The Potential for Real-Time Traffic Crash Prediction

A KEY TO IMPROVING
FREEWAY SAFETY IS THE
ABILITY TO PREDICT
CRASH OCCURRENCE.
THIS FEATURE ADDRESSES
THE PREDICTION OF
CRASH POTENTIAL USING
REAL-TIME LOOP
DETECTOR DATA.
HISTORICAL CRASHES AND
CORRESPONDING
ARCHIVED DATA FROM
LOOP DETECTOR STATIONS
SURROUNDING CRASH
LOCATIONS WERE USED.

INTRODUCTION

The development of freeway crash prediction models using intelligent transportation systems (ITS) archived data could be a substantial advancement in the field of real-time traffic management. Such models not only are expected to improve safety but also may go a long way to improve freeway operations by reducing incident-related congestion.

Because there is a need to use real-time traffic data emanating from loop detectors, the approach differs distinctly from previous studies estimating crash frequencies or rates on a certain freeway section through aggregate measures of flow (such as average daily traffic or hourly volumes).

Although the authors try to establish a relationship between the patterns in precrash data from detectors surrounding the crash location, it is imperative that the time of the historical crashes is known with precision.

This feature proposes a shockwave and rule-based methodology to estimate the time of the crash and then identifies how much time and distance ahead of crash occurrence loop data may be used to predict the impending hazard. The final objective is to predict the possibility of crashes on freeways using real-time loop data.

BACKGROUND

Hughes and Council were the first researchers to explore the relationship between freeway safety and peak period operations using loop detector data. Not only did they indicate that "traffic flow

> consistency" as perceived by the driver may be an important factor in freeway

safety, they also expressed a need to determine the time of the historical crashes accurately to avoid the "cause and effect" fallacy, which might identify some freeway conditions as crash-prone when they actually may be the result of a crash.

Since then, the issue of estimating the time of the crash has been raised in quite a few studies with similar objectives.^{2,3} Although the importance of the issue has been identified, it has not been addressed thoroughly. Lee, Saccomanno and Hellinga carried out visual analyses of speed profiles at surrounding loop detector stations for 234 crashes to determine the actual time of crash occurrence. In a later study, the authors came up with a more systematic approach.⁴

It was argued that the time of the crash may be approximated by the time the shockwave (of backward forming type) hits the loop detector station located immediately upstream of the crash site. It was justified based on the assumption that the shockwave speed on urban freeways is 20 kilometers per hour (12 miles per hour) and, therefore, very low expected errors are involved in the approximation. The problem with such an assumption is that slightly lower shockwave speeds will cause the errors to inflate and the approximations to become questionable.

Research aimed at freeway crash prediction through loop data also was carried out by Oh, Oh, Ritchie and Chang and Golob and Recker. ^{5,6} The data used in these studies were obtained from just one station downstream and/or upstream of the crash location. None of these studies looked at the "progression" of alarming driving conditions with the flow of traffic by analyzing data from a series of stations surrounding the crash location at several time periods leading to the crash.

Despite these shortcomings, the main contribution of these studies is that they demonstrated the possibility of determining crash potential at a certain freeway location (or section) in real time using data from upstream/downstream loop detectors.

This study presents a refined shockwave analysis approach toward determining the time of historical crashes. By analyzing the data from a series of detectors at different

BY MOHAMED ABDEL-ATY, PH.D., P.E., ANURAG PANDE, LIANG Y. HSIA, P.E. AND FATHY ABDALLA, PH.D.

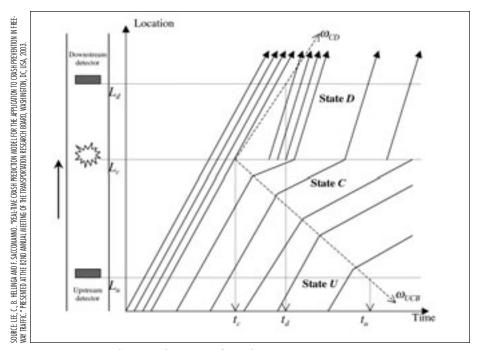


Figure 1. Time-space diagram in the presence of a crash.

time increments, it also examines the possibility that alarming crash-prone conditions on a freeway actually might originate upstream/downstream of the crash location early and "travel" with traffic until they culminate a crash at a certain time.

STUDY AREA AND AVAILABLE DATA

Data from dual loop detectors on Interstate 4 (I-4) in the Orlando, FL, USA, metropolitan area were used in this study. The following data were collected on I-4 every 30 seconds: average vehicle counts, average speed and lane detector occupancy. These data were collected for the three through lanes in both directions and at stations spaced at approximately one-half mile for a 36.25-mile stretch.

The crash data were collected from the Florida Department of Transportation (FDOT) crash database for the years 1999 to 2002. First, the location for all the crashes that occurred in the study area during this period was identified. For every crash, the loop detector station nearest to its location was determined and referred to as the station of the crash.

Because the first objective was to estimate the accurate time of the crashes, loop detector data from the station of the crash, four upstream stations and two downstream stations were collected for a period of 90 minutes around the

reported time of every crash (one hour prior and one-half hour later). For estimating the time of the crashes, data in time series of 30 seconds were used.

Loop detectors are known to suffer from intermittent hardware problems that result in unreasonable speed, volume and occupancy values. Values that included occupancy > 100; speed = 0 or > 100; flow > 25 and flow = 0 with speed > 0 were removed from the raw 30-second data.⁷

IMPACT OF CRASHES ON TRAFFIC FLOW

Crashes are a specific type of incident and generally have a more profound impact on freeway operations. The effects of a crash on traffic flow patterns develop over time both upstream and downstream of the crash. However, the changes in traffic flow characteristics are distinct on loop detectors located in upstream and downstream directions.

In the upstream direction, a queue can be observed, resulting in a significant reduction in lane speed and a significant increase in occupancy. On the other hand, a decrease in lane flow and occupancy can be observed downstream.

The critical aspect for determining the time of the crash is the time elapsed in the progression of the shockwave from the crash location to the upstream loop detector station. In general, this duration (the shockwave speed) and changes observed in the loop data are affected by the severity of the crash; the roadway geometry; the presence of on- and off-ramps; the distance between loop detector stations; and prevailing traffic flow conditions.⁸

The impact of a crash under the assumption of a constant shockwave speed may be shown by a time-space diagram (see Figure 1). L_d and L_u represent the location of detector stations downstream and upstream of the crash site, respectively. The times t_c , t_d and t_u are the time of the crash and the time of the shockwave arriving at downstream and upstream stations, respectively.

It is clear from Figure 1 that if the speed of the backward-forming shockwave is known, the time of the crash can be estimated easily. The times of the shockwave hitting two adjacent upstream stations may be determined by observing when the drops in speed profiles of the two stations occur. The gap between the two arrival times is the time that the backward-forming shockwave takes to travel from the first upstream station to the next upstream station.

TIME OF CRASH ESTIMATION

The first step in estimating the time of the crash was to estimate the speed of the backward-forming shockwave resulting from the crash. The difference between times of shockwave arrival at the two adjacent stations located immediately upstream of the crash location was used. Because the milepost of all loop detectors on I-4 was known accurately, the distance between the two detectors could be used to obtain the shockwave speed (a similar approach was attempted by Lee, Hellinga and Saccomanno).⁹

Once the shockwave speed is known, it is not difficult to determine t_c using the milepost of crash location (also known from the FDOT crash database). The following equation may be used for the estimation:

$$t_u - t_c = \frac{(L_u - L_c)}{\omega_{UC}} \tag{1}$$

All the variables in Equation 1 have the notation used in Figure 1. Due to the underlying assumption that shockwave speed remains constant while it hits the first and second stations in the upstream direction, it was mandatory to validate the results. The critical issue in the validation was that there is no way to know the actual time of the crash (true value) to compare the shockwave model estimates.

The model was validated using the traffic simulation package PARAMICS. A small freeway section on I-4 was simulated and three traffic flow statistics (speed, volume and density) were obtained from locations one-half mile apart on the section just as the loop data are archived for I-4 in real time.

Crashes were configured to occur at various locations between a set of two detectors (for example, very near to the upstream or downstream loop, exactly midway between the loops). The simulation experiment showed that the time of these "artificial" crashes could be accurately estimated using the shockwave method under various scenarios.

Aggregation Across Lanes Versus Using Lane of the Crash

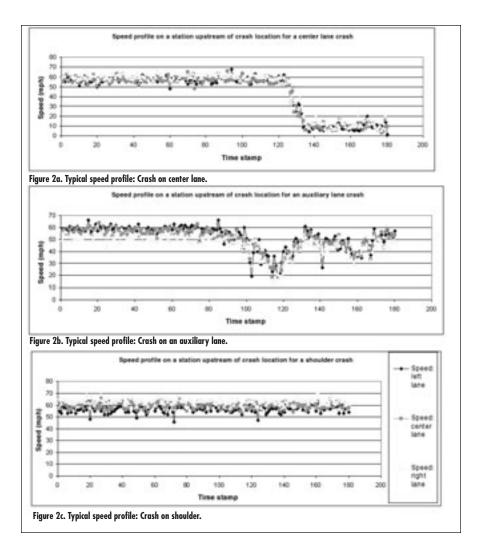
After the methodology was developed and validated as explained above, it could be applied either by aggregating the data across three lanes or by using the data from the specific lane on which the crash had occurred. The lane of the crash was known from the FDOT crash database.

The advantage of using the aggregated data is that the time of the crash could be estimated for a large sample of crashes because the data for at least one of the lanes obviously are available for more crashes than the data for a specific lane.

On the other hand, because the algorithm relies on the impact of the shock-wave hitting at successive upstream stations, sometimes the aggregated data (averaged over three lanes) might dampen this impact, and the drop in speed or rise in occupancy may not be significant enough to be detected by the algorithm as a shockwave hit. Therefore, it was decided to apply the algorithm for the specific lane of the crash for each case.

Crashes at Different Locations

Although the results of the above algorithm were validated on the simulation data, it was necessary to understand



some of the complexities involved before applying it to the real data (for example, for crashes that occur on the median, it is almost impossible to detect any effect on upstream loop detectors).

Because even the "rubberneck" effect dies down before being felt at the station immediately preceding the crash location, the algorithm was examined further and validated by looking at speed and occupancy profiles obtained at stations immediately upstream for randomly selected crashes. These crashes were selected from different roadway locations (such as the three mainstream lanes, median, shoulder, auxiliary lanes) to identify the lanes from a clear pattern of sudden drop in the loop detector speed data could be observed.

The visual inspection of several crash profiles from aforementioned roadway locations led to the formulation of the following rules:

• For crashes on the left, center, or right

- main traffic-stream lanes: Estimate the time of the crash by applying the existing methodology on the data from the respective lane (the lane of the crash).
- For the fourth (right-most) traffic lane or auxiliary lanes: Use time estimated by applying the existing methodology on the data from the right-most lane (lane three).
- Shoulder: No obvious pattern could be observed in the upstream loop data; therefore, it would not be appropriate to modify the time.
- Median: No obvious pattern could be observed in the upstream loop data; therefore, it would not be appropriate to modify the time.

The logic behind the formulation of the aforementioned rules may be understood through careful inspection of Figure 2. It also helps to visualize the trends observed in the speed patterns from the station upstream of three different crash locations. Note that these are the typical speed profiles and most of the other crashes on these roadway locations also depicted similar trends.

"Crash on center lane" (see Figure 2a) represents crashes on the mainstream freeway (lanes equipped with loop detectors). Figure 2b depicts the speed pattern for crashes on the fourth lane (auxiliary lane) on the freeway. Therefore, the impact of a crash occurring on this lane could be captured by observing the drop in speed on the adjacent lane (the rightmost lane equipped with loop detectors).

To represent crashes on the shoulder and median, a shoulder crash was chosen, which shows no visible speed drop pattern in any of the lanes equipped with loop detectors (Figure 2c). The time series shown in Figure 2 has readings obtained from three freeway lanes for a period of 90 minutes (one hour prior to and one-half hour later than the reported time of each crash). Out of these 180 readings, the 120th reading is the reported time of the crash.

After applying this methodology for all crashes having the desirable lane data available, the time of crash was modified accordingly. Due to the unavailability of specific lane data for the required loop stations and the time period for all the crash cases, the sample size was reduced to about one-fourth of the original crash cases. The final sample used in the analysis in the next section was 556 crashes for years 1999–2002 with complete loop detector data available.

USING TRAFFIC PARAMETERS TO PREDICT CRASHES

Methodology

The case-control stratum analysis methodology is adopted to identify the relationship between the traffic parameters measured through loop detectors and crash occurrences while controlling for location, time of day, day of week and season.

In a logistic regression setting, the function of dependent variables yielding a linear function of the independent variables would be the logit transformation.

$$(x) = \ln \left[\frac{\pi(x)}{1 - \pi(x)} \right] = \beta_0 + \beta_1 x$$
 (2)

where:

 $\pi(x) = E(Y|x)$ is the conditional mean of Y (dummy variable representing crash occurrence) given x when the logistic distribution is used. Under the assumption that the logit is linear in the continuous covariate x, the equation for the logit would be $g(x) = \beta_0 + \beta_1 x$. It follows that the slope coefficient, β_1 , gives the change in the log odds for an increase of one unit in x, for example, $\beta_1 = g(x+1) - g(x)$ for any value of x.

Hazard ratio is defined as the exponential of this coefficient, in other words, it represents how much more likely (or unlikely) it is for the outcome to be present for an increase of one unit in x. ¹⁰ It implies that hazard ratio significantly different from one for a particular parameter is an indicator of strong association of that parameter with crash occurrence. It also is noteworthy that a value greater than one signifies that crash risk increases with an increase in the parameter value; a value less than one indicates an increase in crash risk as parameter value decreases.

Data Preparation

Once the methodology for determining the time of the crashes was developed and applied, the data for matched case control analysis were prepared based on the refined/adjusted time of the crash. The methodology used for predicting crashes here is matched case-control logistic regression.

Therefore, if a crash occurred on April 12, 1999 (Monday) at 6:00 p.m. on I-4 eastbound and the nearest loop detector was at station 30, data were extracted from station 30, four loops upstream and two loops downstream of station 30 for one half-hour period prior to the estimated time of the crash for all the Mondays of the year at the same time.

This matched sample design controls for factors affecting overall traffic patterns, such as type of drivers on the freeway. Therefore, this crash will have a loop data table consisting of the speed, volume and occupancy values for all three lanes from loop stations 26–32 (on the east-bound direction) from 5:30 p.m. to 6:00 p.m. for all Mondays of 1999, with one of them being the day of crash.

Because the 30-second data have ran-

dom noise and are difficult to work with in a modeling framework, the 30-second data were combined into 5-minute levels to obtain average and standard deviations. The one-half hour period was divided into six time slices. All the stations were named from "B" to "G", with "B" being the farthest station upstream. It should be noted that "F" is the station of the crash; "G" and "H" are the stations downstream of the crash location.

Similarly, the 5-minute intervals were given IDs from 1 to 6. The interval between the time of the crash and 5 minutes prior to the crash was named slice 1; the interval 5 to 10 minutes prior to the crash was named slice 2; the interval 10 to 15 minutes prior to the crash was named slice 3.

Two effects, average and standard deviation, initially were calculated for speed, volume and occupancy during each time slice at every station. For example, the standard deviation and average speed on the left lane during the 5-minute slice just prior to a crash at the station of the crash would be named "SSLF1" and "ASLF1," respectively.

ANALYSIS AND RESULTS

For each of the seven loop detectors (B to G) and six time slices (1 to 6), there are values of means (AS, AV, AO) and standard deviations (SS, SV, SO) of speed, volume and occupancy for all crash and corresponding non-crash cases. Due to data availability, there were different numbers of non-crash cases for each crash. To carry out matched case-control analysis, a symmetric data set was created (each crash case in the dataset has the same number of non-crash cases as controls) by randomly selecting five non-crash cases for each crash.

Exploratory analysis with the original effects (5-minute standard deviations and average of speed) showed that the hazard ratios for standard deviation of speed all were greater than unity although they all were less than one for the average speeds at stations B–H and time slices 1–6. Therefore, the coefficient of variation in speed (standard deviation divided by average) was a natural choice as a precursor resulting in hazard ratio values substantially greater than one.

A similar logic was applied for volume. Therefore, mean and standard deviation of speed and volume were combined into the variables *CVS* (coefficient of variation of speed) and *CVV* (coefficient of variation of volume), expressed in percentage as (*SS/AS*)*100 and (*SV/AV*)*100.

For example, CVS at station F (the crash location) at time slice 3 (10–15 minutes before the crash) ranged from 0.2 to about 1.82, with the highest percentage of cases at about 0.7 and another peak at about 1.5. Possibly, this indicates two common situations leading to crashes on freeways that involve high variation in speed with low or high average speed.

With five variables (CVV, CVS, AV, SV and AO) at each of the seven loop detectors and six time slices, a stratified conditional simple (one variable at a time) logistic regression analysis was conducted to identify time duration(s) and location of loop detector(s) whose traffic characteristics are significantly correlated with the binary outcome (crash versus non-crash).

This was done by calculating the hazard ratio using proportional hazard regression analysis (*PHREG* of *SAS*) of each of the 210 single variable models; one model for each of the five variables over every station B–H and time slice 1–6. The outcome of these models was the hazard ratio value for these variables at various stations and time slices. The p-value for the test indicates whether the value is significantly different from one.

The hazard ratio is an estimate of the expected change in the risk ratio of having a crash. Therefore, if the output hazard ratio of a variable is significantly different from one (for example, three), increasing the value of this variable by one unit would increase the risk of a crash at station F (station of the crash) by three times. The initial analysis concluded that the variables *CVV* and *CVS* had the most significant hazard ratios, but this was not the case for all the time slices and stations.

Figure 3 and Table 1 show the hazard ratio values for CVS. Note that the ratio increases as the space dimension moves from station B (farthest from the crash site) to station F and then drops slightly at the downstream stations (G and H). By comparing among the six 5-minute

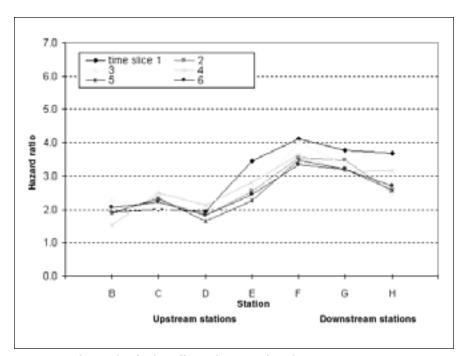


Figure 3. Hazard ratio values for the coefficient of variation of speed.

Table 1. Hazard ratio values corresponding to the coefficient of variation of speed.								
	Station							
	Upstream stations					Downstream stations		
	В	C	D	E	F	G	Н	
Time slice 1	1.923	2.023	1.947	3.446	4.122	3.779	3.676	
Time slice 2	1.920	2.297	1.830	2.530	3.549	3.480	2.525	
Time slice 3	2.268	2.054	1.728	2.347	4.035	3.341	2.437	
Time slice 4	1.524	2.498	2.125	2.820	3.638	3.137	3.158	
Time slice 5	1.898	2.329	1.642	2.273	3.483	3.194	2.593	
Time slice 6	2.054	2.226	1.829	2.449	3.329	3.214	2.709	

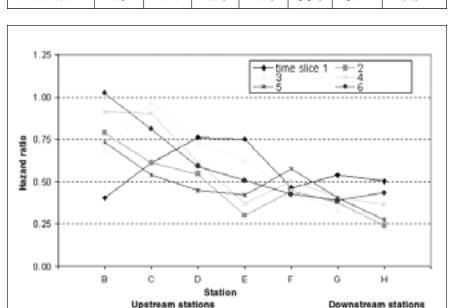


Figure 4. Hazard ratio values for the coefficient of variation of volume.

time slices, at stations F and G there are no significant differences in the hazard ratios of *CVS* between the different time slices (at least the first three time slices). This indicates that an adverse traffic condition before a crash remains for some time. Therefore, this time could be used to anticipate and, hopefully, try to prevent the crash.

Figure 4 illustrates the hazard ratios for CVV. A dropping trend in the hazard ratio of CVV may be observed as the site of the crash approaches. A lower hazard ratio at the station of the crash and beyond indicates that the CVV is significantly correlated with crash occurrence and may be used to detect crashes.

In short, a higher CVS and lower CVV increase the likelihood of crashes. Although this trend is observed starting about 2 miles upstream of the crash location, it is considerably clear about one-half mile upstream and also downstream.

It also is clear that the "ingredients" for a crash start about 15 minutes before the crash. The CVS factor represents high variation in speed relative to the average speed and, surprisingly, the CVV factor represents low variation in volume relative to the average volume. Other factors (AV, SV and AO) were tested but were not found to be significant.

DISCUSSION

A high CVS value has been identified as one of the crash causes on freeways. CVS is defined as the standard deviation divided by the average speed over a 5-minute interval. Lower speed associated with high variance (leading to a high value of coefficient of variation) depicts frequent formation of queues followed by their quick dissipation.

Figure 3 shows a rise in hazard ratio (risk) as the station where the crash occurred approaches. The hazard value is particularly high for station F, but also high for stations E (one-half mile upstream) and G and H (one-half to one mile downstream). Time slice 1 has the highest hazard ratio (0 to 5 minutes before crash occurrence). However, time slice 3 (10–15 minutes before the crash) is particularly high at station F. This figure illustrates a spatial and temporal dimension of increased hazard values

and, therefore, the ability to detect a rising trend at and around certain locations indicating a possibility for safety problems 10–15 minutes later, which would provide enough time to prepare for an impending risk.

Table 1 shows the hazard ratio values for all coefficients of variation. As explained earlier, the value of the hazard ratio signifies the resulting change in odds of observing a crash when the value of a certain parameter is changed by one unit.

For insight into how the crash risk varies approaching the crash location, consider an interval of 10–15 minutes (slice 3) prior to crash occurrence. During this interval, increasing the coefficient of variation by one unit at a location about one-half mile upstream (station E) would increase the odds of a crash by 2.347 times. At the same time, a similar change in *CVS* at the location of the crash (station F) and at a location one-half mile downstream (station G) would increase crash risk by factors of 4.035 and 3.341, respectively.

From an application point of view, if an increasing variation in speed is observed at a certain loop detector station, freeway sections in the vicinity of this station and about one-half mile upstream of it are more likely to experience a crash than any of the downstream sections.

The other factor, the low value of CVV, indicates that high traffic flow with low variability in volume is positively correlated with crash occurrences on freeways. Figure 4 depicts a drop in hazard ratio (risk) as the station where the crash occurred approaches. The hazard value is particularly low for station F and the stations downstream (G and H). A possible interpretation of this criterion might be that in case of high variability in volume, the density changes and, consequently, the gaps between vehicles change, which alerts drivers.

On the other hand, in case of low variability in volume, the density and the gap remain almost fixed in the traffic stream, which causes drivers to relax, thus slowing their reaction time. It also could be that low variability of volume and high traffic flow might sometimes be associated with queues. Queue formation

and shockwaves are a common cause of rear-end crashes on freeways.

CONCLUSIONS

This feature presents a simple statistical approach to predict crashes based on realtime data. The case control logistic regression was used with loop detector data to detect traffic patterns that could produce a high crash potential. Even if the first time slice (0-5 minutes prior to the crash) is excluded due to practical considerations of the time required to act on the information and warn drivers, it was shown that crash-prone conditions in terms of high coefficient of variation in speed and low coefficient of variation in volume are not ephemeral on freeway sections. The hazard ratio values for these variables were significantly different from one around the crash location for three to four time slices (they existed for about 15 minutes), which should provide enough time for prediction (and prevention) of crashes.

This study demonstrated the applicability of loop detector data for predicting freeway crashes. Once a potential crash location is identified in real time, measures for reducing the speed variance may be implemented to reduce the risk. For example, warning messages could be displayed on variable message signs, or strategies to calm speed using variable speed limit techniques could be adopted. However, real-time application still needs thorough investigation.

References

- 1. Hughes, R. and F. Council. "On Establishing Relationship(s) Between Freeway Safety and Peak Period Operations: Performance Measurement and Methodological Considerations." Presented at the 78th Annual Meeting of the Transportation Research Board (TRB), Washington, DC, USA, 1999.
- 2. Lee, C., F. Saccomanno and B. Hellinga. "Analysis of Crash Precursors on Instrumented Freeways." *Transportation Research Record*, No. 1784 (2002): 1–8.
- 3. Lee, C., B. Hellinga and F. Saccomanno. "Real-Time Crash Prediction Model for the Application to Crash Prevention in Freeway Traffic." Presented at the 82nd Annual Meeting of TRB, Washington, DC, 2003.
 - 4. Ibid.
 - 5. Oh, C., J. Oh, S. Ritchie and M. Chang.

"Real Time Estimation of Freeway Accident Likelihood." Presented at the 80th Annual Meeting of TRB, Washington, DC, 2001.

6. Golob, T. F. and W.W. Recker. "Relationships Among Urban Freeway Accidents, Traffic Flow, Weather and Lighting Conditions." *California PATH Working Paper UCB-ITS-PWP-2001-19, Institute of Transportation Studies*. University of California, Berkeley, 2001.

7. Chandra C. and H. Al-Deek. "New Algorithms for Filtering and Imputation of Real Time and Archived Dual-Loop Detector Data in the I-4 Data Warehouse." Presented at the 83rd Annual Meeting of TRB, Washington, DC, 2004.

8. Adeli, H. and A. Karim. "Fuzzy-Wavelet RBFNN Model for Freeway Incident Detection." *Journal of Transportation Engineering*, Vol. 126, No. 6 (2000): 464–471.

9. Lee, Hellinga and Saccomanno, note 3 above.

10. Agresti, A. Categorical Data Analysis, 2nd Edition. John Wiley and Sons Inc., 2002.



MOHAMED ABDEL-ATY,

Ph.D., P.E., is an associate professor of civil engineering at the University of Central Florida (UCF). His main research interests

are in traffic safety, travel demand analysis and ITS. He has published more than 90 papers. He is a member of two TRB committees and the Editorial Advisory Boards of Accident Analysis and Prevention and the ITS Journal. He is the 2003 UCF Distinguished Researcher.



ANURAG PANDE

is a Ph.D. candidate in the Department of Civil and Environmental Engineering at UCF. His research interests include traffic safety analysis, ITS and statis-

tical and data mining applications in transportation engineering. He is a student affiliate of the UCF ITE student chapter.



LIANG Y. HSIA,

P.E., is an engineer administrator with the Florida Department of Transportation. He manages the statewide ITS architecture, ITS, transportation manage-

ment center software, standards and research programs. He received his bachelor of science in architectural engineering from the Chinese Culture University and a master of science in construction from the University of Florida with additional graduate studies in computer science, structural engineering and traffic engineering. He is a member of ITE and ASCE.

FATHY ABDALLA,

Ph.D., obtained his Ph.D. from the University of Central Florida in 2003.

ITE JOURNAL ON THE WEB / DECEMBER 2005