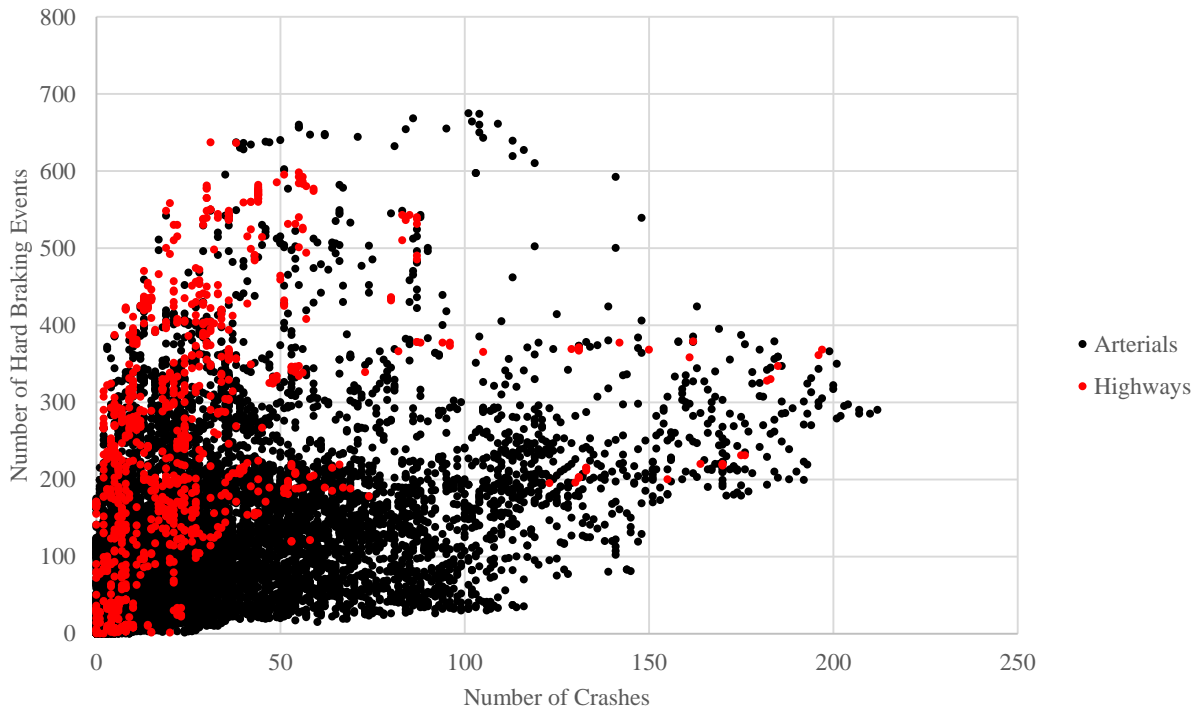
33 Despite relatively high correlation results, Spearman's rho provides no indication of where correlation is  
 34 good and where it is poor. In order to identify differences in the data sets, hot spots were identified for both  
 35 collisions, in FIGURE 4a, and HBEs, in FIGURE 4b. Visually, these maps reveal a critical difference  
 36 between the crash data and the surrogate safety measures. The locations with the most collisions tend to be  
 37 local streets, such as in downtown Quebec City, or on urban arterials like Laurier Boulevard and 1<sup>st</sup> Avenue.  
 38 In contrast, the locations with the most HBEs tend to be on highways, such as Félix-Leclerc, Henry IV, and  
 39 Charest. This is perhaps logical, as a deceleration of 3 m/s<sup>2</sup> is more likely when a driver is traveling at  
 40 highway speeds compared to urban arterials and local streets. This is a crucial discrepancy, as priority sites  
 41 identified through network screening would be very different if using HBEs rather than collision data.



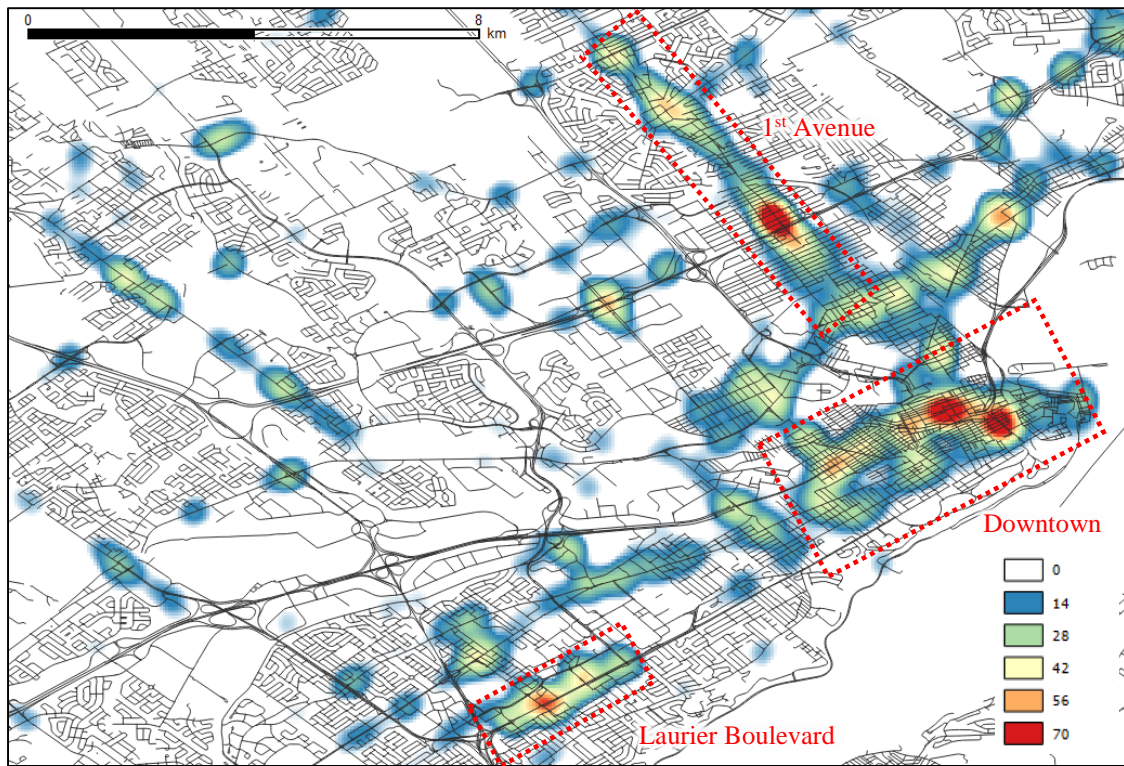
1 **TABLE 2 Spearman Rank Correlation Coefficients**

		Buffer Size		
		100 m	200 m	500 m
<b>Alpha = 0.8</b>				
Threshold	3 m/s <sup>2</sup>	0.497	0.564	0.669
	4 m/s <sup>2</sup>	0.453	0.522	0.639
	5 m/s <sup>2</sup>	0.386	0.474	0.610
<b>Alpha = 0.6</b>				
Threshold	3 m/s <sup>2</sup>	0.477	0.540	0.644
	4 m/s <sup>2</sup>	0.399	0.479	0.603
	5 m/s <sup>2</sup>	0.278	0.360	0.543
<b>Alpha = 0.4</b>				
Threshold	3 m/s <sup>2</sup>	0.408	0.465	0.580
	4 m/s <sup>2</sup>	0.244	0.321	0.500
	5 m/s <sup>2</sup>	0.129	0.171	0.341

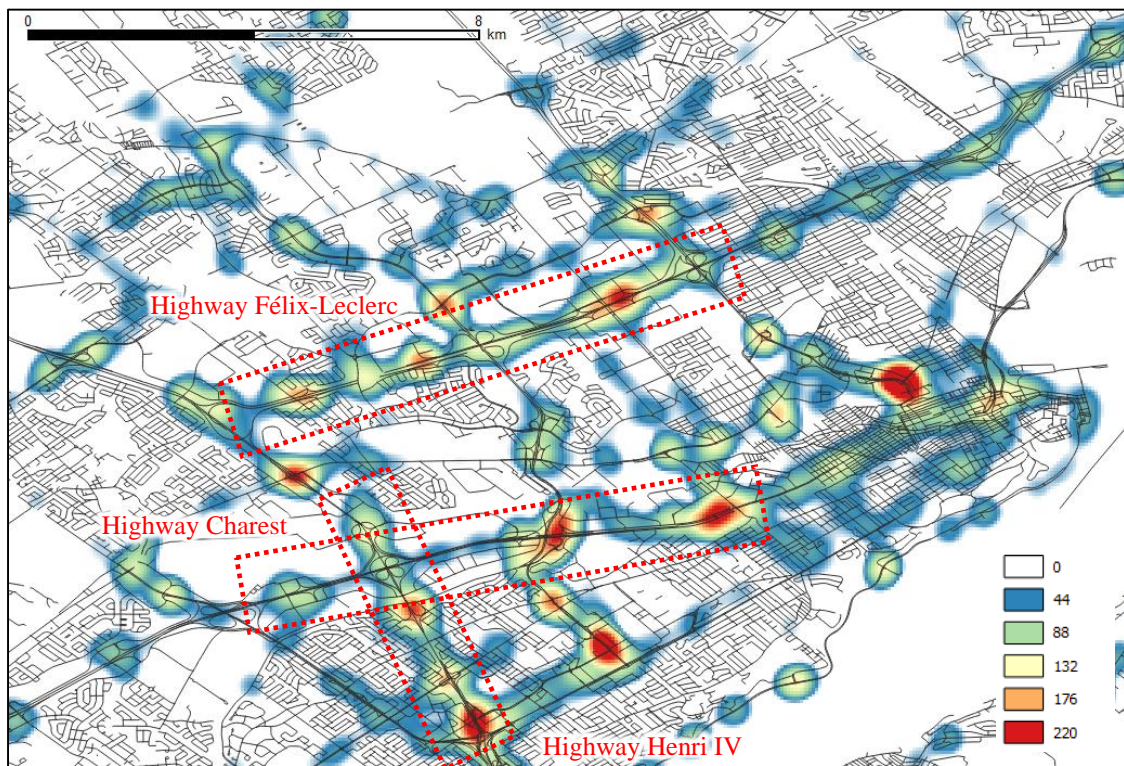
2  
 3 Considering a correlation exceeding 0.644 was observed despite the disparity in the locations of  
 4 the hot spots, it was believed that disaggregating facility types may improve the results. Consider the plot  
 5 of collisions and HBEs provided in FIGURE 3. Locations with many more HBEs than collisions are likely  
 6 to be highways (those that appear as hotspots in FIGURE 4b), while locations with more collisions are  
 7 likely to be local or arterial streets (those that appear as hotspots in FIGURE 4a).



8  
 9 **FIGURE 3 Correlation between number of hard braking events and crashes**



(a)



(b)

1  
2  
3  
4

5  
6  
7

8 **FIGURE 4** Hot spot analysis by collision data (a) and hard braking events (b)

1 Comparing the lists of intersections ranked by collision occurrence and HBEs, the top 1 %  
 2 identified using crash data shared no common sites with the and the top 1 % identified using HBEs. The  
 3 top 5 % in both lists contained only 64 intersections in common out of 778 total (8 % similarity), and the  
 4 top 25 % in both lists contained 999 intersections in common out of 3889 total (26 % similarity). With the  
 5 results presented above and the apparent bimodal nature of the surrogate/crash relationship, a final attempt  
 6 to improve the correlation results was made by separating highways from the local and arterial streets, and  
 7 performing a separate correlation for each facility group. Intersections, collisions, and HBEs were filtered  
 8 according to those corresponding to highways or highway ramps, and those occurring on all other facilities.  
 9 TABLE 3 shows the new rho values based on facility type. Although the correlation for highways was  
 10 lower than for all facility types combined, the correlation on local and arterial streets was improved, and an  
 11 average weighted on proportion of facility type demonstrated an overall increase in correlation. However,  
 12 as the total improvement in correlation was only two percentage points, this result demonstrates that  
 13 although the relationship between HBEs and collision occurrence is dependent on facility type, more facility  
 14 types must be considered if substantial improvements are to be made.

15 **TABLE 3 Spearman Rank Correlation Coefficient for Highways and Local/Arterials**

	Roadway Type			
	All	Local/Arterial	Highway	Weighted
$\rho$	0.644	0.674	0.536	0.664

16

17 **User-Centric Analysis**

18 The results for the braking occurrence model (which included all trips) and the braking severity model  
 19 (which included only those trips with at least on HBE exceeding 3 m/s<sup>2</sup>) are presented in TABLE 4. The  
 20 model results contain only parameters significant at 95% confidence unless otherwise noted.

21 **TABLE 4 Model Results for Occurrence and Severity of Hard Braking Events**

Explanatory variables	Braking Occurrence		Braking Severity	
	Parameter	z stat	Parameter	z stat
Instantaneous Speed	N/A	N/A	0.0847	22.9
Trip Speed	-0.0046*	-1.52*	-	-
Trip Length	-	-	-	-
AM Peak	0.0996	2.71	-	-
PM Peak	0.1655	4.19	0.0956	2.48
Night	-	-	-	-
Weekday	-	-	-	-
Tau 1	-0.2947		1.4784	
Tau 2	N/A		3.0387	
Number of cases	20840		12087	
Log likelihood at convergence	-14167.40		-11090.47	
Log likelihood for constants-only model	-14177.35		-11356.39	

*\*note: coefficient for trip speed is significant at 87% confidence*

1 In the braking occurrence model, only the AM and PM peak period variables were found to be  
2 statistically significant. This indicates that HBEs are more common during peak periods than during other  
3 periods in the day. This result is expected as the congestion experienced in these times should contribute to  
4 increased braking and, therefore, to increased collisions. Trip length was found to have no effect on braking  
5 occurrence, while average trip speed had a negative effect, though it was only significant at 87% confidence.  
6 Other time-of-day variables also failed to show a significant relationship with braking occurrence. In the  
7 braking severity model, the PM peak period variable was again found to be significant and positive. This  
8 indicates that not only do more HBEs occur in the PM peak period, but that they tend to be more severe  
9 (harder braking) than at other times of the day. Additionally, the instantaneous vehicle speed was found to  
10 be positive and significant. Faster vehicles who use braking as an evasive manoeuvre tend to brake more  
11 aggressively or severely (i.e. have a higher deceleration rate). This result is again intuitive, as it is expected  
12 that travel at higher speeds would require more severe evasive actions, and potentially more severe  
13 collisions.

## 14 **CONCLUSIONS**

15 The purpose of this study was to extract surrogate safety measures from the smartphone-collected GPS-  
16 data of regular drivers and to analyze those measures from both an environment-centric and user-centric  
17 perspective. The smoothing parameter and deceleration threshold have a strong influence the number of  
18 HBEs identified. However, buffer size has a much greater influence on the strength of the correlation to  
19 actual collision occurrence. Selecting these three parameters is a crucial step in the presented analysis. The  
20 strongest correlation between HBEs and collisions was 0.669 (threshold = 3 m/s<sup>2</sup>, alpha = 0.8, buffer = 500  
21 m) though reducing alpha to 0.6 only decreased correlation to 0.644 with far fewer observations. This is a  
22 promising result in this early research, as even the most aggregate approach to using HBEs as surrogate  
23 safety measures demonstrated a relatively high correlation.

24 Despite strong correlation, hot spots identified by both methods vary greatly. HBEs tend to  
25 overestimate risk on highway facilities and underestimate risk on local and arterial. There are several  
26 potential explanations for the disagreement. First, regardless of the sample drivers used, it is more likely  
27 that their trips will utilize highway facilities, while the probability of trips in residential neighbourhoods is  
28 low (except for the neighbourhoods where sample drivers live). Therefore, the sample of smartphone GPS  
29 data collected from regular drivers likely over-represents travel (and therefore collision risk) on highways  
30 and underrepresents travel (and risk) on urban residential streets. Disaggregating analysis by facility type  
31 showed potential for improvement, although the analysis must consider more than two facility types if  
32 improvements are to be substantial. Therefore, facility types should be disaggregated before analyses  
33 begins. Second, the use of a constant deceleration threshold is likely biased towards facilities with higher  
34 mean travel speeds. A deceleration rate of 3 m/s<sup>2</sup> is more probable when traveling at highway speeds  
35 compared to urban arterials and local streets. A lower deceleration rate of 2 m/s<sup>2</sup> may be common on  
36 highways but accurately represents evasive manoeuvres on local streets. If facility types are disaggregated,  
37 then the deceleration threshold could be set according to each specific facility type. Thirdly, hard  
38 decelerations may not be the primary evasive manoeuvre in the local and arterial streets. Other surrogate  
39 indicators should be used to capture more evasive manoeuvres in urban environments.

40 The user-centric analysis showed that more HBEs occur during the AM and PM peak periods, and  
41 that braking in PM peak period tends to be more severe. This result is logical as more congestion in these  
42 periods should yield more braking and more collisions. However, more congestion may also lead to reduced  
43 speed and therefore less severe collisions. This potential contradiction should be explored in future work.  
44 As vehicle speed increases, the severity of braking also increases. Intuitively, faster vehicles must decelerate  
45 more rapidly to avoid collisions. Although other spatio-temporal driving behaviours could not be linked to  
46 occurrence of braking events, the limited observations demonstrate that surrogate measures defined using  
47 braking as the primary evasive manoeuvre are not only spatially correlated with actual collision occurrence,  
48 but also make sense intuitively with respect to the behaviours that are related to collision occurrence and  
49 severity (traveling at peak periods and at higher speeds).

1           Limitations of this study include the use of density-based measures (heat maps). Future work should  
2 focus directly on links and/or intersections as the unit of analysis. The lack of supplemental data from  
3 accelerometers may be perceived as a limitation, though using only GPS data may be a strength of this  
4 approach, as it results in much less data that requires processing and collecting only GPS data from  
5 smartphones limits the battery requirements of the application. Additionally, correlations between HBEs  
6 and collisions was high even without this additional data. In the future, more work is needed to determine  
7 if surrogate indicators, such as over speeding or speed variation, more closely correlated with collision  
8 occurrence can be extracted from the GPS data. Additionally, further disaggregating the analyses by facility  
9 type should improve the results. Correlation at the link-level should be considered in addition to the  
10 intersection-level analysis presented above. In order to increase correlation with collision data (to use for  
11 network screening purposes), analysis could be done according to facility type (highway, primary,  
12 secondary, tertiary, local, etc.), and the threshold, smoothing parameter, and buffer size could be adjusted  
13 for each. Regardless of future improvements, hard braking events derived from the smartphone-collected  
14 GPS data of regular drivers show promising potential in the field of surrogate safety.

#### 15           **ACKNOWLEDGEMENT**

16           Funding for this project was provided in part by the Natural Sciences and Engineering Research Council.  
17           The authors recognize Charles Chung, CEO of Brisk Synergies, for his assistance in data preparation and  
18           processing.  
19



## 1 REFERENCES

1. Ellison, A. B., S. Greaves, and M. Bliemer. Examining Heterogeneity of Driver Behavior with Temporal and Spatial Factors. *Transportation Research Record*, no. 2386, 2013, pp. 158-157.
2. Algerholm, N., and H. Lahrmann. Identification of Hazardous Road Locations on the basis of Floating Car Data. *Road safety in a globalised and more sustainable world*, 2012.
3. Lee, C., B. Hellinga, and K. Ozbay. Quantifying effects of ramp metering on freeway safety. *Accident Analysis and Prevention*, no. 38, 2006, pp. 279-288.
4. Kockelman, K. M., and Y.-J. Kweon. Driver injury severity: an application of ordered probit models. *Accident Analysis and Prevention*, Vol. 34, 2002, pp. 313-321.
5. Bagdadi, O. Assessing safety critical braking events in naturalistic driving studies. *Transportation Research Part F*, no. 16, pp. 117-126.
6. Wu, K.-F., and P. P. Jovanis. Defining and screening crash surrogate events using naturalistic driving data. *Accident Analysis and Prevention*, no. 61, 2013, pp. 10-22.
7. Tarko, A., G. Davis, N. Saunier, T. Sayed, and S. Washington. Surrogate Measures of Safety. Transportation Research Board, 2009.
8. Laureshyn, A., K. Astrom, and K. Brundell-Freij. From Speed Profile Data to Analysis of Behaviour. *IATSS Research*, Vol. 33, no. 2, 2009, pp. 88-98.
9. Bahler, S. J., J. M. Kranig, and E. D. Minge. Field Test of Nonintrusive Traffic Detection Technologies. *Transportation Research Record*, no. 1643, 1998, pp. 161-170.
10. Oh, C., J.-s. Oh, and S. G. Ritchie. Real-time estimation of Freeway Accident Likelihood. in *Transportation Research Board Annual Meeting*, Washington, D.C., 2001.
11. Golob, T. F., W. W. Recker, and V. M. Alvarez. Freeway safety as a function of traffic flow. *Accident Analysis and Prevention*, no. 36, 2004, pp. 933-946.
12. Lee, C., F. Saccomanno, and B. Hellinga. Analysis of Crash Precursors on Instrumented Freeways. *Transportation Research Record: Journal of the Transportation Research Board*, no. 1784, 2002, pp. 1-8.
13. Herrera, J. C., D. B. Work, R. Herring, X. Ban, Q. Jacobson, and A. M. Bayen. Evaluation of traffic data obtained via GPS-enabled mobile phones: The Mobile Century field experiment. *Transportation Research Part C*, no. 18, 2010, pp. 568-583.
14. El Faouzi, N.-E., H. Leung, and A. Kurian. Data fusion in intelligent transportation systems: Progress and challenges – A survey. *Information Fusion*, no. 12, 2011, pp. 4-19.
15. Jun, J., J. Ogle, and R. Guensler. Relationships between Crash Involvement and Temporal-Spatial Driving Behavior Activity Patterns Using GPS Instrumented Vehicle Data. in *Transportation Research Board Annual Meeting*, Washington, DC, 2007.
16. Dia, H., and K. Thomas. Development and evaluation of arterial incident detection models using fusion of simulated probe vehicle and loop detector data. *Information Fusion*, no. 12, 2011, pp. 20-27.
17. Sethi, V., N. Bhandari, F. S. Koppelman, and J. L. Schofer. Arterial Incident Detection using Fixed Detector and Probe Vehicle Data. *Transportation Research Part C*, Vol. 3, no. 2, 1995, pp. 99-112.

18. Shen, W., and L. Wynter. Real-Time Road Traffic Fusion and Prediction with GPS and Fixed-Sensor Data. in *15th International Conference on Information Fusion*, Singapore, 2012, pp. 1468-1475.
19. Fazeen, M., B. Gozick, R. Dantu, M. Bhukhiya, and M. C. Gonzalez. Safe Driving Using Mobile Phones. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 13, no. 3, 2012, pp. 1462-1468.
20. Ellison, A. B., S. P. Greaves, and M. C. Bliemer. Driver behaviour profiles for road safety analysis. *Accident Analysis and Prevention*, no. 76, 2015, pp. 118-132.
21. Dingus, T. A., S. G. Klauer, V. L. Neale, A. Petersen, S. E. Lee, J. Sudweeks, M. A. Perez, J. Hankey, D. Ramsey, S. Gupta, C. Bucher, Z. R. Doerzaph, J. Jermeland, and R. R. Knipling. The 100-Car Naturalistic Driving Study, Phase II – Results of the 100-Car Field Experiment. NHTSA, Washington, DC, DOT HS 810 593, 2006.
22. City of Quebec. Mon Trajet. *City of Quebec*, [http://www.ville.quebec.qc.ca/citoyens/deplacements/mon\\_trajet.aspx](http://www.ville.quebec.qc.ca/citoyens/deplacements/mon_trajet.aspx). Accessed May 13, 2015.
23. Miranda-Moreno, L. F., C. Chung, D. Amyot, and H. Chapon. A system for collecting and mapping traffic congestion in a network using GPS smartphones from regular drivers. in *Transportation Research Board Conference Processings*, Washington, DC, 2014.
24. Brisk Synergies. *Brisk Synergies*, <http://www.brisksynergies.com/>. Accessed July 22, 2015.
25. Marchal, F. TrackMatching. 2015. <https://mapmatching.3scale.net/>. Accessed May 1, 2015.
26. OpenStreetMap. About. *OpenStreetMap*, 2015. <http://www.openstreetmap.org/about>. Accessed May 11, 2015.
27. Rakha, H., F. Dion, and H.-G. Sin. Using Global Positioning System Data for Field Evaluation of Energy and Emission Impact of Traffic Flow Improvement Projects. *Transportation Research Record*, no. 1768, 2006, pp. 210-223.