

**Measuring Congestion Using Large-Scale Smartphone-Collected GPS Data in an Urban Road Network**

**Joshua Stipancic**, PhD Candidate, McGill University  
**Luis Miranda-Moreno**, Associate Professor, McGill University  
**Aurelie Labbe**, Associate Professor, McGill University

Paper prepared for presentation  
at the Big Data Applications for Travel Demand Management Session

of the 2016 Conference of the  
Transportation Association of Canada  
Toronto, Ontario

## **ABSTRACT**

Although travel time is an important performance measure for transportation professionals, congestion is perhaps a more important consideration for road users. Congestion is a dynamic phenomenon with variation across both space and time making it a promising application of smartphone-collected GPS data. The purpose of this study is to utilize GPS data collected using a smartphone application and regular drivers to estimate congestion at both the macroscopic and microscopic level across an urban road network. Data is collected using the “Mon Trajet” smartphone application in Quebec City, Canada. This data consists of nearly 50,000 trips collected during 3 weeks from approximately 5000 drivers. The application allowed for the collection of a large number of trips from regular drivers using a system that minimally impacts them or their behaviour. Given the large spatio-temporal dimension, several data issues are identified and corrected using the presented methodology. First, position of the GPS traces is provided in terms of a latitude and longitude and is not linked spatially to the road network. TrackMatching is a commercially available map-matching service used to match GPS data to the OpenStreetMap road network. Speeds are smoothed, and level of congestion is computed using the Congestion Index (CI). CI is computed at the individual link level at time intervals of 1 hour, to yield a detailed picture of congestion both across time and space. Finally, the progression of congestion over time is mapped across the entire road network and congestion variability across different time scales is computed.

## INTRODUCTION

Maintaining adequate levels of service within urban road networks requires accurate and quantitative measures of system performance. Performance measures are necessary for transportation professionals to operate existing transportation networks and plan future facilities (1). For the driving population, network performance influences travel choices (2). Perhaps the most common performance measure, link or route travel time, continues to grow in importance for both road users and practitioners (3). For traffic operators, “travel time is one of the most important measures for evaluating the performance of traffic networks” (2) and is a key parameter defining traffic state (4). For road users, travel time is easily understood (2). Although travel time is easily communicable, it may not be the biggest factor influencing travel decisions. Congestion, or travel time reliability, may be more important to road users and may have a greater impact on travel decisions (1). Road users may prefer a longer route over a shorter route if their travel time on the longer route is more reliable. In other words, people are willing to accept longer travel times if they can be assured that they will usually arrive on time (5). Congestion is a concern to the urban community who depend on the reliability of goods and services, and congestion is related to personal and community wellbeing within the urban environment (1).

Detailed road network congestion information would be beneficial for both transportation professionals and the population at large. Recently, advances in technology and management practices have “increased [the] need for very accurate road traffic information” (6) and have assisted in the development of traffic sensors, including radar, magnetic, and video-based devices. However, because congestion “is a dynamic phenomenon with elements of both space and time” (1), it is a promising application of probe vehicle data. Probe vehicles act “as moving sensors, continuously feeding information about traffic conditions” (6) through continuous instrumentation and tracking (3). Probe vehicles allow for a precise measurement of origin-destination (or route) travel time for a vehicle operating within normal traffic (3). Although several methods for instrumenting vehicles exist, GPS data, collected either by dedicated devices or GPS-enabled smartphones, is a reliable source (7). While probe vehicle studies have traditionally been limited in driver sample size and spatio-temporal coverage due to the labour cost associated with operating probe vehicles (2), the proliferation of GPS-enabled smartphones has the potential to increase the number of drivers sampled, increase temporal coverage to several weeks or months, and increase the spatial coverage to include the entire road network.

Although GPS probe vehicles have been successfully implemented in freeway environments (8), specific consideration for urban environments is required due to the more variable and interrupted traffic flows caused by signalization or geometry (2). Tall buildings in urban centers can completely block GPS signals or create spurious signals through multipathing (1). These issues translate into inaccuracies of the GPS locations with respect to the road network which must be corrected through map-matching. The low frequency or inconsistency of GPS trip data in some links in the network, particularly in the off peak hours, is due to the low penetration rate of the smartphone application. Another issue is how to represent congestion variability across time and space and choosing the most appropriate measures of congestion. Finally, as the application provides large volume of data, methods for automating data analysis and processing are required. Despite increasing congestion levels, smartphone-based systems for measuring and monitoring traffic congestion are still rare in North American cities (9). The purpose of this paper is to present a methodology for computing congestion measures (accounting for spatial and temporal variations of travel time) using GPS data from regular drivers collected through a smartphone application. The three primary objectives are; to process network-wide GPS travel data collected; to quantify congestion on the network scale using the Congestion Index (CI); and, to observe and quantify the hourly propagation and daily variability of congestion.

## LITERATURE REVIEW

Several methods for estimating road traffic state (travel time or congestion) have been explored in the existing literature, and depend predominantly on the type of data that is collected. Using fixed point traffic sensors, the naïve method uses average speeds at a specific point that, when combined with flow, density, and speed relationships, is used estimate traffic conditions (10). The naïve method has been criticized for systematic bias (11) as detector data “only reflect conditions averaged over a fixed time period at a single point in space” whereas link travel time “reflects traffic conditions averaged over a fixed distance and a variable amount of time” (12). In trajectory methods, trajectories of simulated vehicles are constructed based on traffic data observed by several consecutive fixed sensors (13). van Lint and van der Zijpp (13) improved traditional methods by assuming linear speed variation (rather than piecewise-constant variation), which more accurately represents spatially and temporally dependent variation in flow. Coifman (12) constructed trajectories based on several loop detectors speeds to estimate travel time in a freeway environment, demonstrating that estimated travel times were within 10% of actual travel times on average. Liu and Ma (2) fused loop data with signal phase information in urban corridors to estimate travel times generally within 5% of actual. As with the naïve method, trajectory methods are limited because data is collected from a single point in space, and “changes in the traffic stream may be overrepresented or underrepresented” (12).

Vehicle reidentification (VRI) is “the process of matching vehicles from one point on the roadway (one field of view) to the next” (14) based on a ‘reproducible feature’ or vehicle signature (15). When a vehicle is identified at two locations within the network, the travel time between those locations is determined. Vehicle signatures may be captured using license plate recognition (16) or media access control addresses captured from Bluetooth devices within passing vehicles (17), though most Bluetooth detectors have a detection rate of 5% or less (18). Vehicle length (19) and magnetic signature (20) have also been used to define vehicle signature, as presented by Coifman and Cassidy (19) who were able to reidentify 20% of vehicles based on length and Sun et al. (14) who used inductive loops and feature-based colour extracted from video stills to achieve an approximately 90% match rate. Kwong et al. (20) presented a system for VRI based on permanent wireless magnetic sensors installed along an urban corridor, spanning several intersections. The authors estimate a successful matching rate of 65-75% (20).

Due to the issues and assumptions associated with the above techniques, probe vehicles have become a popular method for measuring traffic conditions (8). Continuous tracking over time and space represents a substantial improvement over methods of travel time estimation using fixed sensors. Initial work by D’Este, Zito, and Taylor (1) explored the feasibility of using GPS to collect traffic data and concluded that GPS was “a relatively cheap, efficient and effective means” of collecting traffic data. Techniques for estimating travel time vary depending on the number of probe vehicles utilized. The traditional approach uses the average travel time from a relatively large number of probe vehicles operating within the same time and space. However, the number of probe vehicles is dependent on traffic flow, and the labour requirement associated with a large number of probe vehicles is high (9). Due to these limitations, approaches have been developed to use relatively fewer probe vehicles with statistical adjustments to extrapolate probe vehicle travel time to mean travel time (21). Li and McDonald (3) proposed an approach using only a single probe vehicle, using the driving pattern of the probe vehicle to estimate the difference between the probe vehicle and average traffic conditions (3). Smartphone-collected GPS data enables the use of a large number of probe vehicles without the high labour costs associated with traditional instrumentation, and data coming from regular drivers may better represent typical traffic conditions. These benefits were demonstrated in an earlier study by Stipanovic, Miranda-Moreno, and Saunier (22).

Measures of congestion are typically based on either travel time or speed (23). Congestion may be measured by a change in travel time (24) from a baseline or expected travel time measurement (12). Measures like travel time index (TTI) use the ratio of peak travel time to off-peak travel time to determine congestion at the link level (25). TTI and similar techniques can incorporate historical trend data to separate recurring and non-recurring congestion (12). Skabardonis, Varaiya, and Petty (26) utilized a delay-based approach to separate recurrent from non-recurrent delay. Delay can be measured by comparing the average speed to a free flow speed, where a reduction in speed introduces delay. In Washington State, mean speeds 75% of free flow speed are considered to signify the onset of congestion (25). The Congestion Index (CI) was proposed by Dias et al. (27) as a ratio of actual speed to free flow speed. With GPS probe vehicles, congestion measures based on speed are preferred. Travel time estimation may be influenced by errors in the reported GPS coordinates and on assumptions of the start and end of trip (or link). Considering the relatively recent advent of GPS data in transportation research, several shortcomings remain in the literature. Few studies have considered the rich source of data available from GPS-enabled smartphones. Congestion studies using probe vehicles have primarily focussed on the corridor-level without consideration for estimating travel time or congestion at the network level. Despite successful probe vehicle studies in freeways, additional focus on urban roadways is required.

## METHODOLOGY

### Data Structure

GPS data from the smartphones of regular drivers contains observations describing the entirety of their trip both across time and across the road network. For each trip,  $i$ , logged into a smartphone application, GPS travel data is returned as a series of observations,  $O_{it}$ , such as

$$trip_i = \begin{Bmatrix} O_{i0} \\ O_{i1} \\ \vdots \\ O_{it} \\ \vdots \\ O_{in} \end{Bmatrix} = \begin{Bmatrix} i, c_{i0}, dt_{i0}, x_{i0}, y_{i0}, z_{i0}, v_{i0} \\ i, c_{i1}, dt_{i1}, x_{i1}, y_{i1}, z_{i1}, v_{i1} \\ \vdots \\ i, c_{it}, dt_{it}, x_{it}, y_{it}, z_{it}, v_{it} \\ \vdots \\ i, c_{in}, dt_{in}, x_{in}, y_{in}, z_{in}, v_{in} \end{Bmatrix}$$

where  $i$  is a unique trip identifier,  $O_{it}$  is the an observation in trip  $i$  at time  $t$ ,  $c_{it}$  is a unique coordinate identifier,  $dt_{it}$  is the datetime,  $x_{it}$ ,  $y_{it}$ , and  $z_{it}$  are the latitude, longitude, and altitude, and  $v_{it}$  is the speed. From each trip, several key pieces of trip information include the origin  $(x_{i0}, y_{i0})$  and destination  $(x_{in}, y_{in})$  and start  $(dt_{i0})$  and end times  $(dt_{in})$ . Total travel time can also be computed  $(dt_{in} - dt_{i0})$ . The time between consecutive observations,  $\Delta t$ , is typically between 1 and 2 seconds. Depending on the application used to collect the data, socio-demographic information may also be available. Once a trip has been collected and reported by the user, initial pre-processing of the data using methods including Kalman filtering (28) to reduce variability are typical. The data is then stored in a database from which observations are exported for analysis.

### Map Matching

Although the raw GPS data from a smartphone application is rich in spatio-temporal data, position is provided only in terms of latitude and longitude and is not linked spatially to the road network. Additionally, location variability is expected in the raw data. If the goal is to determine congestion at the link level, then it is necessary to explicitly match each trip to the travelled network links. This process, known as ‘map matching’ ensures that traffic conditions extracted from the trip data are correctly assigned to the links in which the traffic conditions are occurring. TrackMatching is

a commercially available, cloud-based web map-matching software service (29) that matches GPS trip data to the OpenStreetMap (OSM) road network (30). Before GPS data is sent to TrackMatching, the data is into individual trips and formatted according to the software input requirements, including only the coordinate ID, timestamp, latitude, and longitude for each observation. The software returns a new latitude and longitude,  $x'_{it}$  and  $y'_{it}$ , which correspond to a specific OSM link ID,  $l_{it}$ , as shown below.

$$\{c_{it}, dt_{it}, x_{it}, y_{it}\} \rightarrow \text{TrackMatching} \rightarrow \{c_{it}, dt_{it}, x'_{it}, y'_{it}, l_{it}, s_{it}, d_{it}\}$$

$x'_t$  and  $y'_t$  are chosen based on the Euclidean distance from the raw GPS points to the nearest link and on network topology (31). Track Matching also returns the source,  $s_{it}$ , and destination nodes,  $d_{it}$ , which can be used to identify direction of travel along the link. The algorithm generates a set of candidate paths and assigns the trip to the most probable path from origin to destination. After map-matching is completed, each observation corresponds to an exact location within the road network, and the series of links can be used to define the route from origin to destination. However, because important information including speed and datetime is lost d map matching process, the results are merged back with the original data to preserve the complete data set as shown below. The processes of data collection and map matching are illustrated in Figure 1.

$$trip_i = \left\{ \begin{array}{l} i, c_{i0}, dt_{i0}, x'_{i0}, y'_{i0}, z_{i0}, v_{i0}, l_{i0}, s_{i0}, d_{i0} \\ i, c_{i1}, dt_{i1}, x'_{i1}, y'_{i1}, z_{i1}, v_{i1}, l_{i1}, s_{i1}, d_{i1} \\ \vdots \\ i, c_{it}, dt_{it}, x'_{it}, y'_{it}, z_{it}, v_{it}, l_{it}, s_{it}, d_{it} \\ \vdots \\ i, c_{in}, dt_{in}, x'_{in}, y'_{in}, z_{in}, v_{in}, l_{in}, s_{in}, d_{in} \end{array} \right\}$$

## Network Definition

Although the map-matching procedure links each observation to the road network, the use of the OSM network in the TrackMatching algorithm presents a challenge. Ideally, links in the road network would connect two nodes at adjacent intersections. However, because the OSM road network is generated ad-hoc by users, a single link may connect several consecutive intersections. In urban centers, intersection design and operation can significantly impact congestion levels on consecutive links. It is desired to redefine the network such that each link is properly defined between adjacent intersections. Redefining the network requires several steps, which can be completed in any GIS software environment. The process is as follows:

1. Identify all nodes that represent an intersection in the road network. In doing so, nodes that only define network topology are ignored.
2. Split the road network at the identified nodes. Any links connected more than two intersections are broken into several smaller links. Links already properly defined are unchanged.
3. Rename each link according to its original ID and the nodes on either end of the link. Step 2 leaves several links with the same ID. In order for each link to have a unique identifier, the nodes on either end of the link are used to provide a unique ID.
4. Remap the GPS observations to the new network. Travelled links in the GPS trip data are renamed using the same scheme as the mapping data, by concatenating the link ID, source node, and destination node into a unique identifier.

The results of this process are shown in Figure 2.

## Computing Congestion Index

The map matching procedure enables congestion measurement for every link containing sufficient GPS data, providing either a microscopic view of link performance, or a macroscopic view of network performance. As discussed, several measures of congestion based on travel time have been proposed. However, because link travel time is dependent on position, and because the precise latitude and longitude are untrustworthy (and are in fact removed as part of the map matching procedure), a congestion measure based on speed (which is directly available from the GPS data) is preferred. Dias et al. (27) proposed the Congestion Index (CI) as one speed-based congestion measure, calculated as

$$CI = \begin{cases} \frac{\text{free flow speed} - \text{actual speed}}{\text{free flow speed}} & \text{if } CI > 0 \\ = 0 & \text{if } CI \leq 0 \end{cases} \quad (1)$$

This formulation yields CI values ranging between 0 and 1, where 0 is completely uncongested and 1 is completely congested. The first necessary step is calculating the free flow speed on each link,  $L$ . Free flow speed has been defined in numerous ways, though as congestion is generally constrained to the AM and PM peak periods, the speeds observed outside of these times can be used to estimate the free flow speed. For the purpose of this project, the morning peak period was defined as 6:00 to 10:00 AM, and the evening peak from 3:00 to 7:00 PM. The off-peak time,  $T_{off}$ , includes all times outside of these peak periods. Free flow speed on a given link,  $L$ , is calculated as the average of all observed speeds on  $L$  during  $T_{off}$ , or

$$FFS_L = \frac{\sum_i \sum_t v_{it}}{N} \quad \text{if } l_{it} = L \text{ and } t \in T_{off} \quad (2)$$

where  $v_{it}$  is the speed for every observation on link  $L$  during  $T_{off}$ , and  $N$  is the count of those observations. Next, the congestion index for every observation can be computed according to

$$CI_{it} = \begin{cases} \frac{FFS_{l_{it}} - v_{it}}{FFS_{l_{it}}} & \text{if } CI_{it} > 0 \\ = 0 & \text{if } CI_{it} \leq 0 \end{cases} \quad (3)$$

where  $CI_{it}$  is the congestion index for observation  $O_{it}$ ,  $FFS_{l_{it}}$  is the free flow speed on link  $l_{it}$ , and  $v_{it}$  is the observed speed. As congestion levels vary across both distance and time, it is not only necessary to calculate CI at the link level, but also to calculate CI at different time intervals. The peak periods were divided into 60 minute time periods (one per hour) resulting in 8 total time periods. Therefore, the congestion index for link  $L$  during time period  $T$  is calculated as:

$$CI_{LT} = \frac{\sum_i \sum_t CI_{it}}{N} \quad \text{if } l_{it} = L \text{ and } t \in T \quad (4)$$

where  $CI_{it}$  is the congestion index for every observation on link  $L$  during a time period  $T$ , and  $N$  is the count of those observations. To minimize noise, filters were used to set minimum acceptable numbers of trips and observations for CI calculation. For a valid  $CI_{LT}$ ,  $L$  must contain at least 2 trips during time  $T$ , and each of those trips must have at least 2 observations falling on link  $L$ .

## Data Analysis

The final methodological step is to analyze the CI data and report useful metrics. Most simply, it is possible to visualize CI at the link level for any given time interval on any given day using GIS software. Although this type of analysis can certainly shed light on a particular instant in time, a single snapshot of the network, or even several snapshots of the network, are not able to provide general insight or conclusions which would be beneficial to transportation professionals or to the driving public. Because congestion levels vary significantly throughout the day (due to variation in demand) and vary significantly between days (due to variation in demand, non-recurrent incidents, and random variation), it is necessary to quantify and/or visualize that variation, as well as the overall magnitude of congestion.

The first objective was to view the hourly propagation of congestion at the macroscopic scale, for the network as a whole (in general, the hourly variation could also be determined for a specific corridor or link). Congestion does not occur all at once. Instead, it gradually builds and then subsides throughout the peak periods. Similarly, congestion does not occur across all network links simultaneously. In fact, congestion is generated at specific locations in the network at specific times (as trips are generated) and then propagates through the network over time (as those generated trips move from origin to destination). In largely monocentric cities, it is intuitive that the onset of congestion would begin furthest from the city center (where most residents live) earliest in the AM peak period. As time progresses, this congestion generated at the city outskirts should propagate towards the city center (where most residents work). Showing the progression of congestion at the network scale would provide insight into how the network behaves in general. In order to do this, links are binned according to distance from the city center in 1 km increments. Data from all available weekdays is pooled, and  $CI_{LT}$  is computed for each link and each time interval. The average CI within in each distance bin is calculated and plotted. Trend lines are fit, and observations are made based on the macroscopic congestion patterns across the network.

The second objective was to observe the daily variation in traffic congestion at the microscopic, or link level. A detailed understanding of congestion should include not only the average level of congestion for a given link, but should also include a measure of how variable that level of congestion is from day to day. As road users may prefer routes with longer yet consistent travel times, targeting congestion variability may be as important as targeting overall congestion levels. To compute both mean CI and variance in CI, a single hour is chosen for analysis. One  $CI_{LT}$  value is then computed for every available weekday of data for every link. If enough daily observations exist, the mean and variance of  $CI_{LT}$  are computed. Maps can be generated which simultaneously show mean congestion levels (by colour) and variance in congestion (by line thickness). This type of map represents a substantial improvement over a single snapshot by considering temporal (daily) variation in the level of congestion. The most problematic locations in the network can then be observed as those locations with consistently high CI, and highly variable CI.

## DATA DESCRIPTION

GPS travel data was collected in Quebec City, Canada using the Mon Trajet application (32) developed by BriskSynergies(33) . Screenshots from the application are shown in Figure 3. The application, which is available for Apple and Android devices, was installed voluntarily by drivers and allowed them to anonymously log trips into the application. As part of the system developed by BriskSynergies, data is automatically uploaded and stored in a cloud-based platform. In total, approximately 5000 driver participants have logged nearly 50,000 trips using the application. Although this data can be retrieved directly from the platform, the data used in this study is a large sample of the open-source data made available by the City of Quebec. The sample for this study



contained 2413 drivers and 12,724 individual trips during the period between April 28 and May 18, 2014. Over the 21 days sampled, 19.7 million individual data points were logged.

## RESULTS

### Data Processing and Visualization

As stated previously, one substantial limitation is that CI can only be calculated for links that have enough data. Even though the collection campaign in Quebec City was large (almost 13,000 trips were logged through the application) there were no trips on most of the network links. This is especially true of residential links, where data was only available on streets where participants lived. Still, a majority of the major freeways, arterials, and collectors have decent data coverage. These facilities, in fact, are the ones where congestion is the greatest concern, and so meaningful analysis is possible despite the missing data on residential streets.

Data was collected over three complete weeks, resulting in 15 total weekdays of data for analysis. Considering these weekdays in general, each day yields at least one CI measurement for between 2000 and 4000 links. Considering that the Quebec City network has over 50,000 links in total, the data represents about 6% of the network on an average day. For each time interval during the peak periods, there is between 250 and 1750 links for which CI is calculated. This represents between 0.5% and 3.5% of the total road network. As with the total number of expected trips, the total number of GPS trips logged varies with time, first growing to a maximum and then subsiding throughout the peak periods. If the data from all weekdays is pooled together, the results are improved significantly. At least one CI measurement exists for 16,805 links (or 34% of the total road network) and each hour contains between 6600 and 12,000 links with a valid CI (13-24% of the total). CI can be mapped and visualized using any GIS software. A CI map for May 6 between 8:00 and 9:00 AM is presented in Figure 4.

### Hourly Propagation

In order to view hourly propagation of congestion, variation through both time and space must be considered. Each 4-hour peak period can be viewed as having an onset period (lasting one hour), the peak itself (lasting two hours), and a dissipation period (lasting one hour). Plots of CI and distance are provided in Figure 5 for both the AM and PM peak periods. As weekday data was pooled, the congestion profile for each hour is based on 13%-24% of the total links in the network. Starting with Figure 5a, the onset of congestion in the AM peak (6:00 to 7:00 AM) is characterized by relatively consistent CI levels across the network. Peaks in the profile during this time are related to those distances which contain major highways or arterials. From 6:00 AM to 7:00 AM, CI increases across all distances, although the increase is relatively greater as distance to the city center decreases (ranging from CI of 0.12 at the center to 0.08 at a distance of 20 km). Levels of congestion remain relatively consistent between 7:00 and 9:00 AM. During the dissipation period (9:00 to 10:00 AM) the relative profile is the same (with higher congestion in the city center), however, congestion far from the city center has returned levels at or below that of the onset period (CI is approximately 0.1 at the center, but 0.04 at a distance of 20 km). This shows that congestion at the outskirts of the city dissipates earlier than at the center.

Results for the PM peak period, in Figure 5b, are essentially a 'mirror image' of the AM peak. At the onset of congestion (3:00 to 4:00 PM) congestion is relatively high in the city center, but is at or below the baseline CI during the dissipation period (6:00 to 7:00 PM) at the outskirts of the city. During the middle of the peak period (4:00 to 6:00 PM) congestion increases, and is highest in the city center (CI of about 0.13) and lowest at a distance of 20 km (CI of 0.8). During the

dissipation period, congestion is, on average, consistent with distance (although it increases slightly as distance increases) at a CI of between 0.06 and 0.07.

Several general observations were made based on these plots. First, overall congestion levels are nearly equal between the AM and PM peak periods. The baseline CI (during the AM onset and PM dissipation periods) is generally consistent around 0.06, and the maximum CI at the city center was 0.12 in the AM and 0.14 in the PM. In general, CI varies much less at the outskirts than at the city center. In the AM peak, CI returns to or below onset levels during the dissipation period. In the PM, CI begins at or below dissipation levels during the onset period. These plots can help describe how congestion propagates from the outskirts to the city center in the morning, and back from the center to the outskirts in the evening. A simplified description of congestion formation and propagation is presented in Table 1. Obviously, this is an oversimplification of the congestion phenomena. However, the fact that the measure of CI, and the method of computing CI across the network using data from GPS enabled smartphones, is consistent with intuitive characteristics of congestion at the macroscopic level, is extremely promising.

### Daily Variation

In order to understand daily variations in the level of congestion at the microscopic level,  $CI_{LT}$ , was computed for each link, during a single hour (8:00 to 9:00 AM) on each of the 15 weekdays of data. The mean and variance of  $CI_{LT}$  were then computed. In total, there were only 1019 links with enough data during the hour of analysis on all 15 days for which a mean and variance could be computed. It is suspected that some correlation exists between the mean and variance of CI, as it is unlikely that links with low levels of congestion would have much daily variation. In fact, for the 1019, there was a correlation of 0.60 between the variance and the mean. As stated earlier, it would be beneficial to target not only the most congested links or locations, but also to target the locations with the greatest variance in order to improve travel time reliability for the road user. The mean and variance of CI are mapped for the network in Figure 6, where the colour represents the mean CI, and the link thickness represents the variance in CI. Based only on this map, the most troublesome locations in the network can be determined visually. The most critical locations in the network include:

- Southbound and eastbound ramps of the Autoroute Henri-IV/Autoroute Felix-Leclerc Interchange
- Southbound and westbound ramps of the Autoroute Laurentienne/Autoroute Felix-Leclerc Interchange
- Eastbound ramp of the Autoroute Henri-IV and Autoroute Duplessis Interchange
- Eastbound direction of major downtown arterials including Grande Allee O, Boulevard Rene-Levesque O, Chemin Ste-Foy, and Boulevard Charest O

These results are again consistent with intuition about congestion at the microscopic level. The most congested locations in the AM peak are the freeways, ramps, and arterials which bring traffic downtown, and these locations are identified by both the mean and variance in congestion level.

### CONCLUSIONS

The purpose of this paper was to propose measures for representing congestion levels across time and space in an urban road network (Quebec City, Canada) using data collected from the GPS-enabled smartphones of regular drivers. This paper first presented the methodology for processing the GPS data and computing CI. Through map matching and network definition,

observations are explicitly related to properly defined links in the road network. The measure and method for evaluating congestion proved to be relatively easily to compute, and the data analysis showed that results were consistent with the expected behaviour of congestion at both the microscopic and macroscopic levels. Despite some limited spatio-temporal data coverage, enough data was available to calculate and visualize congestion for the majority of major freeways, arterials, and collectors within Quebec City.

When considering hourly propagation of congestion during an average day, the results were consistent with intuitive understanding of how congestion builds and dissipates over time and how congestion propagates across distance between the outskirts of the city and the city center at the macroscopic scale. This analysis clearly showed how CI tends to be both greater and more variable as distance to the city center decreases. In both the AM and PM peak periods, links closest to the city center exhibited, on average, a CI of between 0.12 and 0.14, while links furthest from the city center had CI general below 0.09. Average baseline CI during the AM onset and PM dissipation periods were around 0.06 network wide.

Daily average and variance in CI were computed for a single hour of analysis in the AM peak periods for each link in the network. Although many of the most congested links are also the most variable, the correlation between mean and variance of CI was only 0.06. In other words, there are highly congested links with less variability and more variable links with relatively low levels of congestion. Links with high congestion and low variability represent chronic problems in the network, while links with low congestion and high variability could point towards non-recurrent phenomena (collisions or construction) taking place over the data collection period. Not surprisingly, the most problematic areas were freeway interchanges and arterials which bring motorists east and south into the downtown area of Quebec City. This type of analysis is potentially beneficial in the prioritizing of sites for congestion remediation. Although it may be in the transportation professional's interest to remediate the most congested locations, a greater benefit may be provided to motorists by targeting the sites with the most variability.

Although the methodology and proof of concept were shown to be successful, several items are planned for future research. First, the OSM data is incomplete in some key areas of the network. This is again partially due to the ad-hoc nature of the OSM data. Methods for finding and completing the map itself are required to complete these key corridors. In terms of measuring congestion, although it is acceptable for some links to be without data (specifically residential streets without trip data to process), there are several isolated links in the network without data despite links before and after having data. Methods for filling in this missing data (based on both spatial correlation with other links and temporal correlation with other time periods) would provide a benefit for work in this area. Finally, a greater depth of analysis of the computed CI data is required if this type of work is to be applicable in practice for network planning or congestion remediation. Finally, a software platform for automating the entire process can be built to make this an accessible and practical tool for city planners.

## **ACKNOWLEDGEMENT**

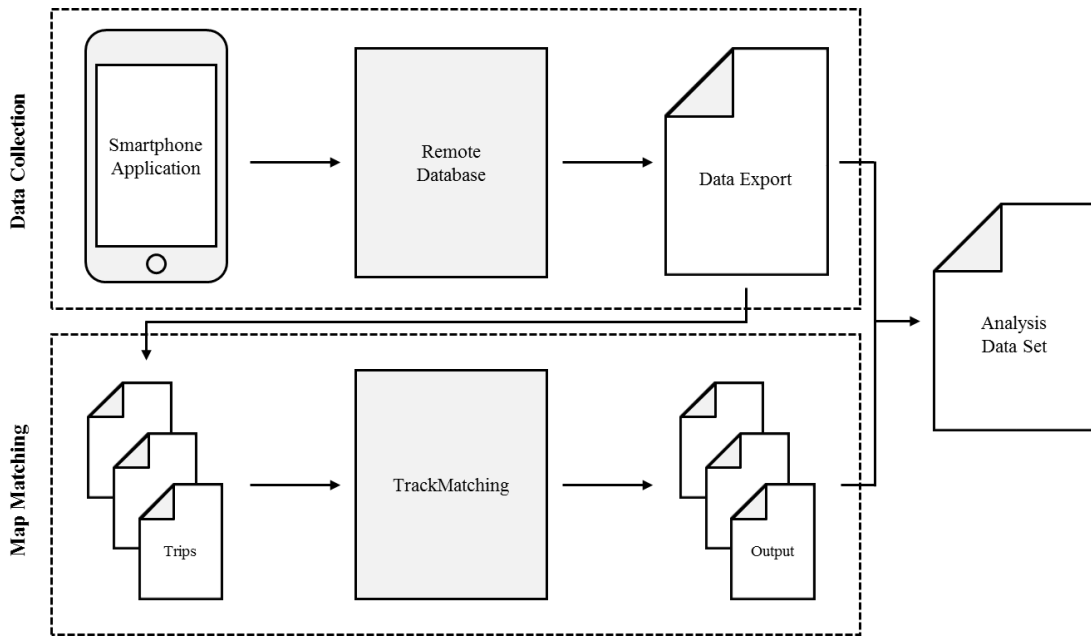
Funding for this project was provided in part by the Natural Sciences and Engineering Research Council. The authors would like to thank the City of Quebec and BriskSynergies for providing the GPS data. The authors recognize Spencer McNee for his assistance in data preparation and processing.

## REFERENCES

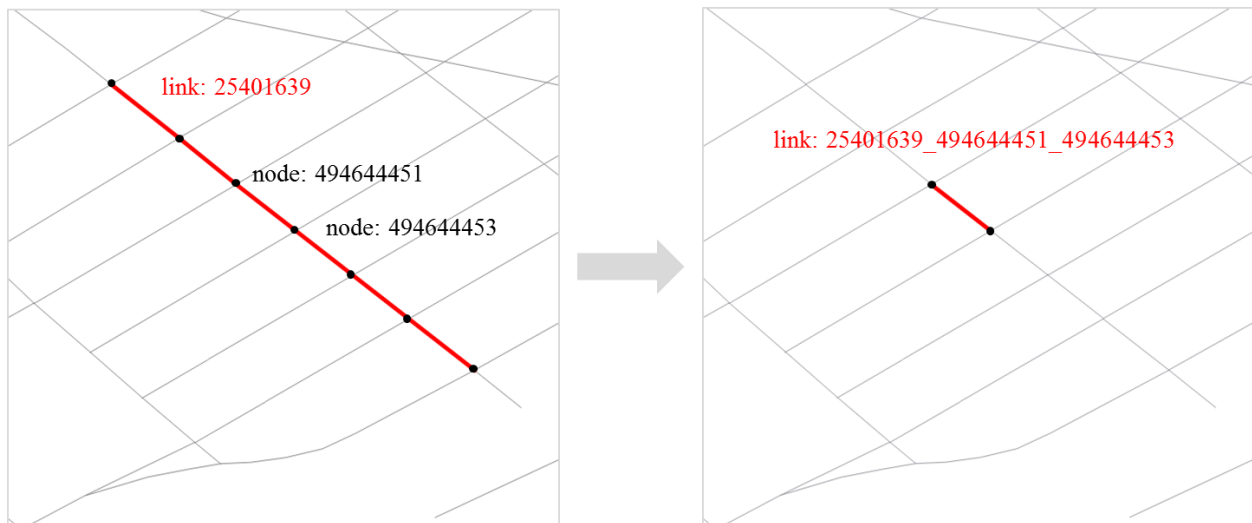
1. D'Este, G. M., R. Zito, and M. A. Tayler. Using GPS to Measure Traffic System Performance. *Computer-Aided Civil and Infrastructure Engineering*, no. 14, 1999, pp. 255-265.
2. Liu, H. X., and W. Ma. A virtual vehicle probe model for time-dependent travel time estimation on signalized arterials. *Transportation Research Part C*, no. 17, 2009, pp. 11-26.
3. Li, Y., and M. McDonald. Link Travel Time Estimation Using Single GPS Equipped Probe Vehicle. in *The IEEE 5m international Conference on Intelligent Transportation Systems*, Singapore, 2002, pp. 932-937.
4. Zhang, H. M. Link-Journey-Speed Model for Arterial Traffic. *Transportation Research Record*, no. 1676, 1999, pp. 109-115.
5. Taylor, M. A.P. Travel through time: the story of research on travel time reliability. *Transportmetrica B: Transport Dynamics*, Vol. 1, no. 3, 2013, pp. 174-194.
6. 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.
7. 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.
8. Quiroga, C. A., and D. Bullock. Travel time studies with global positioning and geographic information systems: an integrated methodology. *Transportation Research Part C*, no. 6, 1998, pp. 101-127.
9. Chen, M., and S. I. Chien. Determining the Number of Probe Vehicles for Freeway Travel Time Estimation by Microscopic Simulation. *Transportation Research Record*, no. 1719, 2000, pp. 61-68.
10. Dailey, D.J. Travel-time estimation using cross-correlation techniques. *Transportation Research Part B*, Vol. 27B, no. 2, 1993, pp. 97-107.
11. Ostrand, M., K. F. Petty, P. Bickel, J. Jiang, J. Rice, Y. Ritov, and F. Schoenberg. Simple Travel Time Estimation from Single-Trap Loop Detectors. *Intellimotion*, Vol. 6, no. 2, 1997, pp. 4-5, 11.
12. Coifman, B. Evaluating travel times and vehicle trajectories on freeways using dual loop detectors. *Transportation Research Part A*, no. 36, 2002, pp. 351-364.
13. van Lint, J., and N. van der Zijpp. Improving a Travel-Time Estimation Algorithm by Using Dual Loop Detectors. *Transportation Research Record*, no. 1855, 2007, pp. 41-48.
14. Sun, C. C., G. S. Arr, R. P. Ramachandran, and S. G. Ritchie. Vehicle Reidentification Using Multidetector Fusion. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 5, no. 3, 2004, pp. 155-164.
15. Coifman, B., and S. Krishnamurthy. Vehicle reidentification and travel time measurement across freeway junctions using the existing detector infrastructure. *Transportation Research Part C*, no. 15, 2007, pp. 135-153.

16. Anagnostopoulos, C., T. Alexandropoulos, V. Loumos, and E. Kayafas. Intelligent traffic management through MPEG-7 vehicle flow surveillance. in *EEE John Vincent Atanasoff 2006 International Symposium on Modern Computing* , 2006.
17. Haseman, R. J., J. S. Wasson, and D. M. Bullock. Real-Time Measurement of Travel Time Delay in Work Zones and Evaluation Metrics Using Bluetooth Probe Tracking. *Transportation Research Record: Journal of the Transportation Research Board*, no. 2169, 2010, pp. 40-53.
18. Haghani, A., M. Hamed, K. F. Sadabadi, S. Young, and P. Tarnoff. Data Collection of Freeway Travel Time Ground Truth with Bluetooth Sensors. *Transportation Research Record*, no. 2160, 2010, pp. 60-68.
19. Coifman, B., and M. Cassidy. Vehicle reidentification and travel time measurement on congested freeways. *Transportation Research Part A*, Vol. 36, 2002, pp. 899-917.
20. Kwong, K., R. Kavaler, R. Rajagopal, and P. Varaiya. Arterial travel time estimation based on vehicle re-identification using wireless magnetic sensors. *Transportation Research Part C*, no. 17, 2009, pp. 586-606.
21. Yang, J.-S. Travel Time Prediction Using the GPS Test Vehicle and Kalman Filtering Techniques. in *2005 American Control Conference*, Portland, 2005, pp. 2128-2133.
22. Stipanovic, J., L. Miranda-Moreno, and N. Saunier. The Who and Where of Road Safety: Extracting Surrogate Indicators From Smartphon Collected GPS Data in Urban Envrionments. in *Transportation Research Board Annual Meeting 2016*, Washington, DC, 2016.
23. Shi, Q., and M. Abdel-Aty. Big Data applications in real-time traffic operation and safety monitoring and improvement on urban expressways. *Transportation Research Part C*, Vol. 58, 2015, pp. 380-394.
24. Palen, J. The Need for Surveillance in Intelligent Transportation Systems. *Intellimotion*, Vol. 6, no. 1, 1997, pp. 1-3.
25. Falcocchio, J. C., and H. S. Levinson. Measuring Traffic Congestion, in *Road Traffic Congestion: A Concise Guide*.: Springer International Publishing, 2015, pp. 93-110.
26. Skabardonis, A., P. Varaiya, and K. F. Petty. Measuring Recurrent and Nonrecurrent Traffic Congestion. *Transportation Research Record*, no. 1856, 2003, pp. 118-124.
27. Dias, C., M. Miska, M. Kuwahara, and H. Warita. Relationship between congestion and traffic accidents on expressways: an investigation with Bayesian belief networks. in *Proceedings of 40th Annual Meeting of Infrastructure Planning (JSCE)*, Japan, 2009.
28. Bachman, C. Multi-Sensor Data Fusion for Traffic Speed and Travel. University of Toronto, Toronto, Masters Thesis 2011.
29. Marchal, F. TrackMatching. 2015. <https://mapmatching.3scale.net/>. Accessed May 1, 2015.
30. OpenStreetMap. About. *OpenStreetMap*, 2015. <http://www.openstreetmap.org/about>. Accessed May 11, 2015.

31. Marchal, F., J. Hackney, and K. W. Axhausen. Efficient Map Matching of Large Global Positioning System Data Sets. *Transportation Research Record*, no. 1935, 2005, pp. 93-100.
32. 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.
33. Brisk Synergies. *Brisk Synergies*, <http://www.brisksynergies.com/>. Accessed July 22, 2015.
34. 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.



**FIGURE 1** Collection and map matching of smartphone-collected GPS data



**FIGURE 2** Redefinition of OSM links

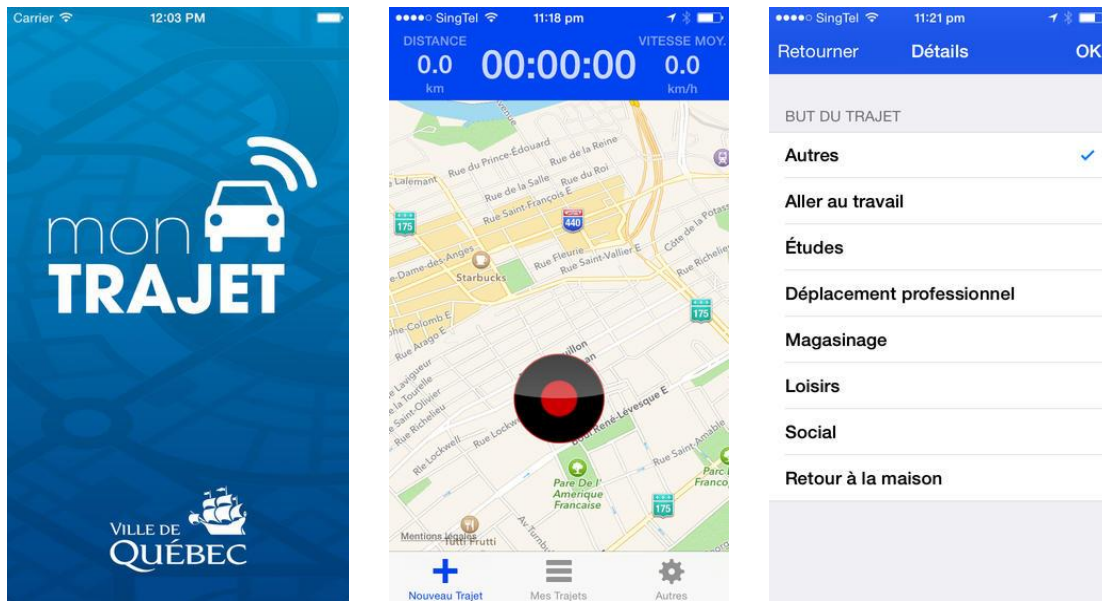


FIGURE 3 Smartphone application interfaces

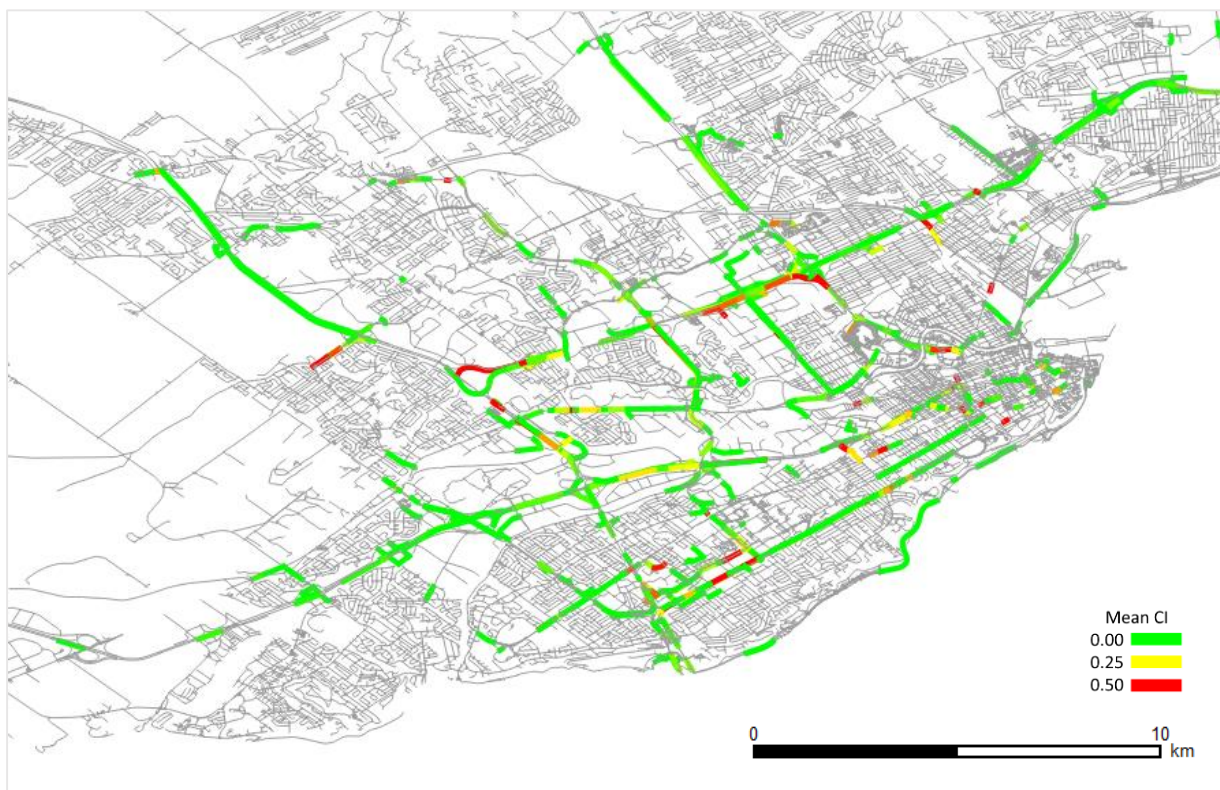
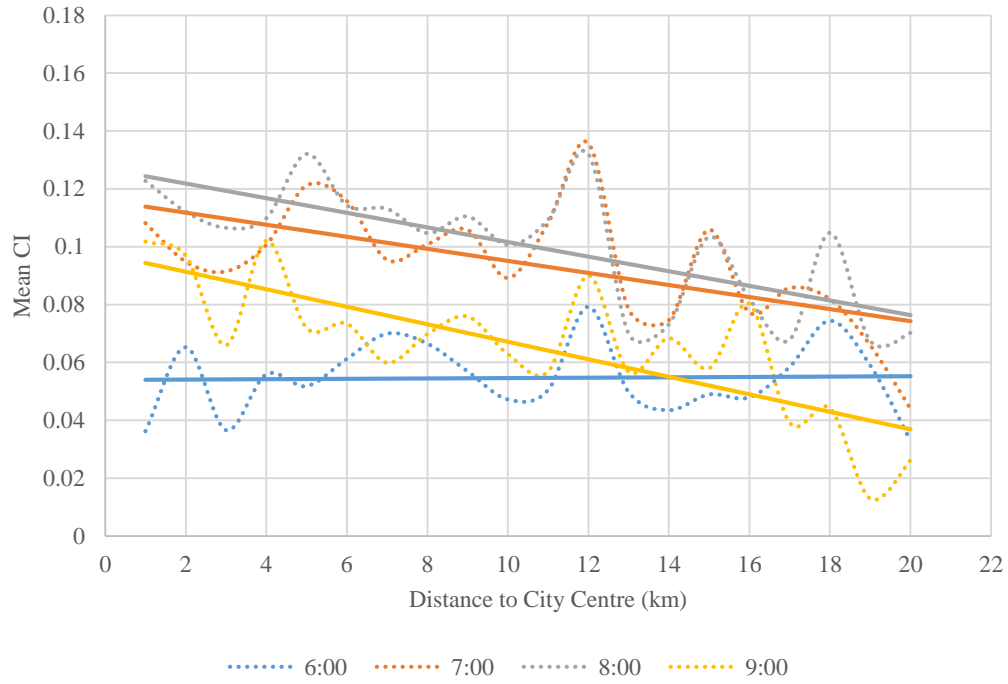
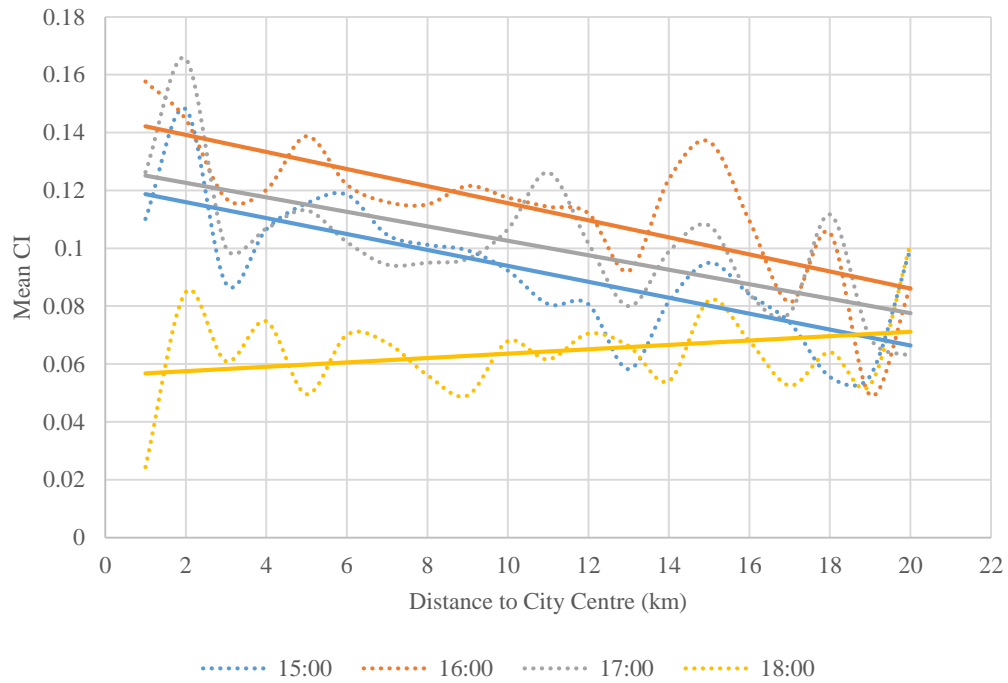


FIGURE 4 CI map for May 6 from 8:00 to 9:00 AM



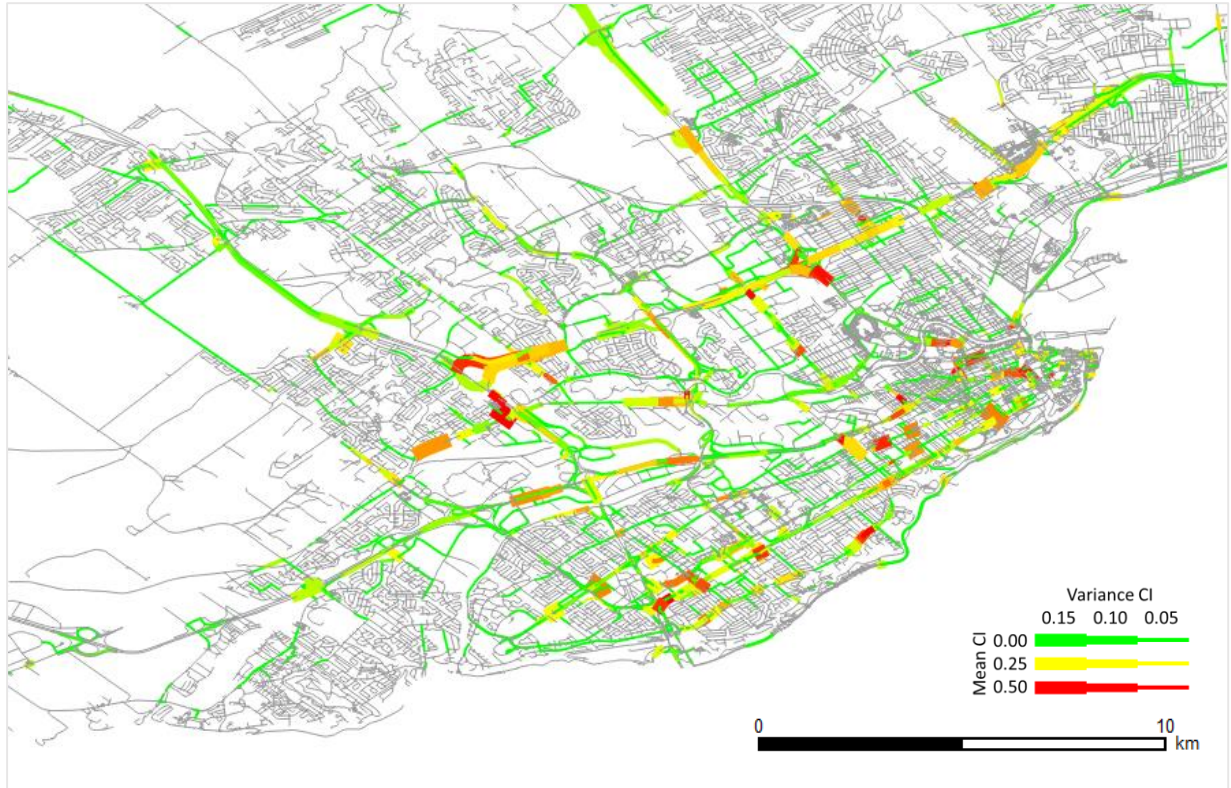


(a)



(b)

**FIGURE 5** Congestion profiles for each hour during AM (a) and PM peak periods (b)



**FIGURE 6** Average and variance of CI on link basis

**TABLE 1** Description of macroscopic congestion formation and propagation in Quebec City

| <b>AM PEAK</b>   |               |   |
|------------------|---------------|---|
| <b>Time</b>      | <b>Period</b> | <b>Description</b>  |
| 6:00 to 7:00 AM  | Onset         | Trips begin to be generated at the outskirts, destined for downtown                                       |
| 7:00 to 9:00 AM  | Peak          | Trips from the outskirts begin to arrive downtown, and more trips are generated throughout the network    |
| 9:00 to 10:00 AM | Dissipation   | Trips are no longer generated at the outskirts, but are still arriving and being generated downtown       |
| <b>PM PEAK</b>   |               |   |
| <b>Time</b>      | <b>Period</b> | <b>Description</b>  |
| 3:00 to 4:00 PM  | Onset         | Trips begin to be generated downtown, destined for the outskirts  |
| 4:00 to 6:00 PM  | Peak          | Trips from downtown begin to arrive at the outskirts, and more trips are generated throughout the network |
| 6:00 to 7:00 PM  | Dissipation   | Trips are no longer generated downtown, but are still arriving at the outskirts                           |