# MIMIC—Multidisciplinary Initiative on Methods to Integrate and Create Realistic Artificial Data

U.S. Department of Transportation
**Federal Highway Administration**

# FOREWORD

Data-driven safety analysis models help State and local agencies quantify safety data, identify high-risk roadway features, and predict the effects of proposed safety measures. However, even when a model performs well overall, it may not accurately represent the interactions between variables for a specific location or crash because the underlying relationships in the real world are unknown. One proposed solution is to generate realistic artificial datasets (RADs) with predetermined safety relationships built into them. Because these relationships are known, the RAD can serve as a testbed, revealing how well a model reflects those underlying cause-and-effect relationships.

This study describes the development of RAD for ramp terminals and speed change lanes at diamond interchanges. A web-based software was developed under the Federal Highway Administration's Exploratory Advanced Research Program. The software provides the ability to generate RAD for multiple years and locations as well as access to pregenerated datasets. This report will be of interest to academics and researchers developing crash modification functions and statistical models to determine how the models best represent real-world relationships.

Brian P. Cronin, P.E.
Director, Office of Safety and Operations.
Research and Development

# TECHNICAL REPORT DOCUMENTATION PAGE

| 1. Report No.<br>FHWA-HRT-23-015 | 2. Government Accession No. | 3. Recipient's Catalog No. | |
|---|---|---|---|
| 4. Title and Subtitle<br>MIMIC—Multidisciplinary Initiative on Methods to Integrate and Create Realistic Artificial Data. | | 5. Report Date<br>January 2023 | |
| | | 6. Performing Organization Code | |
| 7. Author(s)<br>Edara, P. (0000-0003-2707-642X), Sun, C. (0000-0002-8857-9648), Brown, H. (0000-0003-1473-901X), Savolainen, P. (0000-0001-5767-9104), Shankar, V. (0000-0002-6671-2268), Balakrishnan, B. (0000-0002-0994-0213), Shang, Y. (0000-0001-7771-4034), Chakraborty, S. (0000-0003-2022-1735), Adu-Gyamfi, Y. (0000-0002-1924-9792), Li, C. (0000-0002-3237-1477), Aati, K. (0000-0001-8834-7735), Lima, S., Huang, Y. (0000-0002-7346-5293), Mussah, A. (0000-0002-1084-5598), Hopfenblatt, J. | | 8. Performing Organization Report No. | |
| 9. Performing Organization Name and Address<br>University of Missouri-Columbia<br>E2509 Lafferre Hall<br>Columbia, MO 65211 | | 10. Work Unit No. (TRAIS) | |
| | | 11. Contract or Grant No.<br><br>693JJ31950023 | |
| 12. Sponsoring Organization Name and Address<br>United States Department of Transportation<br>Federal Highway Administration<br>HSA Room #E71-324<br>1200 New Jersey Avenue SE<br>Washington, DC 20590 | | 13. Type of Report and Period Covered<br>Final Report; September 2019–October 2022 | |
| | | 14. Sponsoring Agency Code<br>HRSO-2 | |
| 15. Supplementary Notes<br>This project was supported by the Exploratory Advanced Research (EAR) Program. EAR Program oversight was provided by David Kuehn and Jim Shurbutt. Yusuf Mohamedshah served as Federal Highway Administration (FHWA) Contracting Officer's Technical Manager. Additional project guidance was provided by FHWA staff members Carol Tan and Ana Maria Eigen. The project's technical advisory committee consisted of Dean Kanitz (Michigan Department of Transportation (DOT)), Ida van Schalkwyk (Washington State DOT), Ray Shank (Missouri DOT), and John Miller (FHWA Missouri Division). | | | |

16. Abstract
Traditional safety modeling efforts primarily focus on accurately estimating crash frequencies or rates. The true relationships between crashes and potential causal factors are not always easily discernible from safety models. While a model consisting of multiple causal factors may produce accurate estimates of crash measures, it may not accurately explain all causal relationships. Knowing the true cause-and-effect relationships is important while choosing countermeasures to address safety problems. This Exploratory Advanced Research Program project developed a framework to generate realistic artificial data (RAD) datasets that mimic the known causal relationships between contributing factors and crashes. The proposed framework is generic and can be used to generate RAD for other facilities, such as work zones, bicycle/pedestrian facilities, innovative geometric designs, etc. The framework was applied to generate RAD for ramp terminals and speed change lane facilities at diamond interchanges. A web-based software was developed to provide easy access to the RAD dataset. The software provides 196 pregenerated datasets and the option to request custom datasets. Sample RAD datasets were used to test negative binomial and a suite of machine learning models. A model evaluation rubric was developed to evaluate and compare the performance of different models. Additionally, this project developed a second type of RAD dataset—the virtual reality (VR) simulation testbeds for crashes and near-crashes occurring at interchanges. Driving simulator studies offer another source of RAD for evaluating new behavioral and roadway countermeasures. The testbeds were developed using safety critical events recorded in the Strategic Highway Research Program 2 Naturalistic Driving Study data. VR offers an engaging visualization platform to educate the public about interchange crashes and to evaluate different countermeasures. These interventions are well aligned with the USDOT's National Roadway Safety Strategy's Safe System Approach of considering an overlapping set of safety measures—roadway countermeasures, behavioral interventions, enforcement, vehicle safety features, and emergency medical care—to achieve zero roadway fatalities.

| 17. Key Words<br>Crashes, realistic artificial data, safety, synthetic data, visualization, simulator testbeds | | 18. Distribution Statement<br>No restrictions. This document is available to the public through the National Technical Information Service, Springfield, VA 22161.<br>http://www.ntis.gov | |
| 19. Security Classification (of this report)<br>Unclassified. | 20. Security Classification (of this page)<br>Unclassified. | 21. No. of Pages<br>69 | 22. Price<br>NA |

**Form DOT F 1700.7 (8-72)**          Reproduction of completed page authorized

# SI* (MODERN METRIC) CONVERSION FACTORS

## APPROXIMATE CONVERSIONS TO SI UNITS

| Symbol | When You Know | Multiply By | To Find | Symbol |
|---|---|---|---|---|
| **LENGTH** | | | | |
| in | inches | 25.4 | millimeters | mm |
| ft | feet | 0.305 | meters | m |
| yd | yards | 0.914 | meters | m |
| mi | miles | 1.61 | kilometers | km |
| **AREA** | | | | |
| $in^2$ | square inches | 645.2 | square millimeters | $mm^2$ |
| $ft^2$ | square feet | 0.093 | square meters | $m^2$ |
| $yd^2$ | square yard | 0.836 | square meters | $m^2$ |
| ac | acres | 0.405 | hectares | ha |
| $mi^2$ | square miles | 2.59 | square kilometers | $km^2$ |
| **VOLUME** | | | | |
| fl oz | fluid ounces | 29.57 | milliliters | mL |
| gal | gallons | 3.785 | liters | L |
| $ft^3$ | cubic feet | 0.028 | cubic meters | $m^3$ |
| $yd^3$ | cubic yards | 0.765 | cubic meters | $m^3$ |
| NOTE: volumes greater than 1,000 L shall be shown in $m^3$ | | | | |
| **MASS** | | | | |
| oz | ounces | 28.35 | grams | g |
| lb | pounds | 0.454 | kilograms | kg |
| T | short tons (2,000 lb) | 0.907 | megagrams (or "metric ton") | Mg (or "t") |
| **TEMPERATURE (exact degrees)** | | | | |
| °F | Fahrenheit | 5 (F-32)/9 or (F-32)/1.8 | Celsius | °C |
| **ILLUMINATION** | | | | |
| fc | foot-candles | 10.76 | lux | lx |
| fl | foot-Lamberts | 3.426 | candela/$m^2$ | cd/$m^2$ |
| **FORCE and PRESSURE or STRESS** | | | | |
| lbf | poundforce | 4.45 | newtons | N |
| lbf/$in^2$ | poundforce per square inch | 6.89 | kilopascals | kPa |

## APPROXIMATE CONVERSIONS FROM SI UNITS

| Symbol | When You Know | Multiply By | To Find | Symbol |
|---|---|---|---|---|
| **LENGTH** | | | | |
| mm | millimeters | 0.039 | inches | in |
| m | meters | 3.28 | feet | ft |
| m | meters | 1.09 | yards | yd |
| km | kilometers | 0.621 | miles | mi |
| **AREA** | | | | |
| $mm^2$ | square millimeters | 0.0016 | square inches | $in^2$ |
| $m^2$ | square meters | 10.764 | square feet | $ft^2$ |
| $m^2$ | square meters | 1.195 | square yards | $yd^2$ |
| ha | hectares | 2.47 | acres | ac |
| $km^2$ | square kilometers | 0.386 | square miles | $mi^2$ |
| **VOLUME** | | | | |
| mL | milliliters | 0.034 | fluid ounces | fl oz |
| L | liters | 0.264 | gallons | gal |
| $m^3$ | cubic meters | 35.314 | cubic feet | $ft^3$ |
| $m^3$ | cubic meters | 1.307 | cubic yards | $yd^3$ |
| **MASS** | | | | |
| g | grams | 0.035 | ounces | oz |
| kg | kilograms | 2.202 | pounds | lb |
| Mg (or "t") | megagrams (or "metric ton") | 1.103 | short tons (2,000 lb) | T |
| **TEMPERATURE (exact degrees)** | | | | |
| °C | Celsius | 1.8C+32 | Fahrenheit | °F |
| **ILLUMINATION** | | | | |
| lx | lux | 0.0929 | foot-candles | fc |
| cd/$m^2$ | candela/m2 | 0.2919 | foot-Lamberts | fl |
| **FORCE and PRESSURE or STRESS** | | | | |
| N | newtons | 2.225 | poundforce | lbf |
| kPa | kilopascals | 0.145 | poundforce per square inch | lbf/$in^2$ |

*SI is the symbol for International System of Units. Appropriate rounding should be made to comply with Section 4 of ASTM E380. (Revised March 2003)

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| 3D | three dimensional |
| AADT | annual average daily traffic |
| AASHTO | American Association of State Highway and Transportation Officials |
| ADT | average daily traffic |
| AIC | Akaike information criterion |
| ARD | artificial realistic data |
| BIC | Bayesian information criterion |
| CMF | crash modification factor |
| CNN | convolutional neural network |
| DDSA | data-driven safety analysis |
| DOT | department of transportation |
| EDC | Every Day Counts |
| FHWA | Federal Highway Administration |
| EAR | Exploratory Advanced Research |
| GAN | generative adversarial networks |
| HSIS | Highway Safety Information System |
| HSM | *Highway Safety Manual* |
| KNN | k-nearest neighbor |
| LT | left turn |
| MSE | mean squared error |
| NCHRP | National Cooperative Highway Research Program |
| NDS | Naturalistic Driving Study |
| PDO | property damage only |
| RAD | realistic artificial data |
| RSE | relative squared error |
| SCL | speed change lane |
| SHRP | Strategic Highway Research Program |
| SPF | safety performance function |
| SVM | support vector machine |
| USDOT | U.S. Department of Transportation |
| VR | virtual reality |

# EXECUTIVE SUMMARY

Data-driven methods are an important component of transportation safety decisionmaking. The U.S. Department of Transportation's (USDOT) Strategic Plan (2022–2026) stresses the importance of using data-driven methods as part of the overall Safe System Approach toward achieving zero roadway fatalities.[1,2] These methods typically require analytical evaluation of predicted and expected crashes based on geometric and traffic characteristics and other contributing factors. One tool that could facilitate this evaluation is realistic artificial data (RAD).[3,1] RAD can be beneficial to highway safety research in several ways: assessing a new crash estimation method, comparing methods to analyze alternatives, and conducting human factors evaluation of behavioral and roadway countermeasures. Advancing RAD will also enhance the Federal Highway Administration's (FHWA) efforts to encourage practitioners to apply data-driven methods to safety decisionmaking through programs such as Every Day Counts by expanding the number of tools available for safety analysis.[4]

Although artificial data have been used in many diverse applications, such as security, image processing, surveys, cancer genomics, infrared spectroscopy, and geography, their use in transportation has been limited. The main goal of this Exploratory Advanced Research (EAR) Program project was to develop a framework to generate RAD for interchange facilities and to generate datasets using that framework. Even though interchanges are ubiquitous in our transportation network and carry significant traffic volumes, accurate crash data for such facilities are lacking throughout the United States. The scope of this project involves generating RAD for two types of crashes occurring at diamond interchanges—ramp terminal left-turn (LT) crashes and speed change lane (SCL) crashes.

The data generation framework consists of three main steps. The first step identifies a set of contributing factors at the selected interchange facility (e.g., SCL, ramp terminal). Roadway, traffic, and driver contributing factors were synthesized from the literature from each selected facility. Sampling distributions were generated for each of the factors using observed data from Washington and Missouri. The RAD for these factors were then generated by repeatedly sampling the distributions for a given sample size (e.g., 500 sites). Data from other States were also considered. Highway Safety Information System data for interchange crashes were obtained for a 5-yr period. Data from Washington were the most complete and recent for the purposes of developing RAD, although data from California, Illinois, Maine, and Minnesota were also reviewed. In addition to Washington, interchange safety data from Missouri were also used. The Missouri data were acquired from Missouri DOT's Transportation Management System as part of a recently completed *Highway Safety Manual* (HSM) calibration project.[5,6,7]

The second step of the data generation process establishes the effect of each contributing factor on crash frequency. This information was also synthesized from published literature, HSM, and the Crash Modification Factors Clearinghouse.[8] When no reliable information was available for a particular variable, assumptions were made based on analyzing observed data from Washington and Missouri.

The third step of the data generation process quantifies the combined effect of all contributing factors on crash frequency. This quantification was done in two stages. First, the research team

estimated the composite crash measure for a given site based on its roadway and traffic characteristics. They considered both individual effects of each factor and interaction effects between two or more factors in generating the composite measure. A site with a higher composite crash measure was likely to experience a higher crash frequency. In the second stage, the researchers converted the composite crash measure to realistic crash frequency (i.e., counts) using observed crash data. This conversion was done using a hierarchical Poisson approach, with parameters optimized for each level of the hierarchy using observed crash data. The research team adjusted the generated crash data distribution parameters to match the overall distributional shape and crash counts at individual sites. Once the crash counts were generated, they used the crash severity distributions to subdivide the overall crash counts into fatal, injury, and property damage only crashes. In addition to crash severity, crash-specific factors pertaining to the driver (e.g., distraction, age, gender), vehicle type, and roadway (e.g., road condition at the time of crash) were also generated for each crash.

The researchers developed a model evaluation rubric to evaluate the performance of models developed using RAD with a scoring system of 0–100. Table 1 shows the six criteria and the maximum points assigned to each model. Because the primary goal of RAD is to evaluate the ability of models to accurately estimate the cause–effect relationships, model inference is weighed more than other criteria.

**Table 1. Criteria and scores for model evaluation.**

| Criteria | Points |
|---|---|
| Descriptive analysis of data | 10 |
| Model selection | 10 |
| Training and testing data | 10 |
| Overall prediction accuracy | 20 |
| Model inference | 50 |
| Total score | 100 |

The RAD datasets generated for LT and SCL facilities were used to test crash prediction models. Two teams estimated statistical models, while one team developed a series of machine learning models. Statistical models include various forms of negative binomial regression, whereas the machine learning models ranged from a simple ridge regression model to a complex deep learning model TabNet.[9] The model evaluation rubric was applied to the models developed by the three teams. All teams provided basic descriptive statistics of the RAD datasets. Overall scores (out of 100) ranged between 72 and 91, with the main difference in performance appearing in the model inference criteria.

To facilitate the use of RAD, a web-based software was developed to provide access to RAD datasets. Figure 1 is a screenshot of the RAD website's homepage showing the three types of data that are available for each of the two facility types. Figure 2 shows the workflow of the software. The users submit a RAD data request to the web server. The web server will call the RAD generator to produce a set of RAD datasets. Depending on the type of request, either a pregenerated dataset or a custom dataset will be produced. For custom datasets, an email notification with the download URL for the generated data will be sent to the user-provided email.

Source: FHWA.

SHRP = Strategic Highway Research Program; NDS = Naturalistic Driving Study.

**Figure 1. Screenshot. Landing page with main menu options in the RAD software.**



Source: FHWA.

**Figure 2. Graphic. Software development workflow.**

The second type of RAD datasets developed in this project are the virtual reality (VR) simulation testbeds for crashes and near-crashes occurring at interchanges. The testbeds were developed using a four-step process:

1. Analyze Strategic Highway Research Program (SHRP)2 Naturalistic Driving Study (NDS) videos.
2. Develop crash diagrams.
3. Create three-dimensional (3D) modeling of roadway and environment.
4. Reconstruct crash in VR.

Step 1 involves obtaining and analyzing videos of safety-critical events occurring at interchanges. This step was accomplished using SHRP2 NDS data. A total of 114 crash and near-crash events involving left-turning vehicles and 310 events occurring on SCLs were evaluated to develop the testbeds. Both video and kinematic data for safety-critical events were used for reconstruction.

The second step of the crash reconstruction process involves crash diagramming. This task entails drawing detailed trajectories of vehicles involved in the crash event. After drawing the trajectories, road signs are generated. Signs similar to those observed in the crash videos were generated because the actual crash locations were withheld due to privacy concerns.

After extracting the vehicle trajectories and basic signage from the NDS videos, the third step involves creating the roadway and the environment using 3D modeling tools. Coded roadway elements include travel lanes, shoulders, medians, barriers, terrain, overpasses, pavement markings, etc. Environment elements include signage, overall lighting, and foliage next to the highway.The fourth and final step in the testbed development process involves creating a crash simulation. Vehicle information and the trajectories extracted in the first step were overlaid on top of the roadway and environment elements created in the second and third steps to reconstruct the crash. A commonly used simulation engine was used to create the testbeds.

The researchers created a graphical user interface to facilitate the use of simulator testbeds created for LT and SCL crashes. Figure 3 shows a screenshot of the homepage of the user interface. A user has three visualization options (as shown in figure 4):

- An aerial view is a recreated animation of a crash.

- The 360-degree view places the user in the driver's seat of the subject vehicle and provides the driver's perspective of the crash.

- The test-drive view is similar to the 360-degree mode, with the exception that the user actively controls the vehicle.

Although aerial and 360-degree views are not interactive, the test drive mode gives control to the user to drive through the scenario and react to the conditions that led to a crash.

Source: FHWA.

**Figure 3. Screenshot. Main landing page of the simulator testbed user interface.**



Source: FHWA.

**Figure 4. Screenshot. User menu showing three visualization options.**

In summary, this EAR Program project developed synthetic datasets for interchange facilities for the first time. A three-step framework was developed to generate RAD. The developed datasets were then used by state-of-the-art statistical and machine learning approaches for modeling crash frequency and to ascertain the cause-and-effect relationships. A web-based software was developed to provide easy access to the RAD datasets. The software provides 196 pregenerated datasets and the option to submit custom data requests. The RAD dataset is provided in a spreadsheet format similar to the safety datasets obtained from State DOTs. The proposed framework is generic and can be used to generate RAD for other facilities, such as work zones, bicyclist/pedestrian facilities, innovative geometric designs, etc.

This project also extended the idea of RAD by generating realistic simulation testbeds using NDS data. The VR RAD testbeds were developed with two intended purposes. First, the VR animations of crashes and near-crashes can be used for driver education. Because the testbeds were developed using actual crashes documented in the NDS, they provide a realistic experience that is more engaging. For example, outreach activities targeted at teen drivers can use the animations to help provide a realistic, immersive experience of a crash and to encourage safe driving practices in such circumstances. Lower hardware costs, a younger workforce, and

5

investments from technology companies in improving VR experience are all significant reasons to believe that the transportation industry will increasingly embrace VR-enabled training and education.

A second purpose served by the VR testbeds is to assist with human factors research. For example, a driving simulator platform can readily use the testbeds to test the performance of safety countermeasures, such as in-vehicle driver information systems, roadside dynamic message signs, collision avoidance systems, etc. This purpose is well aligned with the USDOT's National Roadway Safety Strategy's Safe System Approach, which considers an overlapping set of safety measures: roadway countermeasures, behavioral interventions, enforcement, vehicle safety features, and emergency medical care.

# CHAPTER 1. BACKGROUND

Data are critical to understanding crash causation, which leads to the optimization of safety countermeasures. The Federal Highway Administration (FHWA) has been an active proponent of data-driven methods for safety decisionmaking. Through the Every Day Counts (EDC) program (EDC-3 and EDC-4), FHWA has been encouraging practitioners to apply a data-driven safety analysis (DDSA) approach to safety decisionmaking.[4] The DDSA advocates for a new line of thinking that relies on predicted and expected safety values using statistical methods. Many States now use DDSA approaches (75 percent per FHWA EDC website) to strategically invest in systemic treatments that target specific crash types rather than chasing high-crash locations and addressing them piecemeal. The U.S. Department of Transportation's (USDOT) *Research, Development, and Technology Strategic Plan (2018–2022)* further reinforced the importance of reliable data and effective analytical tools in achieving the strategic safety goal of zero fatalities.[1] Specifically, the safety data initiative of the systemic safety approach "seeks to develop new and integrated data sources, analysis, and visualization techniques to enhance our understanding of crash risk and our ability to mitigate it."[1]

## PROJECT MOTIVATION

As described by Hauer, artificial realistic data (ARD) dataset can be a useful tool for research on highway safety.[3] Specifically, three areas where the use of ARD could be beneficial are: determining a sample size or assessing a new estimation method, comparing methods to analyze alternatives, and evaluating ways to generate multivariate models. Some of the common challenges in developing safety models include difficulties in identifying causal relationships in the data; the use of average values for variables, variable errors, missing variables, and complex variable dependencies, and the use of simple mathematical functions. ARD offers an innovative way to address these challenges by developing models that not only accurately predict crash frequency but also accurately explain the cause-and-effect relationships between crashes and the independent variables.

Council et al. led the first ARD effort sponsored by FHWA to examine the performance of different modeling methods for cross-sectional studies.[1] Using Highway Safety Information System (HSIS) data from Washington, a dataset was created consisting of 2,400 mi of homogenous segments of 0.02 mi each. Crashes were assigned to each of the segments based on certain known causal relationships. The case study examined single-vehicle-lane-departure crashes occurring on rural two-lane roadways. The crash and roadway data were then provided to a modeler not privy to the assumed causal relationships. The modeler was tasked with estimating regression models and deriving the causal relationships. The model results were then checked against the assumed relationships. This effort is one of the few endeavors in transportation to create synthetic data for safety modeling.

---

[1]Council, F., E. Hauer, B. Lan, D. Harwood, and R. Srinivasan. 2017. *Use of "Artificial Realistic Data" (ARD) to Assess the Performance of Cross-Sectional Analysis Methods in Capturing Causal Relationships Between Individual Roadway Attributes and Safety*. Unpublished Report. Washington, DC: Federal Highway Administration.

**PROJECT OVERVIEW**

In this Exploratory Advanced Research (EAR) Program project, this initial effort by Council et al. was extended by developing synthetic datasets for interchange facilities.[2] Specifically, datasets were generated for crashes occurring at ramp terminals and speed change lanes (SCLs) of diamond interchanges. Diamond interchanges are one of the highly prevalent designs in the United States. These realistic artificial data (RAD) datasets were then used by state-of-the-art statistical and machine learning approaches for modeling crash frequency and to ascertain the cause-and-effect relationships. A web-based software was developed to easily access the RAD datasets. The software provides 196 pregenerated datasets and the option to submit custom data requests. The RAD dataset is provided in a spreadsheet format similar to the safety datasets obtained from State DOTs.

RAD datasets can be used to test the performance of different safety modeling approaches. For example, if a modeler estimated different forms of statistical models using a crash dataset from a particular State DOT, the different models can only be compared using overall goodness-of-fit measures (e.g., prediction accuracy, likelihood value). Since the ground truth cause–effect relationships between independent and dependent variables are seldom known, the models cannot be compared by their ability to extract the true cause–effect relationships. RAD datasets, on the other hand, are created using cause–effect relationships that were established using literature reviews, subject matter expert interviews, and observed safety data from Washington and Missouri. Thus, different models can be compared based on overall goodness of fit as well as on model inference, that is, the model estimated cause–effect relationships versus the ground truth (assumed) relationships. If a model can satisfactorily extract these relationships from RAD data, the user can confidently apply that model to real data (e.g., from a State DOT) and generate reliable crash modification factors (CMFs).

This project also expands the idea of RAD by generating realistic simulation testbeds using Strategic Highway Research Program (SHRP)2 Naturalistic Driving Study (NDS) data. The virtual reality (VR) RAD testbeds were developed with two intended purposes. First, the VR animations of crashes and near-crashes can be used for driver education. Since the testbeds were developed using actual crashes documented in the NDS, they provide an immersive, realistic experience that is engaging. For example, outreach activities targeted at teen drivers can use the animations to help provide a realistic, immersive experience of a crash and to encourage safe driving practices in such circumstances. Lower hardware costs, a younger workforce, and investments from technology companies in improving the VR experience are all significant reasons to believe that the transportation industry will increasingly embrace VR-enabled training and education.

The second purpose of developing the VR testbeds of safety-critical events is to assist with human factors research to improve interchange safety. For example, a driving simulator platform can readily use the testbeds to test the performance of safety countermeasures, such as in-vehicle driver information systems, roadside dynamic message signs, collision avoidance systems, etc.

---

[2]Council, F., E. Hauer, B. Lan, D. Harwood, and R. Srinivasan. 2017. *Use of "Artificial Realistic Data" (ARD) to Assess the Performance of Cross-Sectional Analysis Methods in Capturing Causal Relationships Between Individual Roadway Attributes and Safety*. Unpublished Report. Washington, DC: Federal Highway Administration.

This purpose is well aligned with the USDOT's National Roadway Safety Strategy's Safe System Approach, which considers an overlapping set of safety measures: roadway countermeasures, behavioral interventions, enforcement, vehicle safety features, and emergency medical care.

**STUDY SCOPE**

This project focused on two facilities of a diamond interchange. The first facility is the ramp terminal. Figure 5 shows a diamond interchange with two ramp terminals (one on the south side and one on the north side). Each site refers to one ramp terminal. The crash type of interest is the multivehicle crash occurring between vehicles turning left onto the entrance ramp of the freeway (shown as a left-pointing, curved arrow in the diagram) and the oncoming through vehicle on the crossroad (shown as a downward-pointing, straight arrow). In the RAD dataset, there is no spatial correlation between consecutively numbered sites.



Original photo: Imagery © 2020 Maxar Technologies, map data © 2020 Google®. Modifications by FHWA (see acknowledgments section).

**Figure 5. Map. Example diamond interchange with two ramp terminals.**

The second interchange facility is the freeway SCL. An SCL facility is an uncontrolled terminal between a ramp and a freeway. The schematic in figure 6 shows an entrance SCL measured from the gore point to the end of the taper.[5] Figure 7 shows a real-world example of an entrance SCL segment.



© 2014 American Association of State Highway and Transportation Officials.

**Figure 6. Graphic. Entrance SCLs.[5]**



Original photo: Imagery © 2022 Maxar Technologies, Map data © 2022 Google®. Modifications by FHWA (see acknowledgments section).

**Figure 7. Map. Real-world example of an entrance SCL.**

# CHAPTER 2. LITERATURE REVIEW

The research team reviewed literature from two different domains. First, studies documenting the development and use of synthetic data were reviewed. Second, due to the focus of this project on interchange safety, literature pertaining to the understanding of crash causation at interchanges was examined. While the synthetic data review provided information on available data generation methods, the interchange safety review helped to obtain information on the key independent variables, their impact on crash frequency, and the state-of-the-practice crash prediction models.

## SYNTHETIC DATA

Although the concept of artificial or synthetic data is new in transportation, its use has been demonstrated in other disciplines. A literature review revealed studies have demonstrated the successful development, evaluation, and application of artificial datasets.

### Generation of Synthetic Data

Probabilistic models and deep generative models are two main methods for synthetic data generation. Probabilistic models focus on mimicking the structure and distribution of real data, whereas deep generative models focus on replicating the structure of real data and the information it contains. The Bayesian network and the Markov model are two popular probabilistic models, and deep generative models include variational autoencoder and generative adversarial networks (GAN). In Ping et al., a synthetic data generation tool called DataSynthesizer was proposed using the Bayesian network.[10] Zhang et al. proposed a Bayesian network to generate synthetic high-dimensional data.[11] The Markov chain approach can be applied to generate temporal synthetic data, such as solar states for a smart grid.[12] Islam et al. presented a data augmentation technique to reproduce crash data.[13] GAN have been used to generate synthetic health data and sensor data.[14,15] A method called TGAN was proposed by Xu and Veeramachaneni to synthesize tabular data using GAN.[16]

Ichim provided an overview of several methods of generating synthetic data, such as probability distribution, Latin hypercube sampling, information preserving statistical obfuscation, data shuffling, and multiple imputations.[17] The author proposed the use of a quantile-based bootstrap strategy and tested it using survey data. Bootstrapping has also been used in other studies by Barth et al., Jia and Culver, and Thanathamathee and Lursinsap (2013).[18–20] Other methods that have been used include Hadoop, convolutional neural network (CNN), genetic algorithm, Bayesian Hierarchical model, and n-spheres.[21–25]

### Evaluation of Synthetic Data

The quality of synthetic data is critical to its widespread adoption. One straightforward method to evaluate the quality of synthetic data is to compare the distribution of each variable with the original dataset, but this method does not consider joint distributions of variables. To overcome the limitations, the synthetic and original datasets can be compared by visualizing the joint distribution of high-dimension data with dimension-reduction techniques. The relative performance of two machine learning algorithms on the synthetic dataset and the original dataset

can be used to measure the quality of synthetic data. A good synthetic dataset should preserve the same relative performance as the original dataset.[26]

In this project, a rubric was developed to evaluate the synthetic datasets that were developed. A rubric contains three main components: evaluative criteria; quality definition for those criteria at different levels, and a scoring strategy.[27] The goal is to create a rubric grading system to rank different models based on their performance, which will be helpful for modelers to revise and improve their models. Standard rubrics are created based on expert review and general rules of thumb.[28]

## Applications of Synthetic Data

Patki et al. developed and used a synthetic data vault to create synthetic data for five datasets.[29] A crowdsourced experiment was then performed in which data scientists were asked to create predictive models with both the synthetic data and real data. The results showed that there was no significant difference in the models developed from the synthetic data and real data. Carlucci et al. created a synthetic database of depth images, and experiments to test the database on two publicly available object datasets showed that the features obtained from processing the synthetic data were stronger.[22] In a research study, Soltana et al. developed and tested an approach to generate synthetic test data using synthetic data for citizens' records in a public administration system.[30] The case study demonstrated that the results met the criteria for both logical validity and statistical representation.

von Neumann-Cosel et al. demonstrated the successful application of synthetic datasets in a simulated environment through research in which a lane-tracking algorithm was tested.[31] The use of synthetic images was innovative as lane-tracking algorithms are usually evaluated with real camera data and then verified using ground truth data. The study found that specific output parameters of image processing algorithms could be tested using synthetic images. The implementation of the simulation environment to test the lane-tracking algorithm allowed the process for investigating various scenarios to be automated, thus reducing the required extent of testing on actual roads.

Synthetic datasets have also been implemented successfully in various applications in civil engineering. Jia and Culver investigated several synthetic flow generation methods to develop flow records for hydrological calibration to address the challenges created by the limited availability of historical flow data.[19] The methods were tested using a case study at Buck Mountain Run in Albemarle, VA. The results showed that the best combination of methods was the bootstrapped artificial neural network for low- and medium-flow predictions and a modified drainage area ratio for the highest 10 percent of synthetic flows. In another study, Naess and Claussen used synthetic data to assess the performance of different estimators for the prediction of values for long return periods such as wind loads.[32] Sakshaug and Raghunathan used a Bayesian Hierarchical method to generate synthetic datasets for estimating small areas.[24]

## INTERCHANGE SAFETY MODELING

### State of the Practice Review

Freeway interchanges consist of freeway segments, SCLs, entrance and exit ramp segments, and ramp terminals (intersections). Given the important role freeways play in carrying high traffic volumes, the safety of interchanges is an important concern. American Association of State Highway and Transportation Officials' (AASHTO) Highway Safety Manual (HSM) was updated in 2014 to provide safety performance functions (SPFs) and CMFs for different freeway facilities, including interchanges.[5] However, there are several important limitations related to the use of SPFs and CMFs.[33] The first practical limitation is the amount of data required for calibration. For example, the predictive method for basic freeway segments in the HSM requires a total of 14 data elements, whereas the method for ramps requires 10 elements. Unfortunately, collecting and analyzing data for several of these elements require a high level of effort, as noted in National Cooperative Highway Research Program (NCHRP) Project 17-45 (e.g., length and radii of horizontal curves, length of/offset to median barriers, clear zone width), which can inhibit the effective utilization of these tools.[34] The authors noted challenges in transferability of models to other contexts, which may be reflective of differences in geometric characteristics in the areas where these studies have been conducted. In addition, a recent meta-analysis of SPFs for freeway merge and diverge areas found that existing research in this area has been largely inconsistent.[35] For example, the effect of deceleration length on safety was reported to be significant in some studies and insignificant in others. Collectively, these results reinforce another limitation noted in the HSM in that the SPFs must be calibrated to reflect local driver populations, conditions, and environments.

Another limitation to the use of SPFs and CMFs is specific to the ramp terminal facility type. A review of the literature for on-ramp terminal crashes shows that most studies are constrained by the lack of crash data specific to this type of freeway geometric element. As such, several studies rely on simulations to address the situations of analyzing scenarios that help improve safety and traffic flow through these sections.[36,37]

Determining whether a crash should be located on a ramp terminal is not as straightforward as it seems. The NCHRP 17-45 project, which influenced the production of the HSM chapter on freeway interchanges, includes extensive discussions of the process and criteria for identifying interchange ramp-terminal crashes.[34] In the analysis of this study to generate RAD, the crashes utilized were collected from Missouri and Washington, which have both taken extensive measures to locate and identify such crashes. Washington, in their reporting of freeway crashes for the HSIS database, included an intersection-related variable that makes the process more manageable. Missouri is also able to provide on-ramp terminal crash data at specifically selected locations due to a previously completed research study.[6,7] This project overcame a tremendous data challenge by using Washington and Missouri ramp terminal crash data.

A few studies have attempted to create localized SPFs for interchange ramp terminals by locating and analyzing crashes within specific study areas for specific States.[38–40] An SPF is a calibrated relationship between collision frequency, traffic volume, and other characteristics of a site.[40] Typically, there are several variables employed in these analyses, which are justified by extensive literature in the transportation safety field. These variables include those components

influenced by factors such as traffic volume (exposure), roadway geometry, and traffic signal timing, as described in the literature.

Some of the most common variables utilized in the existing literature (table 2) include annual average daily traffic (AADT) and average daily traffic (ADT) values, clearance intervals as a factor for traffic signal timing, and roadway geometric elements such as roadway surface type, number of lanes, median type and width, and shoulder widths.

**Table 2. Variables used in the safety and operational performance analysis of interchange facilities.**

| Independent variable | Studies that feature the selected variable |
|---|---|
| Segment length | Parajuli et al.;* Le and Porter; Park, Fitzpatrick, and Lord; Claros, Edara, and Sun (2017).[40–43] |
| Speed limit | Wang, Qin, and Noyce;* Chen et al. (2011a); Bonneson and Zimmerman; Fang, Elefteriadou, and Elias; Claros, Edara, and Sun (2016).[38,44–47] |
| Number of lanes | Elefteriadou et al.; Parajuli et al.;* Le and Porter; Claros, Edara, and Sun (2017); Chen at al. (2011a); Claros, Edara, and Sun (2016); Wang et al.; Chen et al. (2011b)[37,40,41,43,44,47–49] |
| ADT | Elefteriadou et al.; Torbic et al.; Le and Porter; Park, Fitzpatrick, and Lord; Chen et al. (2011b), Liu et al.[37,39,41,42,49,50] |
| AADT | Wang, Qin, and Noyce;* Parajuli et al.;* Claros, Edara, and Sun (2017); Chen et al. (2011a); Claros, Edara, and Sun (2016); Wang et al.[38,40,43,44,47,48] |
| Surface type | Chen et al. (2011b); Liu et al.[49,50] |
| Median width | Park, Fitzpatrick, and Lord; Claros, Edara, and Sun (2017); Wang et al.[42,43,48] |
| Lane width | Fang, Elefteriadou, and Elias[46] |
| Shoulder width | Park, Fitzpatrick, and Lord[42] |
| Terminal spacing | Wang, Qin, and Noyce;* Claros, Edara, and Sun (2016)[38,47] |
| Signal timing | Elefteriadou et al.; Wang, Qin, and Noyce;* Bonneson and Zimmerman; Fang, Elefteriadou, and Elias[37,38,45,46] |
| Traffic control type | Torbic et al.;* Claros, Edara, and Sun (2016 and 2017)[39,43,47] |
| Interchange configuration | Torbic et al.;* Parajuli et al.;* Fang , Elefteriadou, and Elias; Claros, Edara, and Sun (2016).[39,40,46,47] |

*Studies used in the generation of SPFs for interchange ramp terminals.

In an extensive study on the safety performance of ramp terminals, Parajuli et al. collected data from 380 ramp terminals in Ontario, QC, Canada.[40] In that study, six different

ramp-terminal-specific SPFs were developed, accounting for the factors of geometry type, traffic control type, and crash severity level.

Wang, Qin, and Noyce presented interesting insights into the effects of yellow or all-red interval timing on crash safety at interchange ramp terminals, as well as terminal spacing and exclusive right turn phases.[38] Their study suggested that the crash frequency will increase with the deficient yellow or all-red intervals and with the increase in terminal spacing.

**Crash Modeling Methods and CMFs**

Many statistical, machine learning, and deep learning methods have been applied in the study of crash modeling. For statistical methods, generalized linear model, Poisson, and zero inflated negative binomial are often employed to model crash frequency.[51] The Poisson distribution has the advantage of simulating the unobserved heterogeneity for a smaller dataset, whereas the negative binomial distribution can simulate a dataset with a lot of zeros and a long tail.[52,53] The Poisson-Gamma model can reflect the skewness of the data and be more tunable with a gamma prior.[54] In this EARP project, a modified version of the Poisson-Gamma model is applied to generate the synthetic data.

Iranitalab and Khattak applied various machine learning methods for predicting crash severity, including k-nearest neighbor (KNN), support vector machine (SVM), decision tree, and random forest.[55] For deep learning methods, Huang et al. used CNN for highway crash detection and risk estimation.[56] A long short-term memory–CNN based model was proposed to predict real-time crash risk on arterials.[57]

The CMF Clearinghouse is an online repository of CMFs for various types of facilities. The research team queried CMFs for interchange-related facilities. The CMF Clearinghouse also cites the studies from which a given CMF is derived. Table 3 provides a summary of these studies, including the independent variables used in the model and the data used for model development. Table 3 lists seven studies. The State data from these studies include Florida, California, Washington, and Texas; Florida data were used by five of the seven studies.

**Table 3. Interchange-related studies listed in CMF Clearinghouse.**

| Study Title | Authors | Data | General Variables (Number and Types) |
|---|---|---|---|
| Safety Evaluation of Geometric Design Criteria for Spacing of Entrance–Exit Ramp Sequence and Use of Auxiliary Lanes | Le and Porter[41] | Digital mapping and satellite imaging applications, primarily Google® Earth™ and Google Maps™; the online interchange database available through the Washington State DOT Interchange Viewer; HSIS database. | 21 variables: AADT, segment length, high occupancy vehicle indicator, crash count (by types), California interstate indicator, Washington interstate indicator. |
| Evaluating the Effects of Freeway Design Elements on Safety | Park, Fitzpatrick, and Lord[42] | Texas DOT geometric database and Texas crashes electronic database. Crash data for 5 yr (1997–2001). | 10 variables: number of lanes, ADT, segment length, shoulder width (inside and outside), lane width, median width, barrier presence indicator, curvature, on-ramp density, crash counts (no PDO). |
| Selecting Optimal Deceleration Lane Lengths at Freeway Diverge Areas Combining Safety and Operational Effects | Chen, Zhou, and Lin[58] | Three-year crash data (2004–2006) obtained from the crash database maintained by Florida DOT. Freeway segments selected from the Florida Highway System. | 6 variables: deceleration lane length, total crash count, average crash frequency, crash index, average percentage of severe crashes, delay. |
| How Lane Arrangements on Freeway Mainlines and Ramps Affect Safety of Freeways with Closely Spaced Entrance and Exit Ramps | Liu et al.[50] | Three-year crash data (2004–2006) obtained from the crash database maintained by Florida DOT. Only crashes that occurred on the deceleration lanes were included in the study. | 10 variables: number of lanes, ADT, posted speed limit (indicator), right shoulder width, arrangements (indicator), road surface condition, land type, road surface type, right shoulder type. |
| Safety Evaluation of Truck-Related Crashes at Freeway Diverge Areas | Wang et al.[48] | Truck-related crash data collected at selected freeway exit ramp segments in the State of Florida. Geometric data and traffic data were collected from the Florida Roadway Characteristics Inventory database. | 12 variables: injury severity (KABCO)[59], deceleration length, number of lanes, shoulder width, median width, speed limit, speed limit difference between highway and ramp, curve (indicator), grade (indicator), AADT, truck AADT, ramp AADT. |
| Identifying Crash Distributions and Prone Locations by Lane Groups at Freeway Diverging Areas | Chen et al.[44] | Interstate freeways in Florida; 3 yr (2004–2006) of crash data for 326 sites. Total of 7,872 crashes were reported with average value of 4.78, 12.82, 10.23, and 15.41 crashes per year. | 10 variables: crash types, number of freeway lanes, lane unbalanced exit ramp (indicator), number of lanes on exit ramp, deceleration length, freeway ADT, exit ramp ADT, right shoulder type. |
| Operational and Safety Performance of Left-Side Off-Ramps at Freeway Diverge Areas | Zhou et al.[60] | Crash data collected at 74 freeway segments in Florida, with 11 sites for left-side off-ramps and 63 sites for the right-side off-ramps. | 8 variables: left-side off-ramp (indicator), log AADT on Freeway, log AADT on ramp, crash counts, number of lanes, posted speed limit, ramp length, deceleration lane length. |

PDO = property damage only.

16

# CHAPTER 3. METHODOLOGY

## SELECTION OF INDEPENDENT VARIABLES

The first step of the data generation process identifies a set of factors expected to contribute to crash occurrence at the selected interchange facility (e.g., SCL, ramp terminal). This step was guided by the literature reviewed in chapter 2. Roadway, traffic, and driver contributing factors were synthesized from the literature. Table 4 lists the 21 factors identified for left-turn (LT) crashes involving crossroad vehicles turning left onto the entrance ramp to the freeway. The factors include exposure (traffic), signal control, geometric design, environment (e.g., lighting), and driver characteristics. Table 5 shows the 19 factors for SCL crashes. The SCL factors were similar to LT factors, except for signal control.

**Table 4. Contributing factors for LT crashes at entrance ramp.**

| Variable | Description |
|---|---|
| aadt | AADT of the crossroad facility |
| left_turn_aadt | AADT of LT movement onto the entrance ramp |
| presence_of_left_turn_lane | Presence of LT lane on the crossroad |
| number_of_left_turn_lanes | Number of LT lanes on the crossroad |
| signal_control_type | Type of LT signal control scheme |
| functional_class | Functional classification of the crossroad facility |
| jurisdiction | Jurisdiction where the site is present |
| no_lanes | Total number of lanes on the crossroad (in both directions) |
| terrain | Terrain of the crossroad at the ramp terminal |
| horizontal_alignment | Horizontal alignment of the crossroad at the ramp terminal |
| intersection_angle | Intersection skew angle |
| median_presence | Presence of median on the crossroad approach |
| channelization_presence | Presence of LT channelization |
| speed_limit | Speed limit on the crossroad approach (in miles per hour) |
| road_surface_cond | Road surface condition when the crash occurred |
| light_condition | Time of the day when the crash occurred |
| visibility | Visibility when the crash occurred |
| gender | Gender of the driver in the left-turning vehicle |
| age | Age of the driver in the left-turning vehicle |
| distraction | Whether the left-turning driver was distracted |
| vehicle_type | Type of vehicle making the LT |

**Table 5. Contributing factors for SCL crashes.**

| Variable | Description |
|---|---|
| SCL_len | Length of SCL segment (in miles) |
| Jurisdiction | Jurisdiction where the site is present |
| No_Lane | Total number of freeway lanes at the start of SCL segment (does not include acceleration lanes) |
| Terrain | Terrain of the SCL segment |
| Horizontal_alignment | Horizontal alignment of the SCL segment |
| Median_Width | Width of the median (in feet) |
| Inside_Shoulder_Width | Width of the shoulder on the left-hand side (in feet) |
| Outside_Shoulder_Width | Width of the shoulder on the right-hand side (in feet) |
| Speed_Limit | Speed limit on the SCL segment (in miles per hour) |
| Freeway_AADT | AADT on the freeway (in one direction) |
| Ramp_AADT | AADT on the entrance ramp |
| Ramp_Truck_AADT | Annual average daily truck traffic on the entrance ramp |
| road_surface_cond | Road surface condition when the crash occurred |
| light_condition | Time of the day when the crash occurred |
| visibility | Visibility when the crash occurred |
| gender | Gender of the driver at fault in the crash |
| age | Age of the driver at fault in the crash |
| distraction | Whether the driver was driving while distracted |
| vehicle_type | Type of vehicle |

Sampling distributions were generated for each of the factors reported in table 4 and table 5 using observed data from Washington and Missouri. When factors were correlated to each other (e.g., AADT and number of lanes, jurisdiction, and terrain), joint sampling distributions were generated. The RAD for these factors were then generated by repeatedly sampling the distributions for a given sample size (e.g., 500 sites).

HSIS data for interchange crashes were obtained for a 5-yr period to help generate RAD. Washington was the most complete and recent (2013–2017) for the purposes of RAD, although data from California, Illinois, Maine, and Minnesota were also reviewed. In addition to Washington, interchange safety data from Missouri were also used in developing RAD. The Missouri data were acquired from Missouri DOT Transportation Management System database as part of a recently completed HSM Calibration project.[6,7]

Interchange schematics and the milepost where the crash occurred were analyzed together to determine if a crash occurred in the SCLs. Crash reports do not indicate whether a crash was a SCL crash. Washington State DOT (WSDOT) publishes schematics of all interchanges and locations online.[61] For this project, drawings of 205 diamond interchanges were analyzed for

Washington. The interchange schematics provide further information regarding the crossroad segment and ramp terminals that made it possible to query all the crashes that occurred at the various possible locations.[61] A similar process was used to extract crash data from 75 diamond interchanges in Missouri. Google® Maps Road Application Programming Interface™ was used to extract speed limit values for the interchange sites.

Extracting LT crash data for ramp terminals was straightforward, whereas extracting crashes for SCL facilities involved additional steps. Figure 8 shows the four SCL facilities at a diamond interchange, two related to the exit and two related to the entrance.[7]



© 2016 Missouri DOT. Modifications: FHWA (see acknowledgments section).

**Figure 8. Graphic. Components of SCLs at an interchange.[7]**

The mileposts of the "gore" and "taper" points that demarcate the SCLs as a portion of the freeway segment are provided in the schematic drawings. To correctly identify the location of crashes as related to the SCL by direction and location, the direction of travel (increasing or decreasing milepost) and location of SCL (merging or diverging from the ramp) were also extracted from the schematics.

The next step in the data collection process involved locating the potential crashes that occurred at the SCLs. Geographic information systems (GIS) software[62] and HSIS shapefile crash data were used to map all the crashes by using their geolocations. From the earlier process of locating ramp terminal crashes, the Global Positioning System coordinates of the ramp terminals gathered in the process were useful in defining a spatial boundary for filtering the most likely candidates of the SCL crashes. Using a 1-mi buffer from the ramp terminal geolocations, an initial filtering was done to reduce the number of potential crashes under study. An example of the process for spatially locating crashes is shown in figure 9.

**Figure 9. Screenshot. Example of process for spatially locating crashes.**

After the potential crashes were identified, the crash data were merged with location data, and a two-stage filter defining the location parameter of the crashes with respect to the SCLs was used to filter and select the crashes that occurred within the footprint of the SCLs. This process utilized the milepost, road inventory, and direction of travel variables to do the matching and selection.

The final processes involved aggregating the crashes by location to get the number of crashes that occurred at each SCL site. After this aggregation was performed, the SCL files were matched to the corresponding segment locations in the "roadlog" file in the same manner as was done for the crashes, by defining the spatial location and filtering the potential one-to-one matches. The result of this multistep process was the proper identification of the crashes that occurred on the speed changes lanes as distinct from the mainline and ramp crashes.

## ESTABLISHMENT OF CAUSE–EFFECT RELATIONSHIPS

The second step of the data generation process establishes the effect of each contributing factor on crash frequency. This information was also synthesized from the literature reviewed in chapter 2, HSM, and the CMF Clearinghouse.[5,8] When no reliable information was available for a particular variable, cause–effect relationships were established based on analyzing observed safety data from diamond interchanges in Washington and Missouri.

## GENERATION OF CRASH DATA

The third step of the data generation process quantifies the combined effect of all contributing factors on crash frequency. This quantification was done in two stages. First, a composite crash score was estimated for a given site based on its roadway and traffic characteristics. Both individual effects of each factor and interaction effects between two or more factors were considered in generating the composite score. This score is transient and has no practical significance other than the fact that a site with a higher composite crash score is likely to experience a higher crash frequency. In the second stage, the composite crash score was converted to realistic crash frequency (i.e., counts) using observed crash data. This conversion was done using a hierarchical Poisson approach with parameters optimized for each level of the hierarchy using observed crash data. Figure 10 shows an example of the generated crash data distribution (12-A) and the observed distribution (12-B). The distribution parameters were adjusted to match the overall shape and crash counts at individual sites. Due to the prevalence of zero crash sites, it is challenging to identify the parameter values that help produce an accurate count of nonzero crashes. In the example shown in figure 10, the RAD process performs reasonably well at generating nonzero crashes, including the very low numbers of high-crash count sites.

Source: FHWA.

A. Generated crash data distribution.



Source: FHWA.

B. Observed crash distribution.

**Figure 10. Graphs. Generated crash data distribution and observed distribution.**

Once the crash counts were generated, crash severity distributions were used to subdivide the overall crash counts into fatal, injury, and property damage only (PDO) crashes. The low crash count values per site did not allow for generating realistic data for all KABCO severities.[59] In addition to crash severity, crash-specific factors pertaining to the driver (e.g., distraction, age, gender), vehicle type, and roadway (e.g., road condition at the time of the crash) were also generated for each crash.

# CHAPTER 4. LT CRASHES AT INTERCHANGE RAMP TERMINALS

RAD datasets were generated for LT crashes at interchange ramp terminals. The datasets are in a tabular format and include separate files for crash contributing factors, crash counts, and individual crash characteristics. A screenshot of the RAD folder is shown in figure 11. The "RAD_input_variables" contains data pertaining to geometric and traffic characteristics, whereas the "RAD_crash_data" contains crash data for each interchange site. In the example shown in figure 11, data for a 5-yr period (assumed to be 2013–2017) was generated. Characteristics of crashes occurring in a particular year are provided in the file with the corresponding year in the title ("2013_crash_characteristics" provides crash severity and driver and vehicle information for 2013). Table 6 describes each of the 17 variables and their format in the input variable and crash datasets. Depending on the variable, its value could be numerical or categorical. Crash characteristics are described in table 7. The crash characteristics include driver, vehicle, road, environment, and crash.



Source: FHWA.

**Figure 11. Screenshot. Files in RAD folder.**

**Table 6. Variables in the RAD dataset for LT crashes at ramp terminals.**

| Variable | Description | Values |
|---|---|---|
| site id | Unique identification number for each ramp terminal site | Numerical |
| aadt | AADT of the Crossroad facility | Numerical |
| left_turn_aadt | AADT of LT movement onto the entrance ramp | Numerical |

23

| Variable | Description | Values |
|---|---|---|
| presence_of_left_turn_lane | Presence of LT lane on the crossroad | Binary: Yes, No |
| number_of_left_turn_lanes | Number of LT lanes on the crossroad | Numerical |
| signal_control_type | Type of LT signal control scheme | Categorical:<br>PP—protected and permitted<br>PO—protected only<br>FYA—permitted |
| functional_class | Functional classification of the crossroad facility | Categorical:<br>rural minor arterial (6)<br>urban principal arterial (14)<br>urban minor arterial (16) |
| jurisdiction | Jurisdiction where the site is present | Categorical:<br>urban, rural |
| no_lanes | Total number of lanes on the crossroad (in both directions) | Numerical |
| terrain | Terrain of the crossroad at the ramp terminal | Categorical:<br>level, rolling, mountainous |
| horizontal_alignment | Horizontal alignment of the crossroad at the ramp terminal | Categorical:<br>tangent, curve |
| intersection_angle | Intersection skew angle | Categorical:<br>90 degrees, <90 degrees |
| median_presence | Presence of median on the crossroad approach | Binary:<br>yes, no |
| channelization_presence | Presence of LT channelization | Binary:<br>yes, no |
| speed_limit | Speed limit on the crossroad approach (in miles per hour) | Numerical |
| Crashes per site per year | Total number of crashes occurring at a site in a year | Numerical |
| Crashes per site per year by severity | Number of crashes by severity for each site in a year | Numerical |

**Table 7. Additional variables specific to individual crashes for LT crashes at ramp terminals.**

| Variable | Description | Values |
|---|---|---|
| crash_id | Unique identification for each crash | N/A |
| gender | Gender of the driver in the left-turning vehicle | Binary: female, male |
| age | Age of the driver in the left-turning vehicle | Categorical: young age (<25 yr), middle age (25–65 yr), old age (>65 yr) |
| distraction | Whether the left-turning driver was distracted | Binary: yes, no |
| vehicle_type | Type of vehicle making the LT | Binary: passenger vehicle, truck |
| road_surface_cond | Road surface condition when the crash occurred | Binary: dry, wet |
| light_condition | Time of the day when the crash occurred | Binary: day, night |
| visibility | Visibility when the crash occurred | Binary: clear, poor visibility |
| severity | Crash severity | Categorical: FI (fatal, disabling, and minor injury) PDO |
| date | Date of the crash (mo/dd/year) | N/A |

# CHAPTER 5. SCL CRASHES AT INTERCHANGES

This chapter describes the RAD datasets generated for SCL facilities at diamond interchanges. The generated datasets are in a tabular format and include separate files for crash contributing factors, crash counts, and individual crash characteristics. The "RAD_input_variables" contains data pertaining to geometric and traffic characteristics, whereas the "RAD_crash_data" contains multivehicle crash data for each interchange site. Characteristics of crashes occurring in a particular year are provided in the file with the corresponding year in the title (e.g., "2013_crash_characteristics" provides crash severity, driver, and vehicle information for 2013). Table 8 describes each of the 15 variables and their format in the input variable and crash datasets. Ten different crash characteristics are described in table 9. As can be expected, some of the input variables for SCLs are different from those reported for ramp terminals. The length of the SCL segment, inside and outside shoulder widths, freeway and ramp AADT, and truck volume are examples of variables unique to SCL datasets.

**Table 8. Variables in the RAD dataset for SCL facilities.**

| Variable | Description | Values |
|---|---|---|
| site id | Unique identification number for each SCL segment site | Numerical |
| SCL_len | Length of SCL segment (in miles) | Numerical |
| Jurisdiction | Jurisdiction where the site is present | Categorical: urban, rural |
| No_Lane | Total number of freeway lanes at the start of SCL segment (does not include acceleration lanes) | Numerical |
| Terrain | Terrain of the SCL segment | Categorical: level, rolling |
| Horizontal_alignment | Horizontal alignment of the SCL segment | Categorical: tangent, curve |
| Median_Width | Width of the median (in feet) | Numerical |
| Inside_Shoulder_Width | Width of the shoulder on the left-hand side (in feet) | Numerical |
| Outside_Shoulder_Width | Width of the shoulder on the right-hand side (in feet) | Numerical |
| Speed_Limit | Speed limit on the SCL segment (in miles per hour) | Numerical |
| Freeway_AADT | AADT on the Freeway (in one direction) | Numerical |
| Ramp_AADT | AADT on the entrance ramp | Numerical |

| Variable | Description | Values |
|---|---|---|
| Ramp_Truck_AADT | Annual average daily truck traffic on the entrance ramp | Numerical |
| Crashes per site per year | Total number of crashes occurring at a site in a year | Numerical |
| Crashes per site per year by severity | Number of crashes by severity for each site in a year | Numerical |

**Table 9. Additional variables specific to individual crashes.**

| Variable | Description | Values |
|---|---|---|
| crash_id | Unique identification for each crash | N/A |
| gender | Gender of the driver at fault in the crash | Binary: Female, male |
| age | Age of the driver at fault in the crash | Categorical: young age (<25 yr), middle age (25–65 yr), old age (>65 yr) |
| distraction | Whether the driver was driving while distracted | Binary: yes, no |
| vehicle_type | Type of vehicle | Binary: passenger vehicle, truck |
| road_surface_cond | Road surface condition when the crash occurred | Binary: dry, wet |
| light_condition | Time of the day when the crash occurred | Binary: day, night |
| visibility | Visibility when the crash occurred | Binary: clear, poor visibility |
| severity | Crash severity | Categorical: FI (fatal, disabling, and minor injury) PDO |
| date | Date of the crash (mo/dd/year) | N/A |

# CHAPTER 6. MODEL EVALUATION RUBRIC

This section presents a rubric that can be used to evaluate the performance of models developed using RAD. The score has a range of 0–100. Table 10 shows the six criteria and the maximum points assigned to each model. The six criteria try to capture different complexities in modeling; the composite score balances tradeoffs in the modeling process. The criteria and points were determined based on a review of literature. Model inference is weighed more than other criteria because the primary goal of RAD is to evaluate the ability of models to accurately estimate the cause–effect relationships.

**Table 10. Criteria and scores for model evaluation.**

| Criteria | Points |
|---|---|
| Descriptive analysis of data | 10 |
| Model selection | 10 |
| Training and testing data | 10 |
| Overall prediction accuracy | 20 |
| Model inference | 50 |
| Total score | 100 |

Grading for each criterion follows a rating procedure. The number of ratings can vary across the criteria and will be described next.

## DESCRIPTIVE ANALYSIS

This criterion pertains to the calculation of basic summary statistics (e.g., mean, standard deviation), scatterplots of response variables and predictors, and correlation matrix. Rating 1, 100 percent of the maximum score of 10 points, is assigned for providing the basic statistics. Rating 2, 70 percent of 10 points, is assigned for partial descriptive statistics. Rating 3, 40 percent of 10 points, is assigned when very limited descriptive statistics are provided. The score for this criterion is computed by multiplying the percentage (based on rank) and 10 points (maximum possible score for the criterion).

## MODEL SELECTION

The selection of independent variables for inclusion in the model is an important step in developing a model. Use of a variable selection process and parameter tuning earns Rating 1 (100 percent of total criterion score). Limited justification of the variable selection process earns Rating 2 (60 percent), and no justification earns Rating 3 (0 percent).

## TRAINING AND TESTING

It is a good practice to split the overall dataset into training and testing datasets to accurately assess model performance. A Rating 1 (100 percent of maximum possible score of 10 points) is assigned when a model performance on a testing dataset is provided. Additional goodness-of-fit measures using cross-validation data folding may also be provided. No points were assigned for models when they were not evaluated using a testing dataset.

## OVERALL MODEL PERFORMANCE

Model performance can be reported using various goodness-of-fit measures. Examples include mean squared error (MSE), confusion matrix, and log-likelihood. Reporting a minimum of two measures of performance is recommended. Five ratings are established for this criterion based on MSE and confusion matrix. Rating 1 (100 percent of maximum possible score of 20 points) for MSE <1 or >90 percent accuracy of confusion matrix. Rating 2 (80 percent) for MSE <1.5 or >80 percent accuracy of confusion matrix. Rating 3 (70 percent) for MSE <2 or >70 percent accuracy of confusion matrix. Rating 4 (50 percent) for MSE <2.5 or >60 percent accuracy of confusion matrix. Rating 5 (30 percent) for MSE <3 or >50 percent accuracy of confusion matrix.

## MODEL INFERENCE

Model inference entails a model's ability to explain the cause–effect relationship between an independent variable and outcome (e.g., crash frequency). As table 10 shows, the highest number of points (50) are assigned to this criterion. Count regression models can generate CMFs that explain the cause–effect relationship. Other surrogate measures (e.g., partial dependence values) may be used for machine learning models. Ratings are awarded based on the ratio of estimated CMF to true CMF. The ratings are as follows:

- Rating 1 (100 percent of total score) when the ratio is between 0.9 and 1.1.
- Rating 2 (90 percent of score) when the ratio is either between 0.5 and 0.9 or between 1.1 and 1.5.
- Rating 3 (70 percent) when the ratio is either between 0.1 and 0.5 or between 1.5 and 1.9.
- Rating 4 (50 percent) when the ratio is either smaller than 0.1 or greater than 1.9.
- Rating 5 (30 percent) is awarded when the CMF values are computed but do not satisfy Rating 4.

# CHAPTER 7. MODEL TESTING USING RAD DATASETS

To illustrate the use of RAD, sample RAD datasets generated for LT and SCL facilities were used to test the performance of crash prediction models. Two teams estimated statistical models, while one team developed a series of machine learning models. The teams developing the models were not privy to the RAD generation procedures to allow for an unbiased estimation of the cause–effect relationships. The detailed model parameters and results of each of the teams are provided as a separate document to FHWA. This chapter first summarizes the models developed by each team and then presents the evaluation results for the models using the rubric presented in chapter 6.

## TEAM 1 MODELS

A variety of models were tested, and the resultant models presented in this report represent the best models in terms of various fit criteria. Some unobserved heterogeneity was determined to exist in the dataset. A few parameters in the models were determined to be random as a result with statistically significant standard deviations. The level of unobserved heterogeneity appears to be limited, however, with no more than three parameters being determined to be random. The intercept in all three models appears to be random. This observed outcome indicates that a basic random effect is seen to be prevalent in all three datasets. Any unobserved heterogeneity over and above the random effect is typically attributed to the number of LT lanes and speed limit (for the freeway speed change dataset). This obtained result indicates that LT lanes appear to be influenced by stochastic effects that are not captured explicitly in the datasets. The speed limit effect is similar in that variations in speed limits are not sufficiently captured by speed limit indicators alone; instead, the indicators captured the underlying processes that represent a continuous distribution of unobserved effects that are at play with respect to speed limits. Such processes may be capturing the effect of driving behavior.

The team presented findings on the statistical significance of the estimated coefficients, the marginal effects of the associated parameters, and the convergent goodness-of-fit measures, such as log-likelihood, Akaike information criterion (AIC), and the Bayesian information criterion (BIC).

## TEAM 2 MODELS

A series of count data models were estimated to investigate the relationship between various site-specific factors and the number of head-on/LT crashes that occurred at each ramp terminal in the RAD.

First, a series of simple negative binomial models were estimated where crashes were regressed against individual variables of interest. The initial analysis focused on understanding the relationship between crashes and AADT. In the research literature, this relationship is generally accounted for in one of two ways: AADT is logged and treated as an offset variable, with its parameter being constrained to be equal to one, or AADT is included as a covariate in logarithmic form. When treated as an offset, an explicit assumption is introduced that crashes will increase proportionately with traffic volume. In contrast, empirical data often show that

crashes tend to increase at a decreasing rate, with this effect often characterizing the effects of increasing congestion and lower speeds at higher ranges of volumes.

As noted in the previous paragraph, when treated as an offset, the parameter estimate on log(AADT) is constrained to be equal to one. In effect, this means that a 1-percent increase in traffic volume would be associated with a 1-percent increase in crashes.

The degree to which the effects of other variables change with respect to how AADT is introduced in the model was also of interest. Several design parameters of interest, such as the type of LT signal phasing that is introduced, the number of travel lanes on the crossroad, and the speed limit, among other factors, may be influenced by AADT. To this end, results from the following three negative binomial models are presented:

- A series of models for crashes versus each individual variable of interest.

- A series of models for crashes versus each individual variable of interest, along with log(AADT) as an offset.

- A series of models for crashes versus each individual variable of interest, along with log(AADT) treated as a covariate.

Forecasting accuracy measures were calculated based on the difference between the predicted crashes and crashes provided in the data for the training dataset. The following measures have been included: the mean absolute deviation, the sum squared error, MSE, the root mean square error, and the standard deviation of errors. CMFs were presented for the statistically significant variables.

**TEAM 3 MODELS**

Seven machine learning models were developed to predict crash counts using the RAD dataset, ranging from a simple ridge regression model to a complex deep learning model TabNet.[9] The seven models are: ridge regression, KNN, SVM, decision tree, random forest, XGBoost,[63] and TabNet.[9] The first five models are commonly used machine learning models found in most introductory textbooks. XGBoost is an implementation of gradient-boosted decision trees. XGBoost is well-known for its speed and good performance in Kaggle (machine learning) competitions. It usually works well for structured or tabular data, which is a good match for the RAD dataset. TabNet is a deep learning model specifically designed for tabular data learning. It uses sequential attention to choose which features to reason from at each decision step and achieves state-of-the-art performance for tabular data modeling. TabNet uses masked self-supervised learning to learn representation for categorical columns and then uses them in the prediction to improve the results.

In this study, the crash count prediction was formulated as a regression problem since the range of crash counts is relatively large for a classification task and is not deterministic. For models that need a validation dataset, a portion of the training dataset was used for validation.

The performance of machine learning models on testing dataset (20 percent of the entire dataset) were reported using MSE, relative squared error (RSE), and R-square. Partial dependence plots

were generated to extract the cause–effect relationships. Partial dependence shows the dependence between the target response and the input feature of interest, marginalizing over the values of all other input features. The partial dependence values were used to compute CMFs and 95 percent confidence interval values for CMFs.

## EVALUATION OF MODELS

The model evaluation rubric proposed in chapter 6 was applied to the models developed by the three teams. All teams provided basic descriptive statistics of the RAD datasets. Samples of these statistics can be found in the appendix. Thus, all teams received the maximum score of 10 points on this criterion. Model selection criterion entails justification of the inclusion of independent variables in the crash frequency models. When sufficient details were provided on the variable selection process and the tuning of parameters, a score of 10 was assigned. A score of 8 was assigned to some models when there was missing information on variable selection. The next criterion of training and testing dataset involved not using the entire RAD dataset for model estimation. A test dataset not used for model estimation can provide a good evaluation of any model overfitting. Teams receiving the data were asked to divide the dataset into training and testing datasets for a robust evaluation. Thus, all models received a score of 10. Overall model performance evaluation involves assessing the ability of the models to accurately predict the crash frequency for a site. The goodness-of-fit measures reported varied between statistical and machine learning models. Statistical models reported AIC, BIC, and likelihood values, whereas machine learning models reported MSE and RSE values. Overall model performance scores are shown in table 11. Finally, the ability of a model to explain the cause–effect relationship between an independent variable and the outcome (i.e., crash frequency) was evaluated. Marginal effects and CMFs were provided for some statistical models. Partial dependence plots and pseudo-CMFs were provided for the machine learning models. The scores for this criterion ranged between 30 and 45 for the different models.

**Table 11. Evaluation of the statistical and machine learning models developed using RAD datasets.**

| Criteria | Maximum Score | Team 1 LT | Team 2 LT | Team 3 LT | Team 1 SCL | Team 2 SCL | Team 3 SCL |
|---|---|---|---|---|---|---|---|
| Descriptive statistics | 10 | 10 | 10 | 10 | 10 | 10 | 10 |
| Model selection | 10 | 8 | 8 | 10 | 8 | 8 | 10 |
| Training and testing data | 10 | 10 | 10 | 10 | 10 | 10 | 10 |
| Overall model performance | 20 | 14 | 16 | 16 | 14 | 14 | 16 |
| Model inference | 50 | 30 | 30 | 35 | 35 | 35 | 45 |
| Total score | 100 | 72 | 74 | 81 | 77 | 77 | 91 |

In summary, the statistical models developed by Teams 1 and 2 had similar performances. The machine learning models developed by Team 3 outperformed the statistical models, especially in the model inference criterion. The CMFs computed from Team 3 models were closer to the true CMFs, thus explaining the assumed cause–effect relationships. The performance of machine learning models on the SCL RAD dataset was found to be better than their performance on the ramp terminal dataset. One possible reason for the superior performance of machine learning

models over count regression models is their ability to better capture the nonlinear relationships between crash frequency and the independent variables.

One other use of RAD is to allow researchers to compare the performance of different models estimated using different observed data. Consider, for example, two studies developing crash frequency models for SCLs—Study A is using observed data from State A, and Study B is using observed data from State B. How to accurately assess the CMFs generated from the two studies? How to determine which CMFs explain the true cause–effect relationships between crash frequency and the corresponding input variables (e.g., freeway AADT, SCL length)? RAD can help answer these questions. The research teams conducting the two studies apply their modeling approaches to the same RAD dataset. The rubric offered in chapter 6 can be used for comparing the model performance along with statistical tests such as goodness of fit, marginal effects, variable variance, etc. Alternatively, the comparison may also be made by only using the model inference criterion, that is, comparing the CMFs generated using the RAD dataset to the known CMFs (i.e., those used to develop RAD). The modeling approach that performs the best on the RAD dataset is more likely to explain the true cause–effect relationship in observed data. Thus, in the above example, if the performance of the model developed in Study A on the RAD dataset is better than the performance of the model developed in Study B, Study A is likely to also produce more reliable CMFs when using observed data. This type of testing was not conducted in this project due to the lack of readily available models (and CMFs) for interchange ramp terminals and SCLs already estimated using data from different datasets.

# CHAPTER 8. SOFTWARE DEVELOPMENT

## RAD SOFTWARE DEVELOPMENT

The software development of this project includes the implementation of a RAD generator and the website that can serve real-time user requests. Figure 12 shows the workflow of the software. The users submit a RAD request to the web server. The web server will call the RAD generator to produce a RAD dataset. Depending on the type of request, either a pregenerated dataset or a custom dataset will be generated. For custom datasets, an email notification with the download URL for the generated data will be sent to the user-provided email.



Source: FHWA.

**Figure 12. Graphic. Software development workflow.**

## WEB DEVELOPMENT

For web development, HTML5, CSS, and JavaScript programming languages are used for front-end user interface implementation, and PHP scripting language is used in the backend to provide web services.

The website developed in this project includes the following seven functionalities:

1. User authentication.
2. Custom RAD request handler.
3. Running time estimator.
4. Multijob scheduler.
5. Email notification.

6. Pregenerated RAD downloads.
7. VR animation and simulator testbeds.

**User Authentication**

To ensure only authenticated users have access to the RAD datasets, a password authentication method is used. Anyone who wants access to the website will be redirected to the login page and must enter their credentials to proceed. Figure 13 shows a screenshot of the login page. The landing page after logging into the website is shown in figure 14.



Source: FHWA.

**Figure 13. Screenshot. User authentication page.**



Source: FHWA.

**Figure 14. Screenshot. Landing page with main menu options in the RAD software.**

## Custom RAD Request Handler

After logging into the website, the user will have the option to submit custom RAD requests to the server. The users need to enter three input parameters (see figure 15): the number of sites, the number of years of RAD, and an email address to receive the dataset once generated. Once the RAD request is submitted to the web server, the RAD request handler calls the RAD generator to create the dataset.



Source: FHWA.

**Figure 15. Screenshot. Custom RAD request handler.**

## Running Time Estimator

After the users submit their RAD request, the running time estimator will provide a notification with the estimated time needed to complete the job to the users on the website. If the users think the job will take too long, they can click the cancel button in the popup window to cancel the RAD request and modify the input parameters to resubmit a new request, which takes less time, or the user can click the OK button in the popup window to proceed if they are satisfied with the estimated running time. Figure 16 shows an example of the popup window message for a custom query.



Source: FHWA.

**Figure 16. Screenshot. Running time estimator for custom query.**

## Multijob Scheduler

There could be multiple users submitting requests at the same time, or one user may submit different requests within a short timeframe. These situations could result in heavy workload for the server and increase the waiting time for users. To solve this problem, a multijob scheduler has been implemented on the server side to facilitate the speed of this RAD software. The multijob scheduler uses multithreading and can handle five requests simultaneously. The schedular uses a first come, first serve mechanism to process jobs. In the future, if the number of requests increases to a situation that the current multijob scheduler could not handle, the server capabilities can be enhanced to address the increasing demand.

## Email Notification

The email notification function allows the users to receive an email stating their requests have finished, and they can download the RAD through the URL provided in the email. This function is convenient because it eliminates wait times and even allows the submission of multiple RAD data requests.

## Download Pregenerated RAD

In addition to submitting custom RAD requests, users can also download pregenerated RAD datasets. There are pregenerated RAD datasets for 16 combinations of number of sites and years. For each combination, five different pregenerated datasets are provided. To download a pregenerated RAD, the users select the number of sites and years they want and then click download, and they will randomly get one dataset from the five pregenerated datasets for the chosen combination. Figure 17 shows the pregenerated RAD window where the users select the two input parameters.



Source: FHWA.

**Figure 17. Screenshot. Window to download pregenerated RAD.**

## VR Animation and Simulator Testbeds

This web page includes crash recreation animation software and data files for different scenarios. A screenshot of the web page is shown in figure 18.

| File Name | Description | Download Link |
|---|---|---|
| Simulator exe file | Animated videos of left turn crashes and near-crashes | Final_RAD_Sim.zip |
| Simulator scenario 1 | Unity files for a daytime near-crash event | Final_RAD_Sim_Scenario_1.zip |
| Simulator scenario 2 | Unity files for a daytime crash event | Final_RAD_Sim_Scenario_2.zip |
| Simulator scenario 3 | Unity files for a night-time near-crash event | Final_RAD_Sim_Scenario_3.zip |

Source: FHWA.

**Figure 18. Screenshot. Web page for VR animation and simulator testbeds.**

# CHAPTER 9. SIMULATOR TESTBED DEVELOPMENT

This chapter presents the approach used to develop VR simulation testbeds for crashes and near-crashes occurring at interchanges. These testbeds serve two main purposes. First, they can be used as a crash visualization tool for public education. Second, the testbeds will enable safety researchers to evaluate human factors countermeasures (behavioral, technology, etc.) to improve interchange safety. The testbeds were developed using a four-step process shown in figure 19.



Source: FHWA.

**Figure 19. Graphic. Four-step process to develop testbeds.**

Step 1 involved obtaining and analyzing videos of safety-critical events occurring at interchanges. This step was accomplished using SHRP2 NDS data. The NDS dataset was obtained for all junction-related crash, near-crash, and baseline events. The dataset consisted of 41,479 events. To identify the LT and SCL events, a series of data reductions were applied to filter the data based on precipitating event, event nature, traffic control, and event severity. The data reduction yielded 114 crash and near-crash events involving left-turning vehicles and 310 events occurring on SCLs. Forward-view videos were reviewed for each of the 114 events to determine the relative location of the events within the interchange footprint (i.e., crossroad, SCL, ramp segment, etc.). Two locations were of particular interest in this study—crossroad crashes involving left-turning vehicles and SCL crashes. Figure 20 and figure 21 show schematics of an LT crash occurring on the crossroad. The subject vehicle could be either the vehicle turning left onto the entrance ramp or the oncoming through vehicle.

Source: FHWA.

**Figure 20. Graphic. LT crash event.**



Source: FHWA.

**Figure 21. Graphic. LT near-crash event.**

Of the 310 events occurring on SCLs, 179 occurred on entrance SCLs and 131 on exit SCLs. The NDS time series data provide event data at 0.1-s intervals, including the distance between follower and leader, relative velocity, and headway. These data were used to manually validate

and further filter the events occurring on entrance and exit SCLs. Figure 22, figure 23, and figure 24 show schematics of three types of SCL events that were reconstructed: near-crash at entrance SCL, crash at exit SCL on freeway, and crash at exit SCL on deceleration lane.

**Figure 22. Graphic. Near-crash event within entrance SCL.[5]**

**Figure 23. Graphic. Crash event within exit SCL on freeway lane.[5]**

**Figure 24. Graphic. Crash event within exit SCL on deceleration lane.[5]**

Both video and kinematic data for safety-critical events were used for reconstruction. The event information included time of day, weather, vehicle status, vehicle trajectory, surrounding traffic, road conditions, signage, and pavement markings.

The second step of the crash reconstruction process involves crash diagramming. This task entails drawing detailed trajectories of vehicles involved in the crash event. Figure 25 shows an example of trajectories of two vehicles, V1 and V2.

**Figure 25. Graphic. Drawing vehicle trajectories in a computer-aided design program.**

After drawing the trajectories, the road signs are generated. Due to privacy concerns, the actual locations of the crashes were not available for the NDS dataset. Signs similar to those observed in the crash videos were generated because the actual crash locations are not known. Figure 26 shows an example of this step.

A. Plan view of vehicle trajectories and overhead sign.

B. Screenshots from NDS videos.

**Figure 26. Graphic and Photographs. Extracting roadway signs from NDS videos.**

After extracting the vehicle trajectories and basic signage from the NDS videos, the third step involves creating the roadway and the environment using three-dimensional (3D) modeling tools. Coded roadway elements include travel lanes, shoulders, medians, barriers, terrain, overpasses, pavement markings, etc. Figure 27 provides an example of a highway and an overpass structure created during this step. Environment elements include signage, overall lighting, and foliage next to the highway.

**Figure 27. Graphic. Example of a highway and an overpass structure.**

The fourth and final step in the testbed development process involves creating a crash simulation. Vehicle information and the trajectories extracted in the first step are overlaid on top

of the roadway and environment elements created in the second and third steps to reconstruct the crash. A commonly used simulation engine was used to create the testbeds. Figure 28 shows a screenshot of simulating a car in the software. Scripts were written to add background traffic to the simulation.



Source: FHWA.

**Figure 28. Screenshot. Simulating vehicles in a simulation-optimized runtime build.**

The researchers created a graphical user interface to facilitate the use of simulator testbeds created for LT and SCL crashes. Screenshots of the user interface are shown in figure 29, figure 30, figure 31, and figure 32. A user has three visualization options, as shown in figure 30. An aerial view is a recreated animation of a crash. The 360-degree view places the user in the driver's seat of the subject vehicle and provides the driver's perspective of the crash. The test-drive view is similar to the 360-degree mode, with the exception that the user actively controls the vehicle. Although the aerial and 360-degree views are not interactive, the test-drive mode gives control to the user to drive through the scenario and react to the conditions.



Source: FHWA.

**Figure 29. Screenshot. Main landing page of the simulator testbed user interface.**

46

**Figure 30. Screenshot. User menu showing three visualization options.**

**Figure 31. Screenshot. Aerial view of an SCL crash.**

**Figure 32. Screenshot. SCL crash shown in 360-degree view.**

# CHAPTER 10. CONCLUSIONS AND CONSIDERATIONS FOR FUTURE RESEARCH

This project generated two types of RAD datasets for safety research. Tabular data similar to what modelers typically use to conduct safety evaluation studies are the first type of RAD. Users have the option to either download pregenerated datasets or to run custom queries for any number of sites and years. The chapter 6 proposed rubric enabled the chapter 7 development of statistical and machine learning models by using these tabular RAD datasets and making a comparison of their relative performance. The strength of different modeling methods in extracting the cause–effect relationships between the dependent and independent variables can be tested by comparing their performance with assumed cause–effect relationships. Although the rubric standardizes the evaluation process and accommodates the differences in performance measures (e.g., goodness of fit) across different modeling methods, it might still be challenging to compare methods by using different measures of effectiveness. To encourage the application of new statistical and machine learning approaches or variations of existing approaches, the RAD datasets can be introduced in graduate courses at universities. Since the datasets are ready to use (i.e., have gone through a quality assurance and a quality control process), researchers do not have to expend effort in obtaining the data, processing, and preparing them for model development. A student competition organized by the Transportation Research Board or other professional transportation societies could further encourage the use of RAD to discover new modeling approaches.

The second type of RAD datasets generated in the study were the VR testbeds. These testbeds, developed using actual crash videos from NDS, offer a realistic and engaging way to support driver education and countermeasure evaluation studies. Improving interchange safety involves implementing effective behavioral and engineering countermeasures. The testbeds provide human factors researchers with a starting point (i.e., a fully developed ramp terminal or an SCL section with signage and traffic). The modeler can easily add interventions to the testbed and conduct human factors evaluations. The developed testbeds are hardware agnostic and can be used across different visualization options—head-mounted devices, driving simulators, 3D projection systems (e.g., CAVE (cave automatic virtual environment)). The use of VR testbeds can be encouraged through a VR hackathon where users are asked to implement novel safety countermeasures, in the simulation, and evaluate their effectiveness.

This project demonstrated the proposed RAD framework for creating new RAD datasets for interchange facilities, although the same framework can be applied to generate artificial data for other roadway facilities. Of particular interest would be those facilities for which it is difficult to obtain accurate and complete real crash data. Examples of such facilities include work zones, alternative intersections (e.g., diverging diamond, J-turns), and bicycle facilities. The RAD software developed in this project can easily be extended to include data for additional facilities. Another direction for future research entails the application of VR testbeds to evaluate some behavioral and roadway countermeasures. For example, a driving simulator experiment can be set up using a testbed developed for SCLs, and the effect of different driver alert systems (e.g., in-vehicle, dynamic message signs) can be evaluated using study participants.

# APPENDIX. DESCRIPTIVE STATISTICS OF SAMPLE RAD DATA

This appendix presents tables of descriptive statistics of crash frequency and various potential contributing factors for a sample RAD dataset of 400 sites and 5 yr. Table 12, table 13, table 14, table 15, and table 16 provide data pertaining to ramp terminal sites, whereas table 17, table 18, table 19, table 20, and table 21 provide data for SCL facilities.

**Table 12. Descriptive statistics of crash frequency (RAD with 400 LT crash sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 0.49 | 0.64 | 0.56 | 0.52 | 0.55 |
| Median | 0 | 0 | 0 | 0 | 0 |
| Standard deviation | 1.04 | 1.40 | 1.20 | 0.97 | 1.31 |
| Minimum | 0 | 0 | 0 | 0 | 0 |
| Maximum | 9 | 11 | 13 | 6 | 12 |

**Table 13. Descriptive statistics of crossroad AADT (RAD with 400 LT crash sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 15,776 | 16,091 | 16,413 | 16,742 | 17,076 |
| Median | 14,074 | 14,354 | 14,642 | 14,934 | 15,233 |
| Standard deviation | 2,959 | 3,018 | 3,078 | 3,140 | 3,203 |
| Minimum | 2,000 | 2,040 | 2,080 | 2,122 | 2,164 |
| Maximum | 40,000 | 40,800 | 41,616 | 42,448 | 43,297 |

**Table 14. Descriptive statistics of left-turn AADT (RAD with 400 LT crash sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 1,376 | 1,597.6 | 1,629.6 | 1,662.2 | 1,695.4 |
| Median | 1,566.8 | 1,403 | 1,431 | 1,459.5 | 1,489 |
| Standard deviation | 297 | 315 | 309 | 315 | 322 |
| Minimum | 160 | 163 | 166 | 169 | 173 |
| Maximum | 4,302 | 4,388 | 4,475 | 4,565 | 4,656 |

**Table 15. Descriptive statistics of independent variables (RAD with 400 LT crash sites).**

| Variable | Description | Values |
|---|---|---|
| presence_of_left_turn_lane | Presence of LT lane on the crossroad | Yes—234 (58.5 percent) <br> No—166 (41.5 percent) |
| number_of_left_turn_lanes | Number of LT lanes on the crossroad | 0—166 (41.5 percent) <br> 1—199 (49.8 percent) <br> 2—35 (8.7 percent) |
| signal_control_type | Type of LT signal control scheme | PP—132 (33 percent) <br> PO—91 (22.8 percent) |

| Variable | Description | Values |
|---|---|---|
| | | FYA—177 (44.3 percent) |
| functional_class | Functional classification of the crossroad facility | Rural minor arterial—12 (3 percent) Urban principal arterial—211 (53 percent) Urban minor arterial—172 (43 percent) |
| jurisdiction | Jurisdiction where the site is present | Urban—383 (95.7 percent) Rural—17 (4.3 percent) |
| no_lanes | Total number of lanes on the crossroad (in both directions) | 2—126 (31.5 percent) 4—221 (55.3 percent) 6—53 (13.2 percent) |
| terrain | Terrain of the crossroad at the ramp terminal | Level—177 (44.2 percent) Rolling—186 (46.5 percent) Mountain—37 (9.3 percent) |
| horizontal_alignment | Horizontal alignment of the crossroad at the ramp terminal | Tangent—293 (73.2 percent) Curve—107 (26.8 percent) |
| intersection_angle | Intersection skew angle | 90 degrees—313 (78.2 percent) <90 degrees—87 (21.8 percent) |
| median_presence | Presence of median on the crossroad approach | Yes—207 (51.8 percent) No—193 (49.2 percent) |
| channelization_presence | Presence of LT channelization | Yes—180 (45 percent) No—220 (55 percent) |
| speed_limit | Speed limit on the crossroad approach (in miles per hour) | 20 mph—10 (2.5 percent) 25 mph—52 (13 percent) 30 mph—30 (7.5 percent) 35 mph—202 (50.5 percent) 40 mph—38 (9.5 percent) 45 mph—26 (6.5 percent) 50 mph—23 (5.8 percent) 55 mph—16 (4 percent) 60 mph—3 (2.5 percent) |

**Table 16. Descriptive statistics of crash frequency (RAD with 400 SCL sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 0.73 | 0.76 | 0.69 | 0.8 | 0.75 |
| Median | 0 | 0 | 0 | 0 | 0 |
| Standard deviation | 2.93 | 2.73 | 2.85 | 3.92 | 2.67 |
| Minimum | 0 | 0 | 0 | 0 | 0 |
| Maximum | 12 | 13 | 13 | 16 | 13 |

**Table 17. Descriptive statistics of freeway AADT (RAD with 400 SCL sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 43,365 | 44,232 | 45,117 | 46,019 | 46,939 |
| Median | 36,237 | 36,961 | 37,700 | 38,545 | 39,224 |
| Standard deviation | 11,876 | 12,113 | 12,355 | 12,602 | 12,855 |
| Minimum | 2,076 | 2,117 | 2,159 | 2,203 | 2,247 |
| Maximum | 193,551 | 197,422 | 201,370 | 205,397 | 209,505 |

**Table 18. Descriptive statistics of ramp AADT (RAD with 400 SCL sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 5,957 | 6,076 | 6,198 | 6,322 | 6,448 |
| Median | 5,388 | 5,496 | 5,606 | 5,718 | 5,832 |
| Standard deviation | 1,204 | 1,228 | 1,253 | 1,278 | 1,303 |
| Minimum | 37 | 37 | 38 | 39 | 40 |
| Maximum | 18,565 | 18,936 | 19,315 | 19,701 | 20,095 |

**Table 19. Descriptive statistics of ramp truck AADT (RAD with 400 SCL sites).**

| Statistic | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|
| Mean | 1,115.5 | 1,137.3 | 1,160 | 1,183.2 | 1,206.9 |
| Median | 754 | 769 | 784 | 799.5 | 815.5 |
| Standard deviation | 333 | 340 | 346 | 353 | 360 |
| Minimum | 3 | 3 | 3 | 3 | 3 |
| Maximum | 7,001 | 7,141 | 7,283 | 7,429 | 7,578 |

**Table 20. Descriptive statistics of roadway geometric variables (RAD with 400 SCL sites).**

| Statistic | SCL Segment Length (mi) | Median Width (ft) | Inside Shoulder Width (ft) | Outside Shoulder Width (ft) |
|---|---|---|---|---|
| Mean | 0.25 | 48.59 | 7.79 | 10.72 |
| Median | 0.27 | 50 | 7.29 | 10.67 |
| Minimum | 0.05 | 10.96 | 3.02 | 6.04 |
| Maximum | 0.67 | 188.95 | 22.98 | 15.79 |

**Table 21. Descriptive statistics of other independent variables (RAD with 400 SCL sites).**

| Variable | Description | Values |
|---|---|---|
| Jurisdiction | Jurisdiction where the site is present | Urban—216 (54 percent)<br>Rural—184 (46 percent) |
| No_Lane | Total number of freeway lanes at the start of SCL segment (does not include acceleration lanes) | 1—15 (3.8 percent)<br>2—208 (52 percent)<br>3—163 (40.7 percent)<br>4—14 (3.5 percent) |
| Terrain | Terrain of the SCL segment | Level—269 (67.3 percent)<br>Rolling—131 (32.7 percent) |
| Horizontal_alignment | Horizontal alignment of the SCL segment | Tangent—299 (74.7 percent)<br>Curve—101 (25.3 percent) |
| Median Barrier | Presence of a median barrier | Yes—184 (46 percent)<br>No—216 (54 percent) |
| Median_Width | Width of the median (in feet) | 50 ft—10 (2.5 percent)<br>55 ft—19 (4.8 percent)<br>60 ft—163 (40.8 percent)<br>65 ft—21 (5.2 percent)<br>70 ft—187 (46.7 percent) |

# ACKNOWLEDGMENTS

# REFERENCES

1. USDOT. 2020. *Research, Development, and Technology Strategic Plan (2018–2022)*. Washington, DC: U.S. Department of Transportation.

2. FHWA. 2022. "Zero Deaths and Safe System" (website). https://highways.dot.gov/safety/zero-deaths, last accessed October 8, 2022.

3. Hauer, E. 2013. "Artificial Realistic Data: A Research Tool." *Theory, Explanation, and Prediction in Road Safety: Promising Directions*. TR E-Circular Number E-C179. Washington, DC: Transportation Research Board.

4. FHWA. 2022. "About Every Day Counts (EDC)" (website). https://www.fhwa.dot.gov/innovation/everydaycounts/about-edc.cfm, last accessed October 7, 2022.

5. AASHTO. 2014. *Highway Safety Manual.* Washington, DC: American Association of State Highway and Transportation Officials.

6. Sun C., P. Edara, B. Claros, A. Khezerzadeh, H. Brown, and C. Nemmers. 2016. *Highway Safety Manual Applied in Missouri – Freeway/Software*. Report No. cmr 16-009. Jefferson City, MO: Missouri Department of Transportation.

7. Sun C., P. Edara, H. Brown, and C. Nemmers. 2016. *Crash Location Correction for Freeway Interchange Modeling: Final Report*. Report No. cmr 16-010. Jefferson City, MO: Missouri Department of Transportation.

8. FHWA. 2022. "Crash Modification Factors" (website). http://www.cmfclearinghouse.org/, last accessed October 7, 2022.

9. Arik, S., and T. Pfister. 2021. "TabNet: Attentive Interpretable Tabular Learning." *Proceedings of the AAAI Conference on Artificial Intelligence* 35, no. 8.

10. Ping, H., J. Stoyanovich, and B. Howe. 2017. "DataSynthesizer: Privacy-Preserving Synthetic Datasets." *Proceedings of the 29th International Conference on Scientific and Statistical Database Management*, no. 42: 1–5. New York, NY: Association for Computing Machinery.

11. Zhang, J., G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. 2017. "PrivBayes: Private Data Release Via Bayesian Networks." *ACM Transactions on Database Systems* 42, no. 4: 1–41.

12. Tushar, W., S. Huang, C. Yuen, J. A. Zhang, and D. B. Smith. 2014. "Synthetic Generation of Solar States for Smart Grid: A Multiple Segment Markov Chain Approach." *IEEE PES Innovative Smart Grid Technologies, Europe*: 1–6. https://doi.org/10.1109/ISGTEurope.2014.7028832, last accessed October 9, 2022.

13. Islam, Z., M. Abdel-Aty, Q. Cai, and J. Yuan. 2021. "Crash Data Augmentation Using Variational Autoencoder." *Accident Analysis and Prevention* 151: 105950.

14. Yale, A., S. Dash, R. Dutta, I. Guyon, A. Pavao, and K. P. Bennett. 2020. "Generation and Evaluation of Privacy Preserving Synthetic Health Data." *Neurocomputing* 416: 244–255.

15. Alzantot, M., S. Chakraborty, and M. Srivastava. 2017. "SenseGen: A Deep Learning Architecture for Synthetic Sensor Data Generation." *2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*: 188–193. https://doi.org/10.1109/PERCOMW.2017.7917555, last accessed October 9, 2022.

16. Xu, L., and K. Veeramachaneni. 2018. "Synthesizing Tabular Data Using Generative Adversarial Networks." *arXiv*:1811.11264 [cs.LG].

17. Ichim, D. 2010. "Quantile-based Bootstrap Methods to Generate Continuous Synthetic Data." *Proceedings of the 2010 EDBT/ICDT Workshops*: 1–10. New York, NY: Association for Computing Machinery.

18. Barth, R., J. Ijsselmuiden, J. Hemming, and E. J. Van Henten. 2017. "Synthetic Bootstrapping of Convolutional Neural Networks for Semantic Plant Part Segmentation." *Computers and Electronics in Agriculture* 161: 291–304.

19. Jia, Y., and T. B. Culver. 2006. "Bootstrapped Artificial Neural Networks for Synthetic Flow Generation with a Small Data Sample." *Journal of Hydrology* 331, no. 3–4: 580–590.

20. Thanathamathee, P., and C. Lursinsap. 2013. "Handling Imbalanced Data Sets with Synthetic Boundary Data Generation Using Bootstrap Re-Sampling and AdaBoost Techniques." *Pattern Recognition Letters* 34, no. 12: 1339–1347.

21. Anderson, J. W., K. E. Kennedy, L. B. Ngo, A. Luckow, and A. W. Apon. 2014. "Synthetic Data Generation for the Internet of Things." *2014 IEEE International Conference on Big Data (Big Data)*: 171–176. New York, NY: Institute of Electrical and Electronics Engineers.

22. Carlucci, F. M., P. Russo, and B. Caputo. 2017. "A Deep Representation for Depth Images from Synthetic Data." *2017 IEEE International Conference on Robotics and Automation (ICRA)*: 1362–1369. New York, NY: Institute of Electrical and Electronics Engineers.

23. Chen, Y., M. Elliot, and J. Sakshaug. 2016. "A Genetic Algorithm Approach to Synthetic Data Production." *Proceedings of the 1st International Workshop on AI for Privacy and Security*: 1–4. New York, NY: Association for Computing Machinery.

24. Sakshaug, J. W., and T. E. Raghunathan. 2010. "Synthetic Data for Small Area Estimation." *Proceedings of International Conference on Privacy in Statistical Databases*: 162–173. Berlin, Germany: Springer-Verlag.

25. Sánchez-Monedero, J., P. A. Gutiérrez, M. Pérez-Ortiz, and C. Hervás-Martínez. 2013. "An n-Spheres Based Synthetic Data Generator for Supervised Classification." In *Proceedings of International Work-Conference on Artificial Neural Networks*: 613–621. Berlin, Germany: Springer-Verlag.

26. Jordon, J., J. Yoon, and M. van der Schaar. 2018. "Measuring the Quality of Synthetic Data for Use in Competitions." *ArXiv* abs/1806.11345 [cs.LG].

27. Pickett, N., and B. Dodge. 2007. "Rubrics for Web Lessons" (website). http://webquest.org/sdsu/rubrics/weblessons.htm, last accessed April 20, 2022.

28. Goodrich, H. 1996. "Understanding Rubrics." *Educational Leadership* 54, no. 4: 14–18.

29. Patki, N., R. Wedge, and K. Veeramachaneni. 2016. "The Synthetic Data Vault." In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*: 399–410. New York, NY: Institute of Electrical and Electronics Engineers.

30. Soltana, G., M. Sabetzadeh, and L. C. Briand. 2017. "Synthetic Data Generation for Statistical Testing." In *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*: 872–882. New York, NY: Institute of Electrical and Electronics Engineers.

31. von Neumann-Cosel, K., E. Roth, D. Lehmann, J. Speth, and A. Knoll. 2009. "Testing of Image Processing Algorithms on Synthetic Data." In *2009 Fourth International Conference on Software Engineering Advances*: 169–172. New York, NY: Institute of Electrical and Electronics Engineers.

32. Naess, A., and P. H. Clausen. 2001. "Combination of the Peaks-Over-Threshold and Bootstrapping Methods for Extreme Value Prediction." *Structural Safety* 23, no. 4: 315–330.

33. Abatan, A., and P. T. Savolainen. 2018. "Safety Analysis of Interchange Functional Areas." *Transportation Research Record* 2672, no. 30: 120–130.

34. Bonneson, J. A., S. Geedipally, M. P. Pratt, and D. Lord. 2021. *Safety Prediction Methodology and Analysis Tool for Freeways and Interchanges.* NCHRP Web-Only Report No. 306. Washington, DC: Transportation Research Board.

35. Papadimitriou, E., and A. Theofilatos. 2017. "Meta-Analysis of Crash-Risk Factors in Freeway Entrance and Exit Areas." *Journal of Transportation Engineering, Part A: Systems* 143, no. 10: 04017050.

36. Portera, A., and M. Bassani. 2021. "Experimental Investigation into Driver Behavior along Curved and Parallel Diverging Terminals of Exit Interchange Ramps." *Transportation Research Record* 2675, no. 8: 254–267. https://doi.org/10.1177/0361198121997420, last accessed October 9, 2022.

37. Elefteriadou, L., C. Fang, R. Roess, and E. Prassas. 2005. "Methodology for Evaluating the Operational Performance of Interchange Ramp Terminals." *Transportation Research Record* 1920, no. 1: 13–24.

38. Wang H., X. Qin, and D. Noyce. 2010. "Development of a Safety Performance Function for Signalized Diamond Interchange Ramp Terminals." *ICCTP 2010: Integrated Transportation Systems: Green, Intelligent, Reliable*: 1034–1046. Reston, VA: American Society of Civil Engineers.

39. Torbic, D. J., D. W. Harwood, D. K. Gilmore, K. R. Richard, and J. G. Bared. 2007. *Interchange Safety Analysis Tool (ISAT): User Manual*. Report No. FHWA-HRT-07-045. Washington, DC: Federal Highway Administration.

40. Parajuli, B., B. Persaud, C. Lyon, and J. Munro. 2006. *Safety Performance Assessment of Freeway Interchanges, Ramps, and Ramp Terminals*. Ottawa, ON, Canada: Transportation Association of Canada.

41. Le, T. Q., and R. J. Porter. 2012. "Safety Evaluation of Geometric Design Criteria for Spacing of Entrance–Exit Ramp Sequence and Use of Auxiliary Lanes." *Transportation Research Record* 2309, no. 1: 12–20. https://doi.org/10.3141/2309-02, last accessed October 9, 2022.

42. Park, B. J., K. Fitzpatrick, and D. Lord. 2010. "Evaluating the Effects of Freeway Design Elements on Safety." *Transportation Research Record* 2195, no. 1: 58–69. https://doi.org/10.3141/2195-07, last accessed October 9, 2022.

43. Claros, B., P. Edara, and C. Sun. 2017. "When Driving on the Left Side is Safe: Safety of the Diverging Diamond Interchange Ramp Terminals." *Accident Analysis and Prevention* 100: 133–142.

44. Chen, H., H. Zhou, J. Zhao, and P. Hsu. 2011a. "Safety Performance Evaluation of Left-Side Off-Ramps at Freeway Diverge Areas." *Accident Analysis and Prevention* 43, no. 3: 605–612.

45. Bonneson, J. A., and K. H. Zimmerman. 2004. "Effect of Yellow-Interval Timing on the Frequency of Red-Light Violations at Urban Intersections." *Transportation Research Record* 1865, no. 1: 20–27.

46. Fang, F. C., L. Elefteriadou, and A. Elias. 2012. "Field Data for Evaluating 2010 Highway Capacity Manual Operational Analysis Methodology for Interchange Ramp Terminals." *Transportation Research Record* 2286, no. 1: 1–11.

47. Claros, B., P. Edara, and C. Sun. 2016. "Site-specific Safety Analysis of Diverging Diamond Interchange Ramp Terminals." *Transportation Research Record* 2556, no. 1: 20–28.

48. Wang, Z., B. Cao, W. Deng, Z. Zhang, J. J. Lu, and H. Chen. 2011. "Safety Evaluation of Truck-Related Crashes at Freeway Diverge Areas." *Proceedings of 90th TRB Annual Meeting Compendium Papers*: 1–15. Washington, DC: Transportation Research Board.

49. Chen, H., Y. Zhang, Z. Wang, and J. J. Lu. 2011b. "Identifying Crash Distributions and Prone Locations by Lane Groups at Freeway Diverging Areas." *Transportation Research Record* 2237, no. 1: 88–97.

50. Liu, P., H. Chen, J. J. Lu, and B. Cao. 2010. "How Lane Arrangements on Freeway Mainlines and Ramps Affect Safety of Freeways with Closely Spaced Entrance and Exit Ramps." *Journal of Transportation Engineering* 136, no. 7: 614–622.

51. Ozturk, O., K. Ozbay, H. Yang, and B. Bartin. 2013. "Crash Frequency Modeling for Highway Construction Zones." In *TRB 92nd Annual Meeting Compendium of Papers*. Washington, DC: Transportation Research Board.

52. Ye, X., K. Wang, Y. Zou, and D. Lord. 2018. "A Semi-Nonparametric Poisson Regression Model for Analyzing Motor Vehicle Crash Data." *PloS One* 13, no. 5: e0197338.

53. Lord, D., and S. R. Geedipally. 2011. "The Negative Binomial–Lindley Distribution as a Tool for Analyzing Crash Data Characterized by a Large Amount of Zeros." *Accident Analysis & Prevention* 43, no. 5: 1738–1742.

54. Shirazi, M., and D. Lord. 2019. "Characteristics-based Heuristics to Select a Logical Distribution Between the Poisson-Gamma and the Poisson-Lognormal for Crash Data Modelling." *Transportmetrica A: Transport Science* 15, no. 2: 1791–1803.

55. Iranitalab, A., and A. Khattak. 2017. "Comparison of Four Statistical and Machine Learning Methods for Crash Severity Prediction." *Accident Analysis & Prevention* 108: 27–36. https://doi.org/10.1016/j.aap.2017.08.008, last accessed October 9, 2022.

56. Huang, T., S. Wang, and A. Sharma. 2020. "Highway Crash Detection and Risk Estimation Using Deep Learning." *Accident Analysis & Prevention* 135: 105392. https://doi.org/10.1016/j.aap.2019.105392, last accessed October 9, 2022.

57. Li, P., M. Abdel-Aty, and J. Yuan. 2020. "Real-time Crash Risk Prediction on Arterials Based on LSTM-CNN." *Accident Analysis & Prevention* 135: 105371. https://doi.org/10.1016/j.aap.2019.105371, last accessed October 9, 2022.

58. Chen, H., H. Zhou, and P. S. Lin. 2012. "Selecting Optimal Deceleration Lane Lengths at Freeway Diverge Areas Combining Safety and Operational Effects." In *TRB 91st Annual Meeting Compendium of Papers DVD*. Washington, DC: Transportation Research Board.

59. Harmon, T., G. B. Bahar, and F.B. Gross. 2018. *Crash Costs for Highway Safety Analysis*. Report No. FHWA-SA-17-071. Washington, DC: Federal Highway Administration.

60. Zhou, H., H. Chen, J. Zhao, and P. Hsu. 2010. "Operational and Safety Performance of Left-Side Off-Ramps at Freeway Diverge Areas." In *TRB 89th Annual Meeting Compendium of Papers DVD*. Washington, DC: Transportation Research Board.

61. WSDOT. 2016. "Interchange Viewer" (website). https://www.wsdot.wa.gov/mapsdata/tools/InterchangeViewer/default.htm, last accessed April 20, 2022.

62. Esri. 2022. *ArcGIS Pro* (software). Version 3.0.

63. Chen, T., and C. Guestrin. 2016. "Xgboost: A Scalable Tree Boosting System." *Proceedings of the 22nd ACM SIGKDD International Conference On Knowledge Discovery And Data Mining*: 785–794. New York, NY: Association for Computing Machinery.

HRSO-2/01-23(WEB)E