

1 Modeling safety-critical events using trucking naturalistic driving data:
2 A driver-centric hierarchical framework for data analysis

3 Miao Cai^a, Mohammad Ali Alamdar Yazdi^b, Amir Mehdizadeh^c, Qiong Hu^c, Alexander Vinel^c, Karen Davis^d,
4 Fadel Megahed^e, Hong Xian^a, Steven E. Rigdon^{a,*}

5 ^aDepartment of Epidemiology and Biostatistics, Saint Louis University, Saint Louis, MO, 63108, United States

6 ^bCarey Business School, Johns Hopkins University, Baltimore, MD, 21218, United States

7 ^cDepartment of Industrial and Systems Engineering, Auburn University, Auburn, AL, 36849, United States

8 ^dDepartment of Computer Science and Software Engineering, Miami University, Oxford, OH, 45056, United States

9 ^eDepartment of Information Systems and Analytics, Miami University, Oxford, OH, 45056, United States

10 **Abstract**

To be done after the draft is ready.

11 **Keywords:** Trucking, Naturalistic driving studies, Safety-critical events

12 **1. Introduction**

13 The World Health Organization (WHO, 2018) estimated that road injury claimed around 1.4 million lives globally
14 in 2016, which was the eighth leading cause of death. Among all types of vehicles on road, large trucks are a concern
15 since they are more frequently involved in catastrophic crashes. In the United States, National Highway Traffic
16 Safety Administration (2017) reported that 4.3% of registered vehicles were large trucks or buses, but they account
17 for 12.4% of fatalities associated with vehicles (Hickman et al., 2018). Truck drivers are often on the road for long
18 routes under on-time demands, complex traffic and weather conditions, with little to no supervision and contact
19 with fellow workers. Therefore, a number of studies have been published to predict and reduce crash risk associated
20 with trucks (Cantor et al., 2010; Chen et al., 2015; Dong et al., 2017).

21 Traditional crash prediction studies collect retrospective reports of crashes in a given road section for a specified
22 time period, match these crash cases with non-crash controls (typically 1 to 4 matching), and then build statistical
23 models, such as logistic regression and neural networks, to study risk factors associated with higher risk of crashes
24 (Blower et al., 2010; Meuleners et al., 2017; Sharwood et al., 2013). This case-control study design is efficient and
25 less time-consuming in trucking safety field since crashes are very rare. However, case-control studies, by nature, are
26 limited in study design. Firstly, it is impossible to estimate and compare the rate of crashes since the number of
27 non-crashes is unknown. Besides, retrospective reports are often subject to recall and report bias: the drivers may
28 not accurately recall the exact conditions at the time of the event; they may intentionally conceal some critical facts
29 to escape from legal punishment (Dingus et al., 2011; Stern et al., 2019).

30 Naturalistic driving studies (NDSs) have been emerging in the past decade thanks to the advancement of
31 technology. An NDS continuously collects driving data (including latitude, longitude, and speed) under real-world
32 conditions using on-board unobtrusive equipment (Guo, 2019). In contrast to retrospective reports, an NDS resembles

*Corresponding Author

Email addresses: miao.cai@slu.edu (Miao Cai), yazdi@jhu.edu (Mohammad Ali Alamdar Yazdi), azm0127@auburn.edu

(Amir Mehdizadeh), qzh0011@auburn.edu (Qiong Hu), alexander.vinel@auburn.edu (Alexander Vinel), davisk4@miamioh.edu
(Karen Davis), fmegahed@miamioh.edu (Fadel Megahed), hong.xian@slu.edu (Hong Xian), steve.rigdon@slu.edu (Steven E. Rigdon)

33 a cohort study: a pre-determined set of drivers are prospectively followed for a certain amount of time. Therefore,
34 NDS comparatively has several advantages. First, NDS collects both crashes and non-crashes, so it is more useful
35 in comparing the rates of events. Second, since vehicle crashes are extremely rare, it may take a huge amount of
36 driving time to have sufficient sample of crashes. Instead, NDS focus safety-critical events (SCEs), which is defined
37 as events that avoid crashes by last-second evasive maneuver (Dingus et al., 2011). SCEs can be 1000 times as high
38 as real crashes and are argued to be good surrogates of crashes (Dingus et al., 2011; Guo et al., 2010). Third, NDS
39 data are collected using programmed instruments or sensors, therefore they are less likely to be subject to human
40 error or manipulation. Lastly, NDS collects data every a few seconds to minutes, and this large-scale high-resolution
41 data provide a promising opportunity to quantifying driving risk (Guo, 2019).

42 However, many issues arise given the characteristics of NDSs. First, the sheer volume of NDS data creates a
43 challenge to data management and aggregation (Mannering and Bhat, 2014). For example, a NDS data set can
44 have billions rows of real-time speeds and locations, and it is important to have scalable and high-performance tools
45 to aggregate these data into units that fit into the framework of statistical modeling. Second, routinely collected
46 NDS data only have vehicle driving data. Crucial environmental variables such as weather and traffic need to be
47 accessed from other sources and merged back to the driving data. Third, even with these data sources, management,
48 and aggregation issues solved, scalable statistical models that account for the characteristics of NDS are needed to
49 analyze the aggregated data.

50 **A brief review of previous NDS analytic studies.**

51 With increasing vehicle and insurance companies collecting NDS data on a regular basis, a scalable and
52 generalizable analyzing framework serves as a pattern for follow-up researchers to better understand NDS data and
53 gain insights into transportation safety. In this paper, we proposed a framework for data collection, aggregation,
54 fusing, and statistical modeling, which is demonstrated in a case study. Although the NDS data used in this study
55 were from large commercial truck drivers, the framework is generalizable to other drivers since the data collected
56 among different drivers are similar.

57 **2. Data sources**

58 The data were collected by a leading freight shipping trucking company (we will name it as Company A for
59 confidentiality reasons) in the United States. From April 2015 to March 2016, Company A equipped all their trucks
60 with in-vehicle data acquisition systems (DAGs) that collect real-time *ping* and *SCEs* data. Details of these two
61 data sources will be introduced in the following subsection. For demonstration purposes, we sampled 497 regional
62 truck drivers who move freights in a region and surrounding states in this study. Apart from these vehicle driving
63 data, demographic variables including age, gender, and race were also provided to the research team. The drivers
64 were anonymized to ensure confidentiality, while a unique identification number was provided for each driver to link

65 the three data sources. The study protocol was reviewed and approved by the Institutional Review Board of Saint
66 Louis University.

67 *2.1. Ping and SCES data*

68 The DAGs ping irregularly (typically every a couple of seconds to minutes) as the truck goes on road. Each ping
69 collects several key variables, including the date and time (year, month, day, hour, minute, and second), latitude
70 and longitude (specific to five decimal places), driver identification number (ID), and speed at that second. In total,
71 13,187,289 rows of ping data were generated by the 497 truck drivers.

72 Apart from ping data, Company A also collected real-time SCES data for all their trucks. In contrast to irregularly
73 collected ping data, SCES were recorded whenever pre-determined kinematic thresholds were triggered. There were
74 9,032 critical events occurred to these 497 truck drivers during the study period. Four types of critical events were
75 recorded in this critical events data, including 3,944 headway, 3,588 hard brakes, 869 collision mitigation, 631 rolling
76 stability.

77 *2.2. Weather*

78 Apart from driver's characteristics and driving condition, weather also poses a threat on truck crashes and injuries
79 (Naik et al., 2016; Uddin and Huynh, 2017; Zhu and Srinivasan, 2011). We obtained historic weather data from the
80 DarkSky Application Programming Interface (API), which allows us to query historic real-time and hour-by-hour
81 nationwide historic weather conditions according to latitude, longitude, date, and time (The Dark Sky Company,
82 LLC, 2019). The variables included visibility, precipitation probability¹, precipitation intensity, temperature, wind,
83 and others.

84 Traffic and road geometry can be collected from Google map API and OpenStreetmap API. However, querying
85 historic traffic data for all our sample pings from Google map will create costs higher than the budget of the research
86 team. The OpenStreetmap API is open-sourced and free platform that provides road geometry data (including
87 speed limit and the number of lanes), but the missing rate (> 50%) is too high to use for sample pings in this study.
88 Therefore, we did not use traffic data or road geometry data in this study.

¹Ideally, historic precipitation at a specific location and time should be yes or not. However, in reality, since the weather stations are distributed not densely enough to record the exact weather conditions in every latitude and longitude in the US, the DarkSky API uses their algorithms to infer the probability of precipitation in each location.

89 **3. Data preparation**

90 *3.1. Data aggregation*

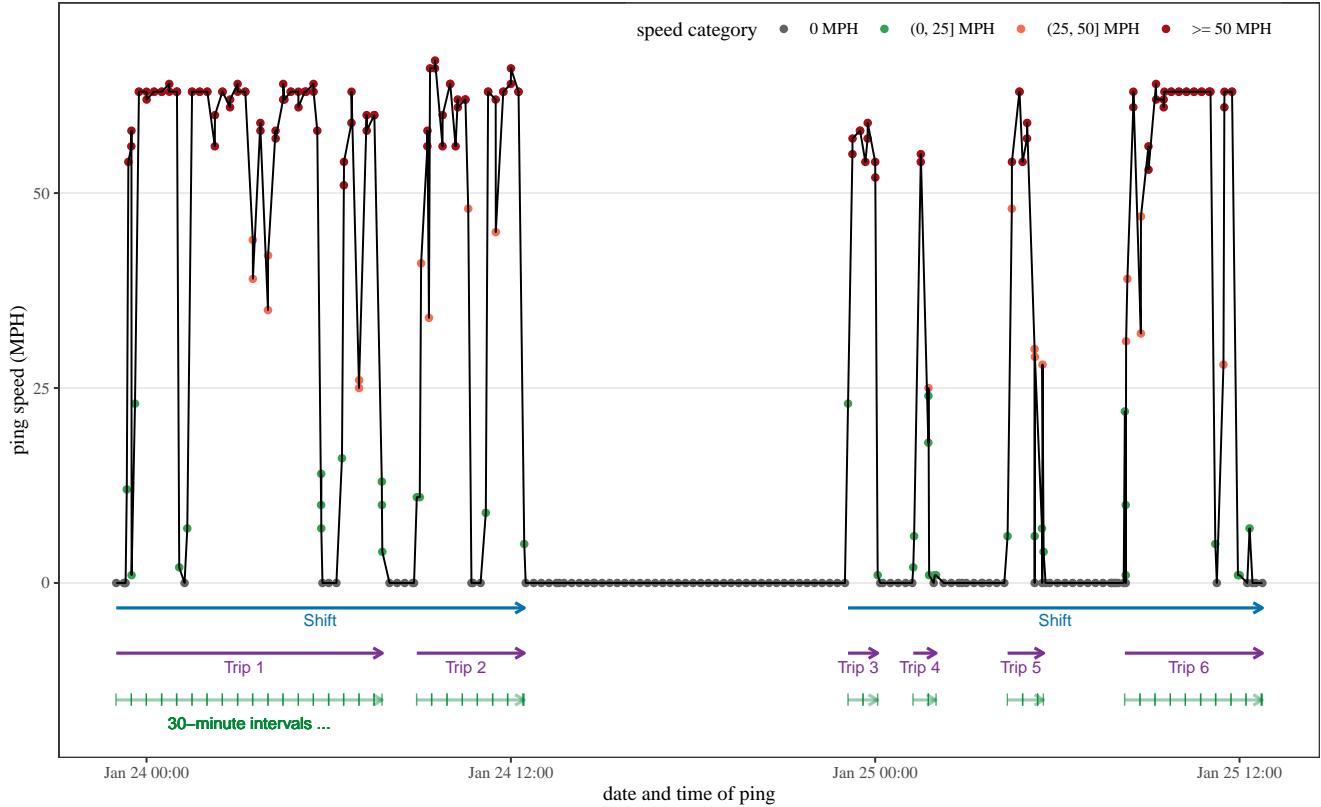


Figure 1: Data aggregation process from pings to shifts, trips, and 30-minute intervals.

91 To shrink the large size of over 10 million ping data, we rounded the GPS coordinates to the second decimal
 92 places, which are worth up to 1.1 kilometers, and we also round the time to the nearest hour. We then queried
 93 weather variables from the DarkSky API using the approximated latitudes, longitudes, date and hour. The weather
 94 variables used in this study include precipitation probability, precipitation intensity, and visibility.

95 For each of the truck drivers, if the ping data showed that the truck was not moving for more than 20 minutes,
 96 the ping data were separated into two different trips. These ping data were then aggregated into different trips. A
 97 **trip** is therefore defined as a continuous period of driving without stop. As Table demonstrates, each row is a trip.
 98 The average length of a trip in this study is 2.31 hours with the standard deviation of 1.8 hours.

99 After the ping data were aggregated into trips, these trips data were then further divided into different shifts
 100 according to an eight-hour rest time for each driver. A **shift** is defined as a long period of driving with potentially
 101 less than 8 hours' stops. The Shift_ID column in shows different shifts, separated by an eight-hour threshold. The
 102 average length of a shift in this study is 8.42 hours with the standard deviation of 2.45 hours.

103 3.2. Cumulative driving time as a measure of fatigue

104 Fatigue has been reported to be the most important predictor to truck crashes, considering that truck drivers are
105 exposed to long routes and lone working environment Stern et al. (2019).

106 Driver's fatigue is difficult to measure in real life. In this study, we attempt to use cumulative driving time in a
107 shift as a proxy measure of the fatigue of truck drivers.

108 4. Methodology

109 4.1. Statistical models

110 Traditional statistical models assume that observations are independent from each other given their predictor
111 variables. However, natural data are almost never independent given the predictor variables. In the example of truck
112 driver's safety events, if we assume the external traffic, weather and driver's socioeconomic status are fixed, truck
113 drivers may exhibit similar driving patterns in multiple trips, and then drivers hired by the same company may
114 share similar culture and safety atmospheres. Therefore, traffic accidents are naturally nested within drivers and
115 drivers are nested within companies. Traditional statistical models that assume independence between observations
116 are not appropriate in this case since objects tend to be similar within a group. Hierarchical models, also known as
117 multilevel model, random-effects model or mixed model, have been developed to allow for the nested nature of data.
118 Instead of assuming independence given predictor variables, hierarchical models assume conditional independence.
119 Hierarchical models are advocated to be the default method since they can produce more precise prediction and
120 more robust results than traditional models.

121 Random-effects models (Han et al., 2018; Pantangi et al., 2019).

Here we model the probability of a critical event occurred using two hierarchical models: logistic and negative binomial (NB) regression models. In the hierarchical logistic regression model, we categorized the number of safety events during the i -th 30-minute interval into a binary variable Y_i with the value of either 0 or 1, where 0 indicated that no critical event occurred during that trip while 1 indicated that at least 1 critical event occurred during the trip. The hierarchical logistic regression model is parameterized as:

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(p_i) \\ \log \frac{p_i}{1 - p_i} &= \beta_{0,d(i)} + \beta_{1,d(i)} \cdot \text{CT}_i + \beta_2 x_2 + \cdots + \beta_k x_k \\ \beta_{0,d(i)} &\sim N(\mu_0, \sigma_0^2) \\ \beta_{1,d(i)} &\sim N(\mu_1, \sigma_1^2). \end{aligned} \tag{1}$$

122 Here $d(i)$ is the driver for interval i , $\beta_{0,d(i)}$ is the random intercept for driver $d(i)$; $\beta_{1,d(i)}$ is the random slope
123 for the cumulative driving time (CT_i) in the shift (the sum of driving time for all previous intervals within that

124 shift) for driver $d(i)$. These random intercepts and random slopes are assumed to have a hyper-distribution with
 125 hyperparameters $\mu_0, \sigma_0, \mu_1, \sigma_1$. x_2, \dots, x_k are other fixed-effect variables including driver demographics (age, gender,
 126 and race), weather (visibility, precipitation intensity and probability), interval specific variables (mean and standard
 127 deviation (s.d.) of speed), and β_2, \dots, β_k are the associated parameters.

Although logistic regression is more robust to outliers of the outcome variable in each 30-interval, it does not fully use the information in the outcome variable since only a binary variable is used. Here we present a hierarchical NB model, with the number of SCEs Y_i^* within the i -th interval as the outcome variable. The hierarchical NB regression model is parameterized as:

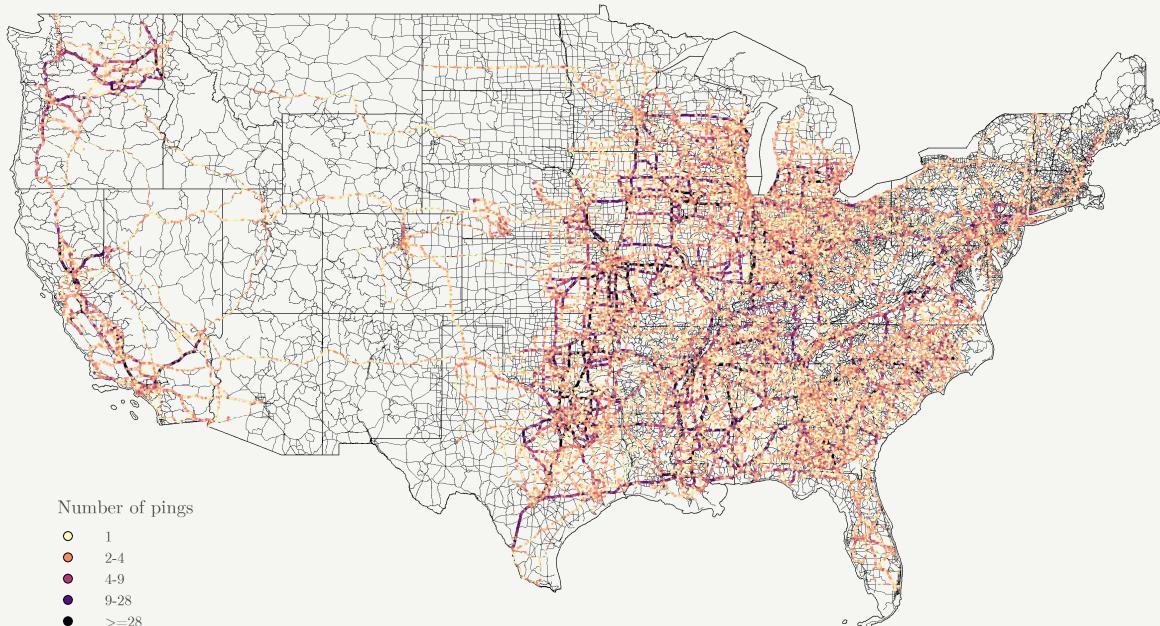
$$\begin{aligned} Y_i^* &\sim \text{NB}(T_i \times \mu_i, \mu_i + \frac{\mu_i^2}{\theta}) \\ \log \mu_i &= \beta_{0,d(i)}^* + \beta_{1,d(i)}^* \cdot \text{CT}_i + \beta_2^* x_2 + \dots + \beta_k^* x_k \\ \beta_{0,d(i)}^* &\sim N(\mu_0^*, \sigma_0^{*2}) \\ \beta_{1,d(i)}^* &\sim N(\mu_1^*, \sigma_1^{*2}). \end{aligned} \tag{2}$$

128 Here T_i is the length of the i -th interval, μ_i is the expected number of SCEs per hour, θ is a fixed over-dispersion
 129 parameter. Other parameters are similar and explained in the previous hierarchical logistic regression model, and we
 130 put a $*$ on the parameter to note the difference between the parameters of the two models.

131 The hierarchical logistic and NB models were estimated using the `lme4` package in R 3.6.2.

Geographical distribution of the moving pings generated by the 496 drivers, 2015-2016

The drivers were employees in large commercial truck company in the United States

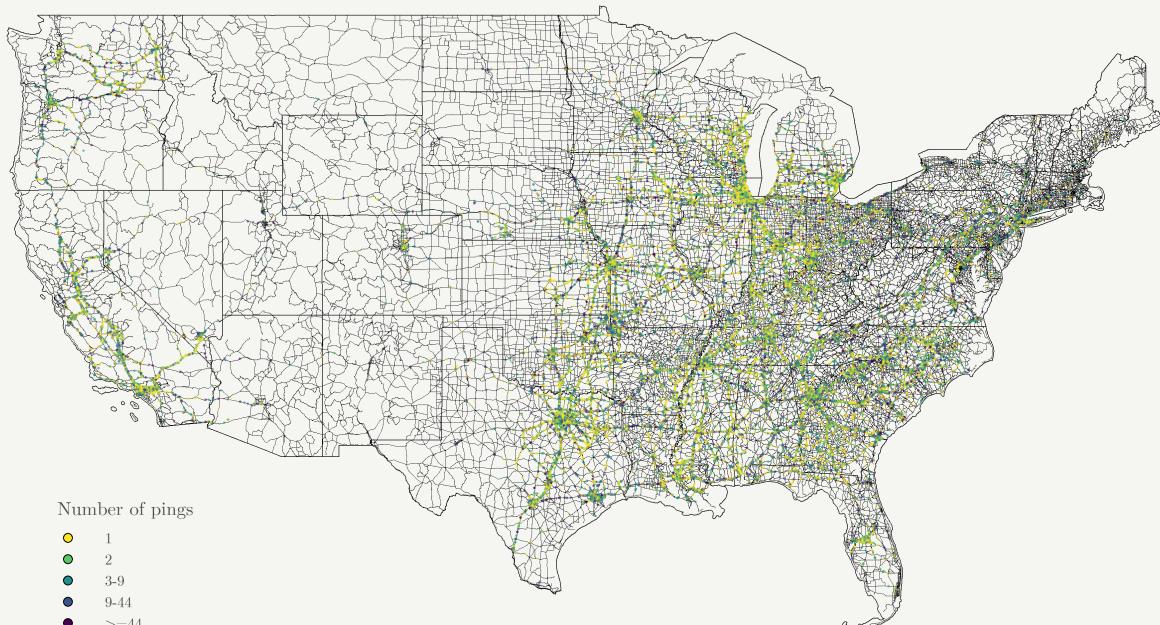


The thin grey line are major highways in the USA. The thicker black lines are state borders.

(a) Active pings

Geographical distribution of the stopped pings generated by the 496 drivers, 2015-2016

The drivers were employees in large commercial truck company in the United States



The thin grey line are major highways in the USA. The thicker black lines are state borders.

(b) Inactive pings

Figure 2: Geographical point patterns of moving and stopped pings generated by the 497 sample drivers.

132 **5. Results**

133 *5.1. Sample description*

134 *5.2. Statistical models*

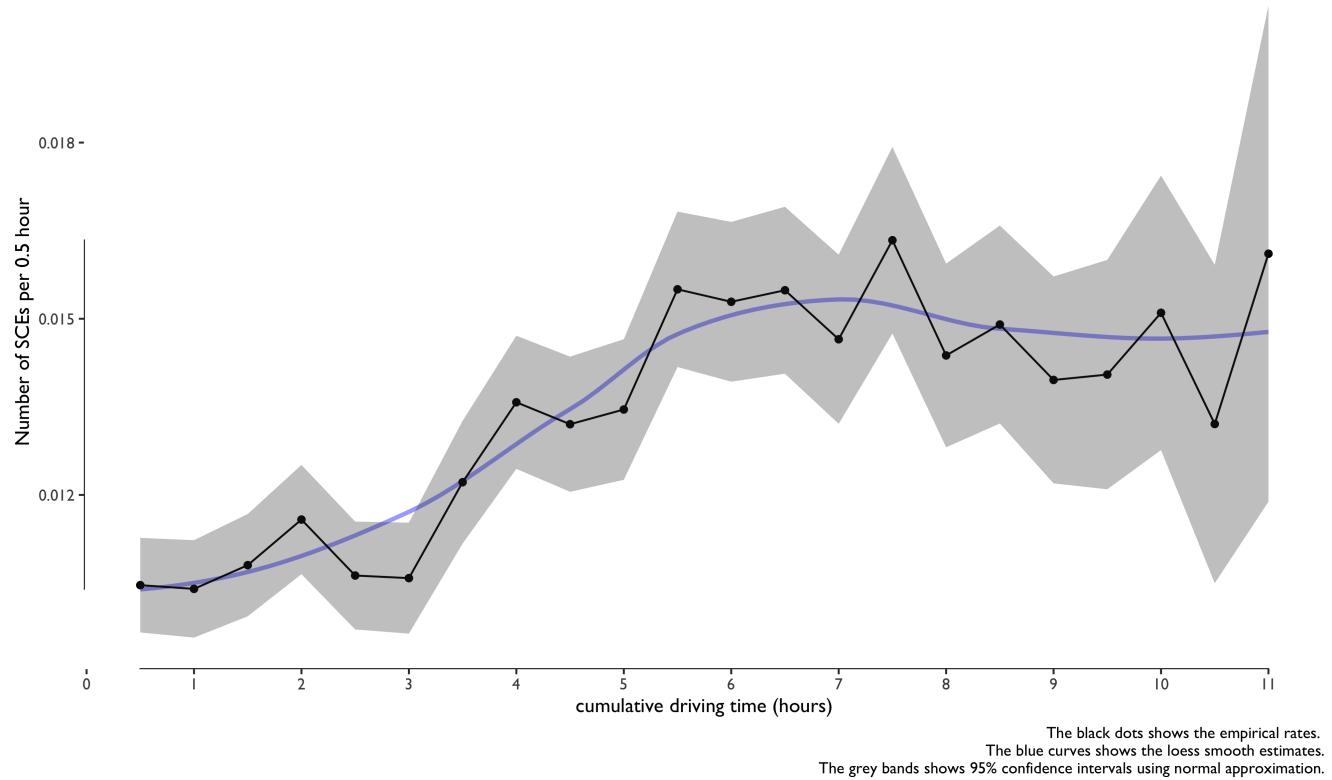


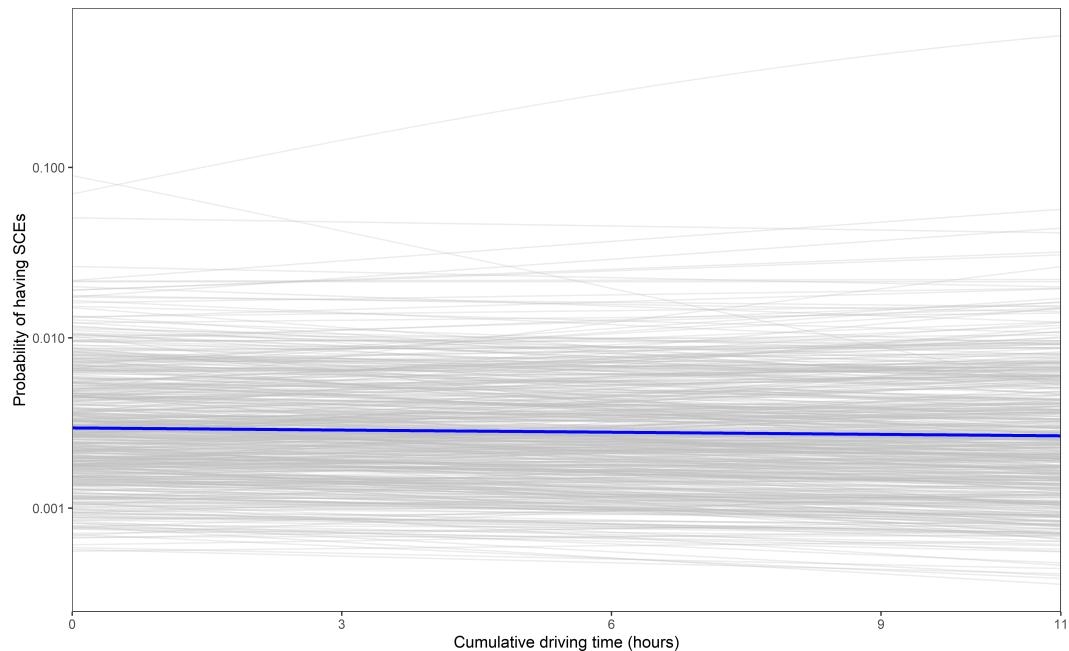
Figure 3: The rate of safety critical events and cumulative driving time

135 **6. Discussion**

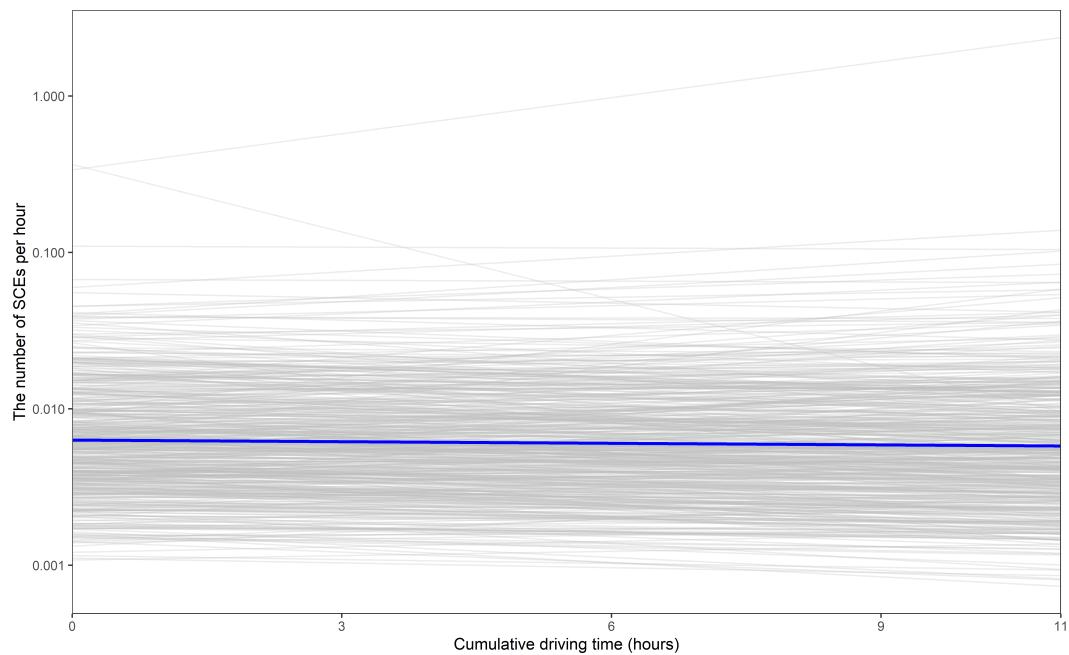
136 **7. Conclusions**

137 **Acknowledgement**

138 This work was supported in part by the National Science Foundation (CMMI-1635927 and CMMI-1634992), the
139 Ohio Supercomputer Center (PMIU0138 and PMIU0162), the American Society of Safety Professionals (ASSP)
140 Foundation, the University of Cincinnati Education and Research Center Pilot Research Project Training Program,
141 and the Transportation Informatics Tier I University Transportation Center (TransInfo). We also thank the DarkSky
142 company for providing us five million free calls to their historic weather API.



(a) Hierarchical Logistics model



(b) Hierarchical negative binomial model

Figure 4: Simulated relationship between cumulative driving time and probability (logistics model)/rate (negative binomial model) of SCEs the 497 sample drivers.

Table 1: Standard and hierarchical logistic and NB models

	logistic (1)	NB (2)	hierarchical logistic (3)	hierarchical NB (4)
Intercept	-4.979*** (0.105)	-7.333*** (0.097)	-5.819*** (0.235)	-8.466*** (0.237)
Cumulative driving	-0.005 (0.004)	-0.004 (0.004)	-0.010 (0.006)	-0.008 (0.007)
Mean speed	-0.0002 (0.001)	-0.0003 (0.001)	0.003*** (0.001)	0.001 (0.001)
Speed s.d.	0.020*** (0.001)	0.017*** (0.001)	0.023*** (0.001)	0.020*** (0.001)
Age	-0.010*** (0.001)	-0.016*** (0.001)	-0.006 (0.004)	-0.007 (0.004)
Race: black	-0.055** (0.025)	-0.124*** (0.026)	0.091 (0.105)	0.093 (0.108)
Race: other	0.238*** (0.042)	0.145*** (0.046)	0.369** (0.179)	0.347* (0.186)
Gender: male	0.288*** (0.050)	0.348*** (0.053)		
Gender: female	0.064 (0.341)	0.061 (0.380)		
Precipitation intensity	0.519 (0.663)	0.418 (0.704)	0.997 (0.670)	0.961 (0.662)
Precipitation probability	-0.175** (0.072)	-0.164** (0.075)	-0.024 (0.074)	0.059 (0.073)
Wind speed	-0.011*** (0.004)	-0.013*** (0.004)	-0.023*** (0.004)	-0.024*** (0.004)
Visibility	-0.029*** (0.005)	-0.043*** (0.005)	0.011** (0.006)	0.010* (0.006)
Interval time	0.015*** (0.002)		0.017*** (0.002)	
Observations	1,019,482	1,019,482	1,019,482	1,019,482
Log Likelihood	-46,303.850	-49,627.630	-43,042.570	-45,961.190
θ		0.036*** (0.001)		
Akaike Inf. Crit.	92,635.690	99,281.260	86,115.150	91,952.390
Bayesian Inf. Crit.			86,292.670	92,129.910

Note:

*p<0.1; **p<0.05; ***p<0.01

143 **References**

- 144 Blower, D., Green, P.E., Matteson, A., 2010. Condition of trucks and truck crash involvement: Evidence from
145 the large truck crash causation study. *Transportation Research Record* 2194, 21–28.
- 146 Cantor, D.E., Corsi, T.M., Grimm, C.M., Özpolat, K., 2010. A driver focused truck crash prediction model.
147 *Transportation Research Part E: Logistics and Transportation Review* 46, 683–692.
- 148 Chen, C., Zhang, G., Tian, Z., Bogus, S.M., Yang, Y., 2015. Hierarchical bayesian random intercept model-based
149 cross-level interaction decomposition for truck driver injury severity investigations. *Accident Analysis & Prevention*
150 85, 186–198.
- 151 Dingus, T.A., Hanowski, R.J., Klauer, S.G., 2011. Estimating crash risk. *Ergonomics in Design* 19, 8–12.
- 152 Dong, C., Dong, Q., Huang, B., Hu, W., Nambisan, S.S., 2017. Estimating factors contributing to frequency and
153 severity of large truck-involved crashes. *Journal of Transportation Engineering, Part A: Systems* 143, 04017032.
- 154 Guo, F., 2019. Statistical methods for naturalistic driving studies. *Annual Review of Statistics and Its Application*
155 6, 309–328.
- 156 Guo, F., Klauer, S.G., Hankey, J.M., Dingus, T.A., 2010. Near crashes as crash surrogate for naturalistic driving
157 studies. *Transportation Research Record* 2147, 66–74.
- 158 Han, C., Huang, H., Lee, J., Wang, J., 2018. Investigating varying effect of road-level factors on crash frequency
159 across regions: A bayesian hierarchical random parameter modeling approach. *Analytic methods in accident research*
160 20, 81–91.
- 161 Hickman, J.S., Hanowski, R.J., Bocanegra, J., 2018. A synthetic approach to compare the large truck crash
162 causation study and naturalistic driving data. *Accident Analysis & Prevention* 112, 11–14.
- 163 Mannering, F.L., Bhat, C.R., 2014. Analytic methods in accident research: Methodological frontier and future
164 directions. *Analytic methods in accident research* 1, 1–22.
- 165 Meuleners, L., Fraser, M.L., Govorko, M.H., Stevenson, M.R., 2017. Determinants of the occupational environment
166 and heavy vehicle crashes in western australia: A case-control study. *Accident Analysis & Prevention* 99, 452–458.
- 167 Naik, B., Tung, L.-W., Zhao, S., Khattak, A.J., 2016. Weather impacts on single-vehicle truck crash injury
168 severity. *Journal of Safety Research* 58, 57–65.
- 169 National Highway Traffic Safety Administration, 2017. A Compilation of Motor Vehicle Crash Data from the
170 Fatality Analysis Reporting System and the General Estimates System.
- 171 Pantangi, S.S., Fountas, G., Sarwar, M.T., Anastasopoulos, P.C., Blatt, A., Majka, K., Pierowicz, J., Mohan,
172 S.B., 2019. A preliminary investigation of the effectiveness of high visibility enforcement programs using naturalistic
173 driving study data: A grouped random parameters approach. *Analytic Methods in Accident Research* 21, 1–12.
- 174 Sharwood, L.N., Elkington, J., Meuleners, L., Ivers, R., Boufous, S., Stevenson, M., 2013. Use of caffeinated
175 substances and risk of crashes in long distance drivers of commercial vehicles: Case-control study. *BMJ* 346, f1140.

- 176 Stern, H.S., Blower, D., Cohen, M.L., Czeisler, C.A., Dinges, D.F., Greenhouse, J.B., Guo, F., Hanowski, R.J.,
177 Hartenbaum, N.P., Krueger, G.P., others, 2019. Data and methods for studying commercial motor vehicle driver
178 fatigue, highway safety and long-term driver health. *Accident Analysis & Prevention* 126, 37–42.
- 179 The Dark Sky Company, LLC, 2019. Dark Sky API — Overview.
- 180 Uddin, M., Huynh, N., 2017. Truck-involved crashes injury severity analysis for different lighting conditions on
181 rural and urban roadways. *Accident Analysis & Prevention* 108, 44–55.
- 182 WHO, 2018. The top 10 causes of death.
- 183 Zhu, X., Srinivasan, S., 2011. A comprehensive analysis of factors influencing the injury severity of large-truck
184 crashes. *Accident Analysis & Prevention* 43, 49–57.