



A novel method for imminent crash prediction and prevention

Zhi Chen^a, Xiao Qin^{b,*}

^a Department of Civil and Environmental Engineering, University of Wisconsin-Milwaukee, 2025 E Newport Ave, NWQ 4515, Milwaukee, WI 53201, United States

^b Department of Civil and Environmental Engineering, University of Wisconsin-Milwaukee, 2025 E Newport Ave, NWQ 4415, Milwaukee, WI 53201, United States



ARTICLE INFO

Keywords:

Crash prediction
Crash prevention
Cell transmission model (CTM)
Variable speed limit (VSL)

ABSTRACT

A crash prediction and prevention method was proposed to detect imminent crash risk and help recommend traffic control strategies to prevent crashes. The method consists of two modules, the crash prediction module and the crash prevention module. The crash prediction module detects crash-prone conditions when the predicted crash probability exceeds a specified threshold. Then the crash prevention module would simulate the safety effect of traffic control alternatives and recommend the optimal one. The proposed method was demonstrated in a case study with variable speed limit (VSL). Results showed that the proposed crash prediction and prevention method could effectively detect crash-prone conditions and evaluate the safety and mobility impacts of various safety countermeasures.

1. Introduction

The development of advanced transportation information systems (ATIS) has made it easier to collect, store, and process traffic data in a real-time fashion. The readily available real-time traffic data offer new opportunities for crash prediction and prevention in terms of traffic control and operations. Many studies have used real-time traffic data to investigate the relationship between crash risk and prevailing traffic conditions. Among all types of traffic sensors, inductive loop detectors have been widely used for real-time crash prediction.

The prevailing traffic circumstances prior to and under which a crash takes place are believed to be one of the major contributors to a crash. Moreover, travel conditions can shift rapidly, and the traffic that a vehicle experienced immediately prior to or at the time of a crash is more relevant than earlier or later traffic conditions. The phenomenon of temporal proximity has been observed and supported in a study that predicted freeway crashes using loop detector data (Abdel-Aty et al., 2004). However, many studies did not consider the traffic conditions occurring right before a crash (e.g. 0–5 min period), citing that preventative actions may take extra time in a real-time crash identification, notification, and prevention system. Therefore, traffic data used in these studies comes from earlier time periods (e.g. 5–10 min before a crash) (Abdel-Aty et al., 2004; Pande and Abdel-Aty, 2006; Hossain and Muromachi, 2012; Sun and Sun, 2015).

The time buffer between traffic data and crash occurrence is also related to the consistency between crash modeling and crash prediction, though it has never been explicitly discussed in previous studies.

Fig. 1 illustrates such consistency by a hypothetical example. The figure shows that one intends to predict the crash risk in the future moment, which is 5 min from now. The traffic conditions in the past 5-min period are known, while those in the future 5-min period are not known. However, crash modeling needs to be conducted in a consistent manner so that resultant crash prediction models can be applied. Initially, the historical crash time is consistent with the hypothesized crash time. Then the 0–5-min period before the crash would be the future 5-min period, and the 5–10-min period before the crash would be the past 5-min period. Therefore, the data from the 5–10-min period before the crash needs to be used for crash modeling so that the crash prediction models can be applied to predict the crash risk in real time, or 0–5-min.

The loop detector spacing can also lead to a lack of consistency, as spacings can vary substantially from site to site and across studies. For example, in one study the spacing ranges from 0.2 to 1.3 mi with an average of 0.5 mi (Xu et al., 2016); in another it ranges from 0.15 to 1.68 mi with an average of 0.5 mi (Xu et al., 2013a); and another example has a range of 0.34 to 2.37 mi with an average of about 1.06 mi (Zheng et al., 2010). Studies have shown that the sensor location may affect the estimation of traffic flow by producing inconsistently biased traffic data (Kwon et al., 2007; Liu and Danczyk, 2009; Danczyk and Liu, 2011; Hong and Fukuda, 2012). The discrepancies in the spatial-tempo domain mean that crash prediction models developed with traffic data collected directly from loop detector stations may be inadequate. Such data issues would undermine the prediction power of developed models. Even when a reliable crash prediction model is available, the issue of deploying effective preventative countermeasures

* Corresponding author.

E-mail addresses: zhichen@uwm.edu (Z. Chen), qinx@uwm.edu (X. Qin).

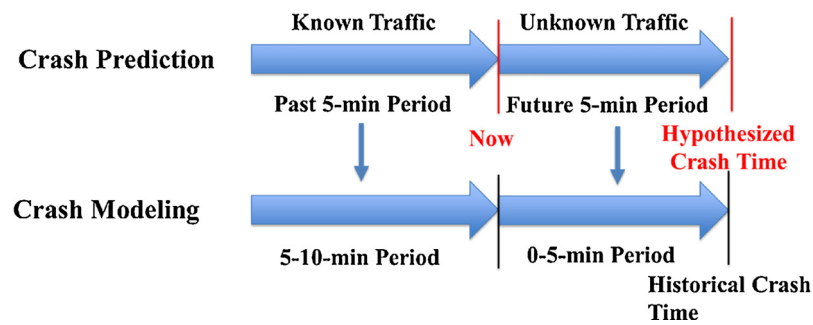


Fig. 1. Consistent time periods for crash prediction and crash modeling.

remains. A performance assessment tool is needed to evaluate the effectiveness of intervening traffic control strategies before their deployment.

The objective of this study is to develop a method for real-time crash prediction and prevention using traffic simulation. Ideally, the method would be able to identify crash-prone conditions by accounting for the spatial-tempo issue of loop detector data, and could efficiently evaluate the performance of traffic control strategy (TCS) alternatives. Inspired by virtual loops extensively applied for vehicle detection, counting, and signal control, the cell transmission model (CTM) was employed to instrument a corridor of highway with virtual detector stations and measure traffic data where physical stations were not available. The paper is organized as follows: Section 2 is the review of relevant literature; Section 3 presents the methodologies of CTM and the binary logistic model; Section 4 describes data collection and processing, as well as CTM calibration and simulation; Section 5 reports on the crash modeling; Section 6 illustrates the crash prediction and prevention method; and lastly, Section 7 presents conclusions and future research.

2. Literature review

Crashes should be more closely related to the traffic conditions occurring during or around the same time of the crash, as opposed to those occurring hours before. One study examined the impact of traffic variables on crash risk using five time slices: 0–5 min before the crash (time slice 1); 5–10 min before the crash (time slice 2); and up to 20–25 min before the crash (time slice 5) (Abdel-Aty et al., 2004). The regression results showed that the traffic variables in time slice 1 are the most statistically significant among all five time slices, which supports the notion that the traffic conditions occurring right before a crash can best model the crash probability. However, most previous studies did not use this time period, citing that extra time was needed to take preventive countermeasures (Abdel-Aty et al., 2004; Pande and Abdel-Aty, 2006; Hossain and Muromachi, 2012; Sun and Sun, 2015). Furthermore, the distance between crash locations and detector locations varies from one case to another, making it impossible to obtain consistent measurements. The aforementioned issues regarding time and distance could undermine the validity and accuracy of real-time crash prediction models.

Ideally, the traffic conditions present at the time of the crash at the crash location should be used in studies that attempt to improve prediction accuracy. Although it is unrealistic to have physical detectors located at every crash location, the development of traffic simulation models has made the virtual detection possible. CTM, a macroscopic traffic flow simulation model that was first proposed by Daganzo (Daganzo, 1994), partitions a highway into continuous cells with user-defined lengths. Under the law of conservation, the traffic density in each cell within the highway evolves and follows the relationships derived from the fundamental diagram.

CTM can well accommodate traffic flow data collected from loop detectors, as they have shown promising results in predicting traffic flows using loop detector data as inputs (Muñoz et al., 2003, 2006,

Sumalee et al., 2011). Muñoz et al. achieved less than 13% of the mean error when simulating density using both CTM and switching-mode model (SMM) (Muñoz et al., 2003), as opposed to density collected from loop detectors. Muñoz et al. improved parameter calibration methods of CTM and SMM (Muñoz et al., 2006); calibrated CTM and SMM produced a 13% and 14% error, respectively, in estimating density, and a 4% and 5% error in estimating flow. Sumalee et al. proposed a stochastic CTM and achieved a 7.9% error in estimating density (Sumalee et al., 2011). CTM is therefore a reliable simulation tool that can generate trustworthy simulated traffic input for predicting crashes. Moreover, well-established traffic flow theories and emerging simulation algorithms provide timely support to the fast development of real-time crash prediction and prevention methods.

The CTM has the capability of simulating traffic control strategies. The CTM has several attractive features (Hadiuzzaman and Qiu, 2013): 1) it is trustworthy in simulating TCS, as it is founded on sound traffic theory; 2) it is parsimonious, as it needs only a few parameters which can be estimated both online and off-line; 3) it requires low computational effort to predict traffic conditions in real-time. Recently, the CTM has been applied to evaluate the safety effects of variable speed limits (VSL). Li et al. developed VSL in CTM and investigated its control strategy to reduce rear-end crash risks near recurrent bottlenecks on a 6-mile long virtual segment (Li et al., 2014b). Later, Li et al. developed a strategy to optimize VSLs on a 29-mile freeway corridor in California (Li et al., 2016). In this study, VSL strategies were optimized to balance the impact on collision risk, injury severity, and travel time.

The relationships between the relatively low number of crashes and the massive volume of real-time traffic data can be sorted out through specific techniques. In general, the approaches for real-time crash prediction can be categorized as either statistical regression models or data mining techniques such as the Kohonen clustering algorithm, neural networks, and the Bayesian network (Pande and Abdel-Aty, 2006; Hossain and Muromachi, 2012; Sun and Sun, 2015). Although data mining methods can accommodate correlation within independent variables for speed, flow, and occupancy (Hossain and Muromachi, 2012), they cannot identify explicit relationships between crash probability and traffic flow variables. Therefore, it is difficult to interpret the crash mechanism and develop effective crash prevention countermeasures. Statistical models, however, can build clear connections between crash probability and traffic flow variables, which is crucial for the development of proactive safety approaches. Among various statistical models used in real-time crash prediction studies, the binary logistic regression is widely used (Abdel-Aty et al., 2005; Zheng et al., 2010; Xu et al., 2013b) because it can easily predict the crash probability given the explanatory variables.

3. Methodology

Crash probability prediction began with using CTM to simulate spatial and temporal traffic during the time period just prior to a crash. The crash occurrence probability was then estimated with simulated traffic conditions using a binary logistic regression model.

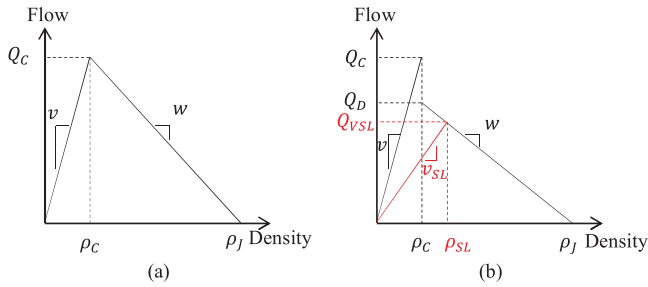


Fig. 2. (a) Triangular fundamental diagram; (b) Fundamental diagram with capacity drop.

3.1. Cell Transmission Model (CTM)

CTM is a macroscopic traffic simulation model proposed by Daganzo (Daganzo, 1994). CTM is a powerful simulation technique which can capture many important traffic phenomena including queue formation and dissipation and shockwave propagation (Daganzo, 1994). CTM is more computationally efficient and easier to configure and calibrate than microscopic simulation models. CTM also operates sufficiently with aggregated traffic data from detector stations. Fig. 2 shows the fundamental diagram with and without a capacity drop for developing CTM.

In CTM, a highway segment is divided into a series of cells. The density of each cell evolves following the conservation law of vehicles. Assuming that Cell i is characterized by the triangular fundamental diagram in Fig. 2(a), where Q_C is the capacity flow, ρ_C is the critical density, ρ_J is the jam density, v is the free-flow speed, and w is the shockwave speed. The density for Cell i without on- or off-ramps is determined by Eq. (1):

$$\rho_i(k + 1) = \rho_i(k) + \frac{T}{l_i}(q_i(k + 1) - q_i(k)) \quad (1)$$

where k is the time step index, $\rho_i(k)$ is the density of Cell i during the k th time step, T is the length of the time step, l_i is the length of Cell i , and $q_i(k)$ is the flow rate into Cell i during the k th time step. The flow rate is determined by the sending and receiving functions. For Cell i , the sending function $S_i(k)$ represents the maximum flow that can be supplied during the k th time step, and the receiving function $R_i(k)$ represents the maximum flow that can be received. The two functions are determined in Eqs. (2) and (3), respectively:

$$S_i(k) = \min(v_i \rho_i(k), Q_{C,i}) \quad (2)$$

$$R_i(k) = \min(Q_{C,i}, w_i(\rho_{J,i} - \rho_i(k))) \quad (3)$$

The flow rate, $q_i(k)$, is determined by:

$$q_i(k) = \min(S_{i-1}(k), R_i(k)) \quad (4)$$

The fundamental diagram changes when the VSL control is deployed, as shown in Fig. 2(b). v_{SL} is the deployed speed limit, and Q_{VSL} and ρ_{SL} are the new capacity and critical density after activating the VSL control. A study by Li et al. (2014b) showed that the sending and receiving functions affected by the VSL control are determined by Eqs. (5) and (6), respectively:

$$S_i(k) = \min(\min(v_i, v_{SL,i}) * \rho_i(k), Q_{VSL,i}) \quad (5)$$

$$R_i(k) = \min(Q_{VSL,i}, w_i(\rho_{J,i} - \rho_i(k))) \quad (6)$$

A phenomenon called ‘‘capacity drop’’ represents the discharge flow rate dropping below capacity after the congestion forms (Hall and Agyemang-Duah, 1991; Cassidy and Rudjanakanoknad, 2005). Accounting for capacity drop helps to better simulate traffic conditions. Capacity drop is accounted for by adopting the fundamental diagram in Fig. 2(b) where Q_D is added to the triangular fundamental diagram. The capacity drops from Q_C to Q_D at the onset of congestion. Similar to the

study by Li et al. (Li et al., 2014b), the modified sending and receiving functions are formulated in Eqs. (7) and (8), respectively:

$$S_i(k) = \begin{cases} v_i \rho_i(k), & \text{if } \rho_i(k) \leq \rho_{C,i} \\ Q_{D,i}, & \text{if } \rho_i(k) > \rho_{C,i} \end{cases} \quad (7)$$

$$R_i(k) = \begin{cases} Q_{C,i}, & \text{if } \rho_i(k) \leq \rho_{C,i} \\ w_i(\rho_{J,i} - \rho_i(k)), & \text{if } \rho_i(k) > \rho_{C,i} \end{cases} \quad (8)$$

3.2. Binary logistic regression model

Eq. (9) shows how the probability of a crash event is formulated in a binary logistic regression model:

$$p(X_i) = \frac{g(X_i)}{1 + e^{g(X_i)}} \quad (9)$$

where $p(X_i)$ represents the crash probability given $X_i = (x_{i,1}, x_{i,2}, \dots, x_{i,k})$, a set of k explanatory variables for sample i , and $g(X_i)$ is a linear combination of the following variable set:

$$g(X_i) = \beta_0 + \beta_1 * x_{i,1} + \beta_2 * x_{i,2} + \dots + \beta_k * x_{i,k} \quad (10)$$

where $(\beta_0, \beta_1, \beta_2, \dots, \beta_k)$ are the corresponding coefficients for $(x_{i,1}, x_{i,2}, \dots, x_{i,k})$.

The parameters $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_k)$ can be estimated by maximizing the following log-likelihood function:

$$\ln L(\beta, X_i) = \sum_{i=1}^n [(\beta_0 + \beta_1 * x_{i,1} + \dots + \beta_k * x_{i,k}) - \ln(1 + e^{-(\beta_0 + \beta_1 * x_{i,1} + \dots + \beta_k * x_{i,k})})] \quad (11)$$

4. Data description and processing

Three data sources were consulted to develop a comprehensive approach: a) 1-min time interval traffic information from the WisTransPortal V-SPOC (Volume, Speed, and Occupancy) application suite (Parker and Tao, 2006); b) crash data from the web-based query and retrieval facility for Wisconsin Department of Transportation crash data and from reports archived in the WisTransPortal data management system; and c) weather information (e.g. snow, rain) from the Road Weather Information System (RWIS) in WisTransPortal.

4.1. Study site and CTM setup

A 4.15-mile corridor on I-94 East in Waukesha, WI was selected as the study site. The site was selected based on the following criteria: spacing of loop detector stations, traffic data quality, and crash sample size. The selected roadway corridor, as shown in Fig. 3, has three lanes with one on-ramp and one off-ramp. The corridor consists of three segments, S_1 , S_2 , and S_3 , which are 1.77-mile, 0.79-mile and 1.59-mile long, respectively. Segment S_2 starts at the end of the off-ramp and ends at the beginning of the on-ramp. The posted speed limit was 65 MPH in S_1 , and 55 MPH in S_2 , and S_3 . Other roadway characteristics such as lane width and shoulder width did not change along the corridor.

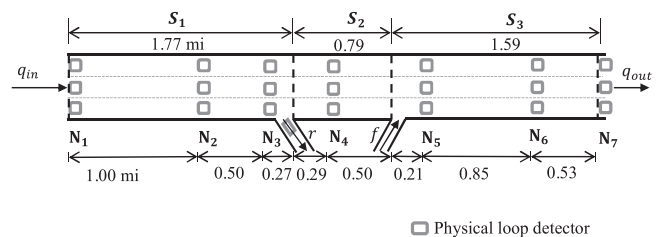


Fig. 3. Layout of physical loop detector stations.

The corridor was instrumented with seven mainline loop detector stations: N_1, N_2, \dots, N_7 . The stations are referred to as physical stations so as to differentiate them from the virtual detectors introduced later. The seven stations space between 0.50 and 1.00 mile, with an average of 0.69 mile and a standard deviation of 0.20 mile. One loop detector station was located on the off-ramp, but no stations exist on the on-ramp. The traffic flow of the on-ramp can be imputed based on the conservation of vehicles using the flows from the nearest upstream and downstream detector stations.

The corridor was divided into 41 virtual cells for CTM simulation, and the cell length is uniform within each of the three segments. Segment S_1 has 17 cells with a length of 0.104 mile; segment S_2 has 8 cells with a length of 0.098 mile; segment S_3 has 17 cells with a length of 0.099 mile. A virtual detector station was instrumented at the boundaries of cells, so there were 42 virtual detector stations and spacing between consecutive virtual stations averaged 0.1 mile with negligible variation. The off-ramp was located at the end of the 17th cell, while the on-ramp was located at the beginning of the 26th cell.

The virtual stations were set up at cell boundaries, similar to physical detector stations, to measure flow, speed, and density. Virtual stations were expected to capture traffic conditions at locations closer to the crash site.

Crashes that occurred at the study site from 2012 to 2014 were included. Any crash that happened within one hour after a crash occurrence was considered a secondary crash and was subsequently removed as indicated in (Hirunyanitiwattana and Mattingly, 2006). Crashes with missing times were excluded, as crash time is required to retrieve the traffic data.

A critical component of developing a crash prediction model is the knowledge of the traffic conditions experienced by the vehicle right before a crash; therefore, it is important to pinpoint the exact time in which a crash occurs. Crash times are sometimes rounded to the nearest 5-minute time stamp, and are therefore not reliable (Golob and Recker, 2003; Kockelman and Ma, 2004). Crash times in this study were carefully reviewed, and no rounding issue was found. Crashes were then randomly sampled and compared to the abrupt changes in traffic conditions based on which crash times could be identified (Abdel-Aty et al., 2005; Zheng et al., 2010). The validation result was positive, and the crash times from the database were used as the actual crash occurrence times.

4.2. CTM calibration

A fundamental diagram is required to operate the CTM simulation. Differing roadway characteristics (e.g., horizontal curves, distances to on-/off-ramps, posted speed limits) mean different cells could have varying traffic patterns, which lead to different fundamental diagrams. Thus, one fundamental diagram was calibrated using the traffic data collected from each mainline detector station. The fundamental diagram was based on the flow-density plot. The flows and speeds were collected from the loop detector stations, while the densities were determined by flow and speed.

The calibration algorithm in Dervisoglu et al. (Dervisoglu et al., 2009) was adopted with modifications to calibrate the fundamental diagram. The full description of the algorithm is summarized as follows:

- 1 Estimate the free-flow speed, v , using the least-squared method with flow-density pairs in the free-flow conditions. Since the speed limits of the segments are 65 MPH and 55 MPH, data points with speeds exceeding 55 mi/h in segment S_1 and 45 mi/h in segments S_2 , and S_3 were deemed to be in free flow conditions.
- 2 Find the maximum measured flow rate, q_{max} , as the capacity, Q_C . Critical density is determined by $\rho_C = \frac{Q_C}{v}$. Few and unsustainable observations with extremely high flow rates, a phenomenon of capacity overestimation, were observed. The formula to compute the

Table 1
Fundamental Diagram Parameters by Physical Station.

Station	v (mi/h)	ρ_C (veh/mi) ^a	ρ_J (veh/mi) ^a	Q_C (veh/h) ^a	Q_D (veh/h) ^a	w (mi/h)
N_1	67.0	106.1	486.0	7111	6890	18.1
N_2	68.4	104.6	588.4	7152	6816	14.1
N_3	66.5	106.7	472.2	7095	6603	18.1
N_4	59.8	97.0	799.0	5796	4989	7.1
N_5	60.8	113.9	779.9	6924	6671	10.0
N_6	58.0	118.0	460.4	6839	6703	19.6
N_7	60.1	114.8	375.5	6903	6683	25.6

^a Parameters are for three lanes.

nominal capacity (in veh/h/lane) of freeways in HCM 2010 was adopted, as opposed to using the high flow rates (Transportation Research Board, 2010):

$$Capacity = \begin{cases} 2400 \text{ veh/h/lane, if } FFS \geq 70 \text{ mi/h} \\ 2400 - 10 \times (70 - FFS) \text{ veh/h/lane, if } FFS < 70 \text{ mi/h} \end{cases} \quad (12)$$

The capacity was then determined by taking the minimum of Q_C and the nominal capacity given by Eq. (12).

- 3 Estimate the shockwave speed, w , and the jam density, ρ_J , using the least-squared method with flow-density pairs exceeding the critical density. The flow rate after the capacity drop was set as the value on the fitted flow-density line at the critical density.

Following the modified algorithm, fundamental diagram parameters were obtained for each physical detector station as shown in Table 1. Note that ρ_C , ρ_J , Q_C and Q_D are for three lanes. The magnitude of the capacity drop is from 2.0% to 6.9% for all physical stations except N_4 which has a 13.9% capacity drop rate. The set of fundamental diagram parameters calibrated for one physical station was assigned to cells near that station.

4.3. CTM simulation

The simulation time step in CTM needs to be chosen so that the Courant–Friedrichs–Lewy (CFL) condition (Courant et al., 1967) is fulfilled. A vehicle cannot travel across more than one cell during one simulation step in the CFL condition, i.e., $v_i * \Delta t \leq l_i$ where v_i is the free-flow speed, Δt is the simulation time step, and l_i is the cell length. A 5-sec time step was used ($\Delta t = 5s$) based on the lengths of cells.

Entering flow and exiting flow of the highway corridor are required to run the CTM. The four flow inputs were required for the study site, including in-flow, q_{in} , out-flow, q_{out} , off-ramp flow, r , and on-ramp flow f (as shown in Fig. 3). The 1-min flow data collected from the first physical station, N_1 , and the last physical station, N_7 , in the 0–5 min period prior to a crash/non-crash were used as the in-flow and out-flow of the corridor. A linear interpolation method was applied to generate the 5-s in-flow, out-flow, on-ramp flow and off-ramp flow data. A CTM was then run to simulate how traffic in cells along the corridor evolves at each time step within the 5-min time interval.

In addition to the flow data, initial densities of cells at the beginning of the simulation interval are also needed for the CTM simulation. The initial density of a cell was obtained from the station’s density data as long as the cell had one loop detector station. Densities of cells between two such cells were interpolated using the following approach:

- 1 Compute the density change rate as the ratio of the difference in densities of two cells with two consecutive loop detector stations and the distance between them: $\nabla \rho = \frac{\rho_{d,0} - \rho_{u,0}}{x_d - x_u}$, where $\nabla \rho$ is the density change rate; $\rho_{d,0}$ and $\rho_{u,0}$ are densities of cells having the downstream and upstream detector stations, respectively; x_d and x_u are the locations of the beginnings of the two cells, that is, the

locations of the two detector stations.

- Determine the initial density of one cell between those two cells by the following: $\rho_{i,0} = \rho_{u,0} + \nabla\rho^*(x_i - x_u)$, where x_i is the location of the beginning of one cell between the two cells.

5. Crash modeling

The simulated traffic data were collected from the virtual upstream and downstream stations to the cell location of each crash/non-crash in the prior 0-5-min period. The time period of 0–5 min prior to a crash was used in order to account for the temporal issue of physical station data, as the simulated traffic data in the future 5-min period would be employed for crash prediction. More details will be illustrated in Section 6. 0.2 mi was selected as the distance from the crash cell location to its virtual upstream and downstream stations. One virtual upstream station and one virtual downstream station that are both 0.2 mi (i.e., two cells) away from the crash cell location were identified as stations from which to collect the simulated traffic data.

The spacing between virtual upstream and downstream stations for the 0.2-mi distance setting is 0.5 mi, which is not larger than the smallest spacing between physical stations. Therefore, the 0.2-mi distance setting provides traffic data from stations with both uniform and short distances from the crash (non-crash) location. The feasibility of uniform and close distances can be tested by comparing the performance of two different models: Model V, which is developed with virtual station data in the 0.2-mi distance setting; and Model P, which is developed with physical station data.

The 5-s traffic data from the two selected virtual stations were aggregated into the 5-min interval for each crash and non-crash case and converted into traffic flow variables in Table 2. Due to the inter-correlation between the three traffic parameters of flow, density, and speed, traffic variables related to density and speed were kept to avoid serious correlations between candidate variables.

Three additional groups of traffic variables were considered aside from mean and standard deviation of density and speed, roadway characteristics, and weather factors that have been frequently used in previous studies (Abdel-Aty et al., 2004; Abdel-Aty and Pande, 2006; Pande and Abdel-Aty, 2006; Abdel-Aty et al., 2012; Xu et al., 2012,

2016). The first group is related to the time-series difference in density and speed; the second group is related to the difference between downstream and upstream density and speed; the third group is related to the traffic state of the location.

The time-series difference is the difference between the density or speed in the next 5-s and that in this 5-s. Variables such as AvgTsdDen_u (average time-series absolute difference in 5-s density at the upstream station) and StdTsdDen_u (standard deviation of time-series difference in 5-s density at the upstream station) were calculated by Eq. (13) and (14), and AvgTsdSpd_u (average time-series absolute difference in 5-s speed at the upstream station) and StdTsdSpd_u (standard deviation of time-series difference in 5-s speed at the upstream station) were calculated in the same way,

$$AvgTsdDen_u = \frac{\sum_{t=1}^{59} |Den_{u,t+1} - Den_{u,t}|}{59} \tag{13}$$

$$StdTsdDen_u = \sqrt{\frac{\sum_{t=1}^{59} \left[(Den_{u,t+1} - Den_{u,t}) - \frac{\sum_{t=1}^{59} (Den_{u,t+1} - Den_{u,t})}{59} \right]^2}{59 - 1}} \tag{14}$$

where $Den_{u,t}$ is the 5-s upstream density at time step $t = 1, 2, \dots, 60$ (60 5-s in one 5-min interval). This variable group measures the traffic trend over time. The average absolute time-series difference in density or speed measures the traffic stability over time, and a large value indicates that the traffic is very unstable. The standard deviation of time-series difference in density or speed measures the consistency of traffic changes, and a large value indicates that the traffic changes are very fluctuant over time.

The second group is related to the difference between downstream and upstream density and speed. Variables such as AvgDiffDen_{d-u} and StdDiffDen_{d-u} were computed by Eqs. (15) and (16), and AvgDiffSpd_{d-u} and StdDiffSpd_{d-u} were calculated in the same way,

$$AvgDiffDen_{d-u} = \frac{\sum_{t=1}^{60} (Den_{d,t} - Den_{u,t})}{60} \tag{15}$$

Table 2
Candidate Variables.

Variable	Description
AvgDen _u	Average 5-s density at the upstream station (veh/mi)
AvgSpd _u	Average 5-s speed at the upstream station (mi/h)
StdDen _u	Standard deviation of 5-s density at the upstream station (veh/mi)
StdSpd _u	Standard deviation of 5-s speed at the upstream station (mi/h)
AvgTsdDen _u	Average time-series absolute difference in 5-s density at the upstream station (veh/mi)
AvgTsdSpd _u	Average time-series absolute difference in 5-s speed at the upstream station (mi/h)
StdTsdDen _u	Standard deviation of time-series difference in 5-s density at the upstream station (veh/mi)
StdTsdSpd _u	Standard deviation of time-series difference in 5-s speed at the upstream station (mi/h)
AvgDen _d	Average 5-s density at the downstream station (veh/mi)
AvgSpd _d	Average 5-s speed at the downstream station (mi/h)
StdDen _d	Standard deviation of 5-s density at the downstream station (veh/mi)
StdSpd _d	Standard deviation of 5-s speed at the downstream station (mi/h)
AvgTsdDen _d	Average absolute time-series difference in 5-s density at the downstream station (veh/mi)
AvgTsdSpd _d	Average absolute time-series difference in 5-s speed at the downstream station (mi/h)
StdTsdDen _d	Standard deviation of time-series difference in 5-s density at the downstream station (veh/mi)
StdTsdSpd _d	Standard deviation of time-series difference in 5-s speed at the downstream station (mi/h)
AvgDiffDen _{d-u}	Average difference between 5-s downstream and upstream density (veh/mi)
AvgDiffSpd _{d-u}	Average difference between 5-s downstream and upstream speed (mi/h)
StdDiffDen _{d-u}	Standard deviation of difference between 5-s downstream and upstream density (veh/mi)
StdDiffSpd _{d-u}	Standard deviation of difference between 5-s downstream and upstream speed (mi/h)
FF	1 = if the location is in the free-flow state; 0 = otherwise
BN	1 = if the location is in the bottleneck front state; 0 = otherwise
BQ	1 = if the location is in the back-of-queue state; 0 = otherwise
CT	1 = if the location is in the congestion state; 0 = otherwise
Curve	1 = Horizontal curve section; 0 = otherwise
Rain	1 = if the weather is rainy; 0 = otherwise
Snow	1 = if the weather is snowy; 0 = otherwise

$$StdDiffDen_u = \sqrt{\frac{\sum_{t=1}^{60} [Den_{d,t} - Den_{u,t}] - AvgDiffDen_{d-u}]^2}{60 - 1}} \tag{16}$$

where $Den_{d,t}$ is the 5-s downstream density at time step $t=1, 2, \dots, 60$. This variable group indicates the difference between traffic conditions upstream and those downstream from the crash location. A large average difference in density or speed implies that the upstream traffic conditions are very different from the downstream traffic conditions. A large standard deviation of the differences implies that the traffic difference is not very consistent. Although the average absolute difference in upstream and downstream traffic parameters appears to have a significant relationship with the crash occurrence in (Xu et al., 2014, 2016), the average of the regular difference rather than of the absolute difference was considered because the sign may carry crucial information.

The third group is associated with the traffic state at the crash/non-crash location. The average density was used to measure the level of traffic congestion at the virtual upstream and downstream station (Yeo et al., 2013). Traffic is congested if the average density is greater than the critical density; otherwise, traffic is in free flow. The traffic state was determined based on the combination of the upstream and downstream traffic conditions:

- 1 Free Flow (FF): when both upstream state and downstream state are free flow;
- 2 Bottleneck front (BN): when upstream is congested and downstream is free flow;
- 3 Back of queue (BQ): when upstream is free flow and downstream is congested; and
- 4 Congested traffic (CT): when both upstream and downstream are congested.

The CTM cannot run for crash cases that have missing physical detector data, so after such crashes were removed, a total of 113 crashes remained crash modeling. 2260 non-crash cases with a 20:1 non-crash to crash case ratio were randomly selected from 1,578,240-min intervals in 2012–2014 at one out of 41 cells. Only the non-crash cases that are not within 2 h from any crash were selected. The 5-min traffic data consisting of data from five 1-min intervals were retrieved from physical stations for non-crash cases in the same way that data were retrieved for crashes. The data were employed to generate simulated traffic data using the CTM. Candidate variables for all non-crash cases were obtained as well. The final dataset consists of 113 crash cases and 2260 non-crash cases.

Table 3 shows the distribution of crash and non-crash cases by traffic state. Most crashes happened in the FF state, while the fewest happened in the BN state. The ratio of crash cases to non-crash cases indicates the crash probability in each state, and a larger ratio suggests a more crash-prone state. As expected, the ratios in the BN, BQ and CT states were considerably higher than those of the FF state.

Traffic patterns may vary in different traffic states, so the traffic flow variables could have distinct distributions across traffic states. For example, Xu et al. (2012) observed varying speed differences between upstream and downstream stations for different traffic states. The hypothesis was tested by dividing the whole dataset into subsets by traffic state. The distributions of traffic flow variables across traffic states were

Table 3
Case Frequency by Traffic State.

Traffic State	Crash	Non-Crash	Ratio
FF	62	1,978	1:31.9
BN	5	90	1:18
BQ	15	95	1:6.3
CT	31	97	1:3.1

Table 4
Number of Significant Runs for Candidate Variables.

Variable	FF	BN	BQ	CT
AvgDen _u	10	0	0	10
AvgSpd _u	10	0	0	0
StdDen _u	10	0	0	10
StdSpd _u	10	0	0	0
AvgTsdDen _u	10	1	0	1
AvgTsdSpd _u	10	0	1	5
StdTsdDen _u	10	0	0	1
StdTsdSpd _u	10	0	0	5
AvgDen _d	10	0	8	0
AvgSpd _d	10	0	9	5
StdDen _d	10	0	0	0
StdSpd _d	10	0	0	0
AvgTsdDen _d	10	1	10	2
AvgTsdSpd _d	10	1	7	0
StdTsdDen _d	10	0	6	5
StdTsdSpd _d	10	0	1	1
AvgDiffDen _{d-u}	0	0	0	10
AvgDiffSpd _{d-u}	0	0	0	10
StdDiffDen _{d-u}	10	0	2	1
StdDiffSpd _{d-u}	10	1	0	2
Curve	10	0	10	0
Rain	0	0	0	0
Snow	9	0	0	0

compared using a t-test. The comparison results show that all traffic variables have different distributions over all four states; and most traffic variables have different distributions in any two states, indicating that it would not be appropriate to develop a single model for all states without considering the interaction between the traffic variables and traffic states.

Crash-prone variables could vary in different traffic states. Data subsets for different states were used to identify statistically significant variables in each state. In each traffic state, the significance of each candidate variable was identified by developing a binary logit model for that variable only. A 10-fold modeling procedure was conducted to avoid spurious significance; the dataset for one traffic state was randomly split into ten subsets, and all variables' significance was checked for any nine out of the ten data subsets. Table 4 reports the number of significant runs for all candidate variables based on the 10% significance level. A variable was identified as truly significant and was kept for further modeling if it was significant in at least eight out of ten runs. Correlations between significant variables in each traffic state were examined. Candidate models were developed with a maximum number of uncorrelated significant variables for each, and the model with the smallest AIC was selected as the optimal model.

Table 5 presents the modeling results by traffic states. The table

Table 5
Modeling Results of Crash Prediction Model by Traffic State.

Variable	Estimate	Standard Error	P-value
FF			
Intercept	-4.586	0.246	< 0.001
StdTsdDen _d	0.460	0.084	< 0.001
StdTsdSpd _d	0.942	0.256	< 0.001
Snow	1.182	0.495	0.017
BN			
Intercept	-2.415	0.466	< 0.001
BQ			
Intercept	-3.762	0.912	< 0.001
StdTsdDen _d	0.425	0.168	0.011
Curve	2.710	0.842	0.001
CT			
Intercept	-2.642	0.865	0.002
AvgDen _u	0.00824	0.00391	0.035

shows that different traffic states have varying contributing variables. The coefficients of $StdTsdDen_d$ and $StdTsdSpd_d$ for the FF state are positive, indicating that the crash risk increases as density and speed at downstream stations are more fluctuant. This is a logical finding because large variations in time-series changes in density and speed reflect turbulent traffic conditions that could increase crash potential. The Snow indicator has a positive sign, implying that snow contributes to crash occurrence in free flow traffic. However, this is not significant in the other traffic states, possibly because drivers tend to drive faster in free flow traffic than in the other states. No variables show significance for the BN state, possibly due to the small sample size.

The positive signs of $StdTsdDen_d$ and Curve for the BQ state indicate that the fluctuant time-series density at the downstream station near the curve would contribute to crash occurrence. The curve indicator shows significance only in this state; this could be because vehicles from the upstream free-flow traffic need to slow down to accommodate slow-moving traffic during congestion at downstream stations, and presence of a curve may worsen the deceleration. $AvgDen_u$ is significant and has a positive coefficient in the CT state. The finding indicates that crash risk increases with the increase in density at the upstream station. Upstream traffic is already congested at the upstream station in the CT state, which would increase upstream density and make the small distance headway even smaller, leading to higher crash likelihood.

Separate models by traffic states were combined into one model, Model V, to assess the impact of two distances on the prediction performance. Model V include the indicator of BN, BQ, and CT state (FF is the reference state), along with interaction terms of traffic states and other variables. Interaction terms were constructed as the interaction of one traffic state and its significant variables, as identified in Table 6. For example, $StdTsdDen_d$ is significant in the FF state, and $FF \times StdTsdDen_d$ is then the interaction term in the combined model. The modeling results show that main effects of both BN and CT states are statistically significant, while the main effect of BQ state is not significant. All interaction terms remain significant and their signs remain the same.

A crash prediction model, Model P, was developed in the same way with observed traffic data collected from physical stations for comparison with Model V. The prediction accuracy of the two models was checked by conducting the 10-fold cross-validation with significant variables from each model. The 10-fold cross-validation method first randomly partitions the dataset into ten equally sized subsamples. A single subsample is used as the validation dataset, and the other nine are used as training datasets. A model was then fitted with significant variables given the training dataset, and was then used to predict the crash probability of observation in the validation dataset. This procedure was repeated ten times, with each of the ten subsamples used exactly once as the validation dataset.

Based on the validation results, ROC (receiver operating characteristic) curves for these two models are plotted in Fig. 4, and the AUC (Area Under Curve) values are 0.80 for Model V and 0.78 for Model P, respectively. The ROC curve is a plot of sensitivity against 1-specificity for different thresholds of predicted crash risk. The

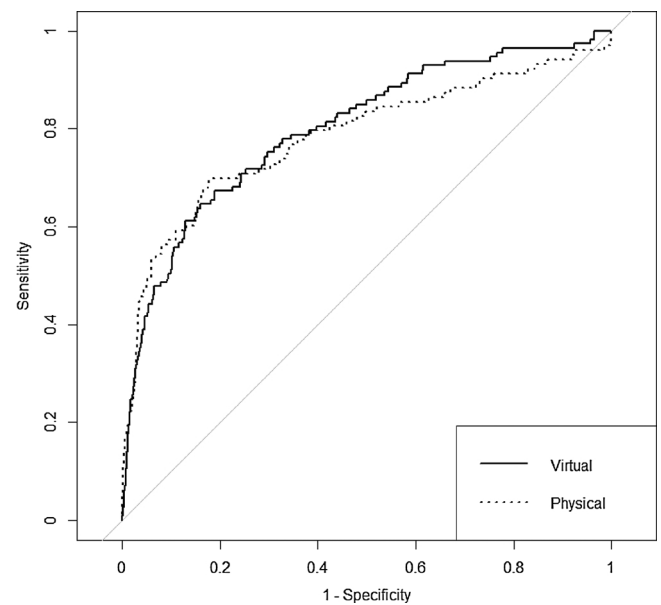


Fig. 4. ROC curves for models with different data sources.

sensitivity represents the proportion of correctly predicted crash cases among all crash cases, or the prediction accuracy of crash cases, while specificity represents the proportion of correctly predicted non-crash cases among all non-crash cases. 1-specificity is the proportion of incorrectly predicted non-crash cases among all non-crash cases, which is also called the false alarm rate. A higher sensitivity along with a lower 1-specificity is preferred. The AUC value represents the total prediction accuracy, and a higher value is favored. Model V provides a higher AUC than Model P. It suggests that simulated traffic data from uniformly and closely spaced virtual stations can provide better model performance by considering the spatial issue of physical station data.

A pre-specified crash probability threshold was determined to classify crashes from non-crashes based on Model V. Both the sensitivity and specificity needs to be balanced to achieve overall desirable classification performance. Equal weights were assigned to the sensitivity and specificity, and the pre-specified crash probability threshold that yielded the maximum weighted summation was 0.0482. The yielded sensitivity and specificity are 0.646 and 0.839, respectively. It means 64.6% (around 73 out of 113) crashes and 83.9% (around 1896 out of 2260) non-crashes can be correctly identified by Model V. This crash probability threshold will be applied in the next Section.

6. Crash prediction and prevention method

In this study, a crash prediction and prevention method was proposed. The purpose is to identify crash-prone traffic conditions in real time and to evaluate TCS alternatives for effectiveness in reducing crash risk. Fig. 5 presents the working process.

The method consists of a crash prediction module and a crash prevention module. The crash prediction module takes the real-time data as the input. It first simulates the traffic in the future 5-min period using CTM and predicts the crash risk for that period based on simulated traffic data. If the predicted crash risk exceeds the pre-specified threshold, the crash prevention module will be activated. Several candidate TCS alternatives may be considered to reduce the crash risk. Each TCS alternative will be simulated in CTM to produce what traffic conditions would be in the future 5-min period if that TCS is deployed. The predicted crash risk is estimated based on the simulated traffic data, and the safety impact of that TCS is evaluated. The optimal TCS is chosen based on established criteria.

In the real-time crash prediction module, traffic conditions during the future 5-min period need to first be simulated using CTM. The initial

Table 6
Results of the Combined Model, Model V.

Variable	Estimate	Standard Error	P-value
Intercept	-4.542	0.238	< 0.001
BN	2.126	0.524	< 0.001
CT	1.899	0.897	0.034
$FF \times StdTsdDen_d$	0.447	0.083	< 0.001
$FF \times StdTsdSpd_d$	0.946	0.255	< 0.001
$FF \times Snow$	1.168	0.494	0.018
$BQ \times StdTsdDen_d$	0.551	0.083	< 0.001
$BQ \times Curve$	3.196	0.657	< 0.001
$CT \times AvgDen_u$	0.00824	0.00392	0.035

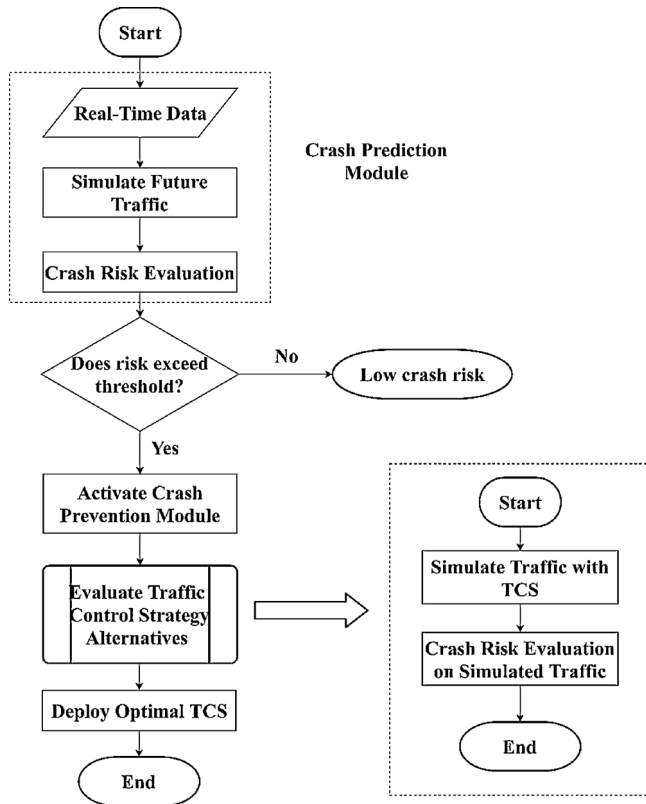


Fig. 5. Process of the crash prediction and prevention method.

densities of all cells were estimated with densities from the seven physical stations at the current moment. The flow inputs, including inflow, q_{in} , off-ramp flow, r , and on-ramp flow f (as shown in Fig. 2) in the future 5-min period were required for CTM simulation and were estimated using the k-nearest neighbor (k-NN) approach. The k-NN approach has been applied in a number of studies to forecast traffic flow rates and has shown promising results (Oswald et al., 2001; Smith et al., 2002; Clark, 2003; Habtemichael and Cetin, 2016).

The past 30 min was considered to be the most recent time period. Flows in the recent time period were considered as the subject flow set. All flow sets during the same time period from last 90 days were considered as candidate flow sets and were matched with the subject flow set. The ten nearest matches with the ten smallest distances were selected. The distance is determined by the following:

$$D(X^m, Y) = \sqrt{\sum_{i=1}^{30} (x_i^m - y_i)^2}, \quad m = 1, \dots, 90 \quad (17)$$

where $X^m = (x_1^m, \dots, x_{30}^m)$ is the m th candidate flow set of 30 1-min flow points; $Y = (y_1, \dots, y_{30})$ represents the subject flow set. The flow in the future 5-min period is calculated as the weighted average of flows in the next 5-min period for those matched flow sets by the following:

$$Y^F = \frac{1}{10} \sum_{k=1}^{10} \frac{(D_k)^2}{\sum_{k=1}^{10} (D_k)^2} X^{k,F} \quad (18)$$

where $Y^F = (y_1^F, \dots, y_5^F)$ represents the estimated flow set in the future 5-min period, D_k is the k th smallest distance for k th nearest matched flow sets among those 10 nearest matched sets, and $X^{k,F} = (x_1^{k,F}, \dots, x_5^{k,F})$ is the flow set in the next 5-min period for k th nearest matched flow sets.

After the required flows are estimated, they are used to run the CTM to simulate traffic in the future 5-min period. Simulated traffic is then used to predict the crash risk of each cell. Simulated traffic data for each cell is collected from its virtual upstream and downstream stations, both of which are 0.2 mi away, and is then converted into variables as

presented in

The predicted crash risk of Cell i is estimated as

$$P_i = \frac{e^\pi}{1 + e^\pi} \quad (19)$$

$$\pi = -4.542 + 2.126*BN + 1.899*CT + 0.447*(FF \times StdTsdDen_d) + 0.946$$

$$*(FF \times StdTsdSpd_d) + 1.168*(FF \times Snow) + 0.551$$

$$*(BQ \times StdTsdDen_d) + 3.196*(BQ \times Curve) + 0.00824$$

$$*(CT \times AvgDen_u)$$

Crash-prone traffic conditions are detected when the predicted crash probability exceeds an established threshold. If crash-prone conditions are detected, the crash prevention module will be activated. The safety impacts of various TCS are then evaluated. The optimal traffic control strategy is then deployed to improve the safety condition.

The proposed crash prediction and prevention method was applied to the study site for demonstration. The VSL control was chosen due to its effectiveness in reducing crashes (Abdel-Aty et al., 2006; Lee et al., 2006; Lee and Abdel-Aty, 2008). The method can be extended to evaluate the safety and mobility impacts of other TCSs such as ramp metering, HOV lane control, hard shoulder running and queue warning whose effectiveness can be successfully simulated in CTM (Gomes and Horowitz, 2006; Kim and Yeo, 2013; Li et al. 2017). Fig. 6 presents the layout of VSL signs along the study corridor. Eight coordinated VSL signs are marked from VSL 1 to VSL 8 and all spacings between adjacent VSL signs are 0.50 mi. Each 0.50-mi spacing consists of five uniform 0.10-mi cells, so there are 35 cells between VSL 1 and VSL 8.

The VSL control strategy proposed in this study was to gradually reduce the posted speed limits of activated VSL signs until a target speed reduction was achieved. When the predicted crash probability of one cell in the future 5-min interval exceeds the pre-specified threshold, the nearest upstream VSL sign will be activated. The pre-specified crash probability threshold was set to be 0.0482 because it provided desirable classification performance with the maximum summation of sensitivity and specificity.

Several parameters need to be decided to develop an effective VSL control, including target speed drop, speed change rate, and maximum speed difference between adjacent VSL signs. Two target speed drop alternatives were proposed: 10 MPH and 20 MPH speed reduction. The target speed limit would be 55 MPH after a 10 MPH speed reduction and 45 MPH after a 20 MPH speed reduction, with an initial speed limit of 65 MPH. The speed change rate determines how fast the VSL sign should change the posted speed limit. A large speed change rate may introduce significant traffic disturbances, whereas a small speed change rate could fail to achieve the target speed limit in a reasonable time period. VSL signs were coordinated to create smooth speed changes between consecutive links. The maximum speed difference between adjacent VSL signs needs to be satisfied. The speed change rate was set to be 10 MPH per 30 s, meaning that the posted speed limit reduces by 10 MPH and stays for 30 s until the next speed change. The maximum speed difference between consecutive VSL signs was set to be 10 MPH. The values for these two parameters have been proven to produce a satisfactory performance for the VSL control (Li et al., 2014a).

Once the crash prevention module is initiated, the proposed VSL

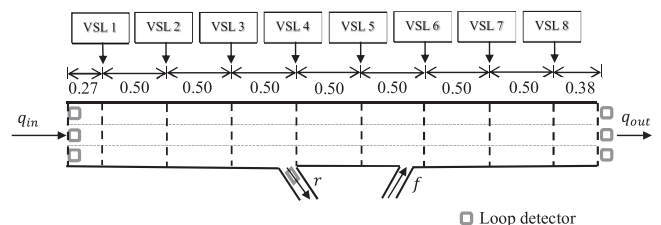


Fig. 6. Layout of VSL signs along the corridor.

strategy with two speed drop alternatives would be simulated in the CTM for 5 min, and then simulated traffic would be used to assess the safety effects and mobility effects. The safety effect is measured as

$$R = \sum_{i=1}^I r_i$$

$$r_i = \begin{cases} p_i - p_{thre}, & \text{if } p_i > p_{thre} \\ 0, & \text{if } p_i \leq p_{thre} \end{cases} \quad (20)$$

where R is the crash risk of the corridor, r_i is the crash risk of Cell i , p_i is the predicted crash probability of Cell i and can be estimated using Eq. (19) given the simulated traffic flow, p_{thre} is the threshold of predicted crash probability for crash classification, which is 0.0482. The mobility effect is measured by the Total Travel Time (TTT).

The proposed crash prediction and prevention method was tested on the 113 crash cases and 2260 non-crash cases that were used for developing crash prediction models. Five minutes before each crash occurrence was equivalent to the “current moment”; the 0-5-min interval before its crash time was equivalent to the “future 5-min period”; the 30-min interval before the “current moment” was equivalent to the recent time period. The flows were estimated using the k-NN approach and were then applied to simulate the traffic in the “future 5-min period”. The crash risk of each cell was re-predicted using Eq. (19) based on the simulated traffic. The crash prevention module was activated when the crash risk of any cell exceeded the threshold. One control strategy would be deployed among three alternatives: 1) Non-activated VSL, 2) VSL control with 10 MPH reduction, and 3) VSL control with 20 MPH reduction. The non-activated VSL strategy would not change the traffic conditions and therefore would not change the crash risk. The control strategy that can provide the smallest crash risk would be deployed.

The proposed VSL strategy consists of three VSL alternatives. The prevention method evaluated the safety impacts of all alternatives and recommended the best one, which is different from the method proposed by Kononov et al. (2012). Kononov et al. determined the target speed limit of VSL based on crash potential that is defined as the product of observed traffic density and the square of observed traffic speed. The target speed limit was determined in such a way that the resultant crash potential would be lower than a pre-specified threshold. However, the approach of Kononov et al. is not applicable in this study as there is no clear relationship between the predicted crash probability and the target speed limit in Eq. (19).

The effectiveness of the crash prevention module was evaluated based on the relative change in R , and TTT. The relative change in the three measures is estimated by

$$\Delta M = \frac{\sum_{k=1}^K M_{k,CPS} - \sum_{k=1}^K M_{k,Non}}{\sum_{k=1}^K M_{k,Non}} \times 100\% \quad (21)$$

where ΔM is the percentage of relative change in one measure (i.e., R , or TTT), $M_{k,CPS}$ is the measure of case k with the crash prevention module, and $M_{k,Non}$ is the measure of case k without the crash prevention module.

Table 7 shows the safety and mobility effects by control strategy. The crash prevention module was triggered 351 times, including 65 times out of 113 crash cases and 286 times out of 2260 non-crash cases.

Table 7
Safety and Mobility Effects by Deployed Control Strategy.

Control Strategy	Crash Cases			Non-Crash Cases			Total		
	N	ΔR	ΔTTT	N	ΔR	ΔTTT	N	ΔR	ΔTTT
Non-activated VSL	55	0%	0%	216	0%	0%	271	0%	0%
VSL:10 MPH Reduction	6	-21.9%	5.2%	29	-27.3%	6.2%	35	-26.3%	6.0%
VSL:20 MPH Reduction	4	-100%	0.1%	41	-36.2%	2.2%	45	-36.4%	2.0%
Total	65	-4.3%	0.5%	286	-10.3%	0.7%	351	-8.9%	0.7%

Among the ten crash cases and 216 non-crash cases where the VSL was activated, improved safety results were observed. The VSL with 10 MPH reduction was activated for six crash cases and 29 non-crash cases; and the VSL with 20 MPH reduction was activated for four crash cases and 41 non-crash cases. On average, the crash prevention module reduced crash risk by 8.9% and increased mobility by 0.7% for all 351 cases. The proposed crash prevention module seems promising in improving safety without compromising mobility.

When considering only the cases where the VSL was activated, the average decrease in the crash risk is 26.3% for 10 MPH reduction and 36.4% for 20 MPH reduction, respectively. In a field evaluation on the VSL system implemented on Interstate 5 (I-5) in Washington, notable safety impact of the VSL has been observed (Pu et al., 2017). The observational before-after EB study suggested that the VSL system yielded 29% reduction in total crashes with a standard deviation of 5%. The findings from the field data support that the proposed method based on simulated traffic data is a viable alternative to safety analysis and evaluation.

7. Conclusions

Conventional real-time freeway crash prediction models identify crash-prone traffic conditions based on live feeds from loop detectors. It is common practice to use traffic data from the 5-10-min period prior to a crash, as this ensures sufficient time for taking the proper precautions. However, the phenomenon of time proximity suggests that traffic conditions occurring within the 0-5-min period of a crash are more relevant when it comes to predicting crashes. Moreover, a crash can happen between two detector stations where traffic information is not available, and the actual traffic conditions at the crash site may deviate from those captured by loop detector stations. Therefore, crash patterns derived from loop detector locations, as opposed to crash locations, are inadequate in accounting for varying distances between crashes and detectors. CTM-simulated traffic data were introduced in this study to fill the spatial and temporal gaps inherent in the observed traffic data collected from physical loop detector stations. Based on the traffic flow theory, CTM can predict traffic conditions anywhere at any time from its virtual detectors.

A real-time crash prediction model was developed with data from a corridor of I-94 in Wisconsin. The corridor was divided into a series of cells to create a uniform and close layout of virtual detector stations. Traffic data simulated from virtual upstream and downstream stations with consistent spacings was used for crash modeling to account for the spatial gap in physical station data. The simulated traffic data in the 0-5-min period prior to the crash/non-crash were used for crash modeling, and the traffic in the future 5-min period were simulated for crash prediction. In this way, the temporal issue of physical station data was also taken into consideration.

Simulated traffic data collected from one virtual upstream station and one virtual downstream station were used for crash modeling. The modeling results showed that varying variables are significantly related to the crash occurrence in different traffic states. Observed traffic data collected from physical stations were also employed for crash modeling. The prediction performance of the two crash prediction models was compared, showing that the simulated traffic data would improve

prediction performance by accounting for the spatial-tempo issue of physical station data.

A crash prediction and prevention method based on simulated traffic data was proposed to detect crash-prone conditions and help select the desirable TCS for crash prevention. The proposed method was tested in a case study with VSL strategies for demonstration, and results showed that the proposed crash prediction and prevention method could effectively detect crash-prone conditions and evaluate the safety and mobility impacts of various TCS alternatives before their deployment.

In future studies, a lane-specific CTM can be developed to provide simulated traffic on a lane-by-lane basis, therefore advancing the crash prediction performance with lane-specific traffic data. Future studies should also test other traffic control strategies, such as more flexible VSL control strategies and ramp metering.

References

- Abdel-Aty, M., Pande, A., 2006. Atms implementation system for identifying traffic conditions leading to potential crashes. *IEEE Trans. Intell. Transp. Syst.* 7 (1), 78–91.
- Abdel-Aty, M., Uddin, N., Pande, A., Abdalla, F., Hsia, L., 2004. Predicting freeway crashes from loop detector data by matched case-control logistic regression. *Transp. Res. Rec.: J. Transp. Res. Board* 1897, 88–95.
- Abdel-Aty, M., Uddin, N., Pande, A., 2005. Split models for predicting multivehicle crashes during high-speed and low-speed operating conditions on freeways. *Transp. Res. Rec.: J. Transp. Res. Board* 1908, 51–58.
- Abdel-Aty, M., Dilmore, J., Dhindsa, A., 2006. Evaluation of variable speed limits for real-time freeway safety improvement. *Accid. Anal. Prev.* 38 (2), 335–345.
- Abdel-Aty, M.A., Hassan, H.M., Ahmed, M., Al-Ghamdi, A.S., 2012. Real-time prediction of visibility related crashes. *Transp. Res. Part C: Emerg. Technol.* 24, 288–298.
- Cassidy, M.J., Rudjanakanoknad, J., 2005. Increasing the capacity of an isolated merge by metering its on-ramp. *Transp. Res. Part B: Methodol.* 39 (10), 896–913.
- Clark, S., 2003. Traffic prediction using multivariate nonparametric regression. *J. Transp. Eng.* 129 (2), 161–168.
- Courant, R., Friedrichs, K., Lewy, H., 1967. On the partial difference equations of mathematical physics. *IBM J. Res. Dev.* 11 (2), 215–234.
- Daganzo, C.F., 1994. The cell transmission model: network traffic. *Transp. Res. Part B-Methodol.* 29 (2), 79–93.
- Danczyk, A., Liu, H.X., 2011. A mixed-integer linear program for optimizing sensor locations along freeway corridors. *Transp. Res. Part B: Methodol.* 45 (1), 208–217.
- Dervisoglu, G., Gomes, G., Kwon, J., Horowitz, R., Varaiya, P., 2009. Automatic calibration of the fundamental diagram and empirical observations on capacity. In: *Proceedings of 88th Transportation Research Board Annual Meeting*. Washington, DC.
- Golob, T.F., Recker, W.W., 2003. Relationships among urban freeway accidents, traffic flow, weather, and lighting conditions. *J. Transp. Eng. ASCE* 129 (4), 342–353.
- Habtemichael, F.G., Cetin, M., 2016. Short-term traffic flow rate forecasting based on identifying similar traffic patterns. *Transp. Res. Part C: Emerg. Technol.* 66, 61–78.
- Hadiuzzaman, M., Qiu, T.Z., 2013. Cell transmission model based variable speed limit control for freeways. *Can. J. Civ. Eng.* 40 (1), 46–56.
- Hall, F.L., Agyemang-Duah, K., 1991. Freeway capacity drop and the definition of capacity. *Transp. Res. Rec.* 1320, 91–98.
- Hirunyanitwattana, W., Mattingly, S.P., 2006. Identifying secondary crash characteristics for California highway system. In: *Proceedings of the Transportation Research Board 85th Annual Meeting*. Washington, DC.
- Hong, Z., Fukuda, D., 2012. Effects of traffic sensor location on traffic state estimation. *Procedia-Soc. Behav. Sci.* 54, 1186–1196.
- Hossain, M., Muromachi, Y., 2012. A Bayesian network based framework for real-time crash prediction on the basic freeway segments of urban expressways. *Accid. Anal. Prev.* 45, 373–381.
- Kockelman, K.K., Ma, J., 2004. Freeway speeds and speed variations preceding crashes, within and across lanes. In: *Proceedings of 83rd Transportation Research Board Annual Meeting*. Washington, DC.
- Kwon, J., Petty, K., Varaiya, P., 2007. Probe vehicle runs or loop detectors?: Effect of detector spacing and sample size on accuracy of freeway congestion monitoring. *Transp. Res. Rec.: J. Transp. Res. Board* 2012, 57–63.
- Lee, C., Abdel-Aty, M., 2008. Testing effects of warning messages and variable speed limits on driver behavior using driving simulator. *Transp. Res. Rec.: J. Transp. Res. Board* 2069, 55–64.
- Lee, C., Hellinga, B., Saccomanno, F., 2006. Evaluation of variable speed limits to improve traffic safety. *Transp. Res. Part C: Emerg. Technol.* 14 (3), 213–228.
- Li, Z., Li, Y., Liu, P., Wang, W., Xu, C., 2014a. Development of a variable speed limit strategy to reduce secondary collision risks during inclement weathers. *Accid. Anal. Prev.* 72, 134–145.
- Li, Z.B., Liu, P., Wang, W., Xu, C.C., 2014b. Development of a control strategy of variable speed limits to reduce rear-end collision risks near freeway recurrent bottlenecks. *IEEE Trans. Intell. Transp. Syst.* 15 (2), 866–877.
- Li, Z., Liu, P., Xu, C., Wang, W., 2016. Optimal mainline variable speed limit control to improve safety on large-scale freeway segments: optimal mainline variable speed limit. *Comput.-Aided Civ. Infrastruct. Eng.* 31 (5), 366–380.
- Liu, H.X., Danczyk, A., 2009. Optimal sensor locations for freeway bottleneck identification. *Comput.-Aided Civil Infrastruct. Eng.* 24 (8), 535–550.
- Muñoz, L., Sun, X., Horowitz, R., Alvarez, L., 2003. Traffic density estimation with the cell transmission model. In: *Proceedings of the 2003 American Control Conference*. Denver, CO. pp. 3750–3755.
- Muñoz, L., Sun, X., Horowitz, R., Alvarez, L., 2006. Piecewise-linearized cell transmission model and parameter calibration methodology. *Transp. Res. Rec.: J. Transp. Res. Board* 1965, 183–191.
- Oswald, R.K., Scherer, W.T., Smith, B.L., 2001. *Traffic Flow Forecasting Using Approximate Nearest Neighbor Nonparametric Regression*. Center for Transportation Studies, University of Virginia.
- Pande, A., Abdel-Aty, M., 2006. Comprehensive analysis of the relationship between real-time traffic surveillance data and rear-end crashes on freeways. *Transp. Res. Rec.: J. Transp. Res. Board* 1953, 31–40.
- Parker, S.T., Tao, Y., 2006. Wistransportal: a wisconsin traffic operations data hub. In: *Proceedings of the 9th International Conference on Applications of Advanced Technology in Transportation*. Chicago, Ill.
- Smith, B.L., Williams, B.M., Oswald, R.K., 2002. Comparison of parametric and non-parametric models for traffic flow forecasting. *Transp. Res. Part C: Emerg. Technol.* 10 (4), 303–321.
- Sumalee, A., Zhong, R.X., Pan, T.L., Szeto, W.Y., 2011. Stochastic cell transmission model (SCTM): a stochastic dynamic traffic model for traffic state surveillance and assignment. *Transp. Res. Part B: Methodol.* 45 (3), 507–533.
- Sun, J., Sun, J., 2015. A dynamic Bayesian network model for real-time crash prediction using traffic speed conditions data. *Transp. Res. Part C: Emerg. Technol.* 54, 176–186.
- Transportation Research Board, 2010. *HCM 2010: Highway Capacity Manual*. Transportation Research Board, Washington, D.C.
- Xu, C., Liu, P., Wang, W., Li, Z., 2012. Evaluation of the impacts of traffic states on crash risks on freeways. *Accid. Anal. Prev.* 47, 162–171.
- Xu, C., Tarko, A.P., Wang, W., Liu, P., 2013a. Predicting crash likelihood and severity on freeways with real-time loop detector data. *Accid. Anal. Prev.* 57, 30–39.
- Xu, C., Wang, W., Liu, P., 2013b. A genetic programming model for real-time crash prediction on freeways. *IEEE Trans. Intell. Transp. Syst.* 14 (2), 574–586.
- Xu, C., Liu, P., Wang, W., Li, Z., 2014. Identification of freeway crash-prone traffic conditions for traffic flow at different levels of service. *Transp. Res. Part A: Policy Pract.* 69, 58–70.
- Xu, C., Liu, P., Wang, W., 2016. Evaluation of the predictability of real-time crash risk models. *Accid. Anal. Prev.* 94, 207–215.
- Yeo, H., Jang, K., Skabardonis, A., Kang, S., 2013. Impact of traffic states on freeway crash involvement rates. *Accid. Anal. Prev.* 50, 713–723.
- Zheng, Z., Ahn, S., Monsere, C.M., 2010. Impact of traffic oscillations on freeway crash occurrences. *Accid. Anal. Prev.* 42 (2), 626–636.