

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 ^a*Department of Epidemiology and Biostatistics, Saint Louis University, Saint Louis, MO, 63103, United States*

6 ^b*Carey Business School, Johns Hopkins University, Baltimore, MD, 21218, United States*

7 ^c*Department of Industrial and Systems Engineering, Auburn University, Auburn, AL, 36849, United States*

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

9 ^e*Department of Information Systems and Analytics, Miami University, Oxford, OH, 45056, United States*

10 **Abstract**

Naturalistic driving studies produce high-resolution, large-scale, and real-world driving data sets, but there is no consistent data aggregation and analysis framework for this type of data. Using routinely collected naturalistic driving data from 497 commercial truck drivers, this study proposes a driver-centric framework for data cleaning, aggregation, and statistical modeling. We aggregated the real-time driving ping data to shifts, trips, and 30-minute intervals according to the driving patterns. Safety-critical events (SCEs), driver demographics, and weather data from a third-party data provider were then merged to the aggregated 30-minute intervals. Driver-centric hierarchical logistic and negative binomial (NB) models with driver-level random intercepts and random slopes for cumulative driving time were proposed to predict safety-critical events.

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

12 **1. Introduction**

13 The World Health Organization [37] estimated that road injury claimed around 1.4 million lives globally in 2016,
14 which was the eighth leading cause of death. Among all types of vehicles on road, large trucks are a concern since
15 they are more frequently involved in catastrophic crashes. In the United States, National Highway Traffic Safety
16 Administration [27] reported that 4.3% of registered vehicles were large trucks or buses, but they account for 12.4%
17 of vehicle-related fatalities [17]. Truck drivers are often on the road for long routes under on-time demands, complex
18 traffic and weather conditions, with little to no supervision and contact with fellow workers. Therefore, trucking
19 safety is an important research topic and a number of studies have been published to predict and reduce crash risk
20 associated with trucks [3, 6, 9].

21 Traditional crash prediction studies collect retrospective police reports of crashes in a given road section for a
22 specified time period, match these crash cases with non-crash controls (typically 1 to 4 matching), and then build
23 statistical models (such as logistic regression and neural networks) to study the risk factors associated with higher
24 risk of crashes and predict real crashes [2, 24, 30]. This case-control study design is efficient and less time-consuming
25 in the field of trucking safety since crashes are very rare. However, case-control studies, by nature, are limited in
26 study design since a) it is impossible to estimate and compare the rate of crashes since the number of non-crashes is
27 unknown, b) retrospective reports are often subject to recall and report bias: the drivers may not accurately recall

*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)

28 the exact conditions at the time of the event, c) the drivers may intentionally conceal some critical facts to escape
29 from legal punishment [8, 33].

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 [12]. In contrast to retrospective reports, an NDS resembles a
33 cohort study: a pre-determined set of drivers are prospectively followed for a certain amount of time. Therefore,
34 NDS has several advantages. First, NDS collects both crashes and non-crashes, so it is more useful in comparing the
35 rates of events. Second, since vehicle crashes are extremely rare, it may take a huge amount of driving time to have
36 sufficient sample of crashes. Instead, NDS focus safety-critical events (SCEs), which is defined as events that avoid
37 crashes by last-second evasive maneuver [8]. SCEs can be 1000 times as high as real crashes and are argued to be
38 good surrogates of crashes [8, 13, 18, 19]. Third, NDS data are collected using programmed instruments or sensors,
39 so they are less likely to be subject to human error, recall bias, or misinformation. Lastly, NDS collects data every a
40 few seconds to minutes, and this large-scale high-resolution data provide a promising opportunity to quantifying
41 driving risk [12].

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 [21]. For example, a NDS data set can have billions rows of real-time
44 speeds and locations, and it is important to have scalable and high-performance tools to aggregate these data into
45 units that fit into the framework of statistical modeling. Second, routinely collected NDS data only have vehicle
46 driving data. Crucial environmental variables such as weather and traffic need to be accessed from other data sources
47 and merged back to the driving data. Third, even with these data sources, management, and aggregation issues
48 solved, there is a lack of consensus on choosing the statistical models that are both sufficiently complex to account
49 for the characteristics of NDS and computationally feasible to fit the large-scale data. With increasing companies
50 collecting NDS data on a regular basis, a scalable and generalizable analyzing framework can serve as a pattern for
51 researchers to better understand NDS data and gain insights into trucking and transportation safety.

52 This paper aims to propose and showcase a generalizable data analytic framework (data collecting, aggregating,
53 fusing, and driver-centric statistical modeling) that accounts for the features of NDS data. To achieve this aim, we
54 addressed the following questions:

- 55 (A) How should we aggregate the high-resolutional NDS data into statistically analyzable units?
56 (B) Where are the third-party data sources available to transportation data analytic studies?
57 (C) What are the risk factors associated with risky driving behavior among the sample truck drivers?

58 The remainder of this paper is organized as follows. Section 2 provides a brief literature review on previously
59 published studies that use NDS data sets. Then, Section 3 presents our NDS data and other third-party data
60 sources. Section 4 demonstrates how we aggregate the ping data into shifts, trips, and 30-minute intervals, and

61 merge different data sources. Section 5 details the driver hierarchical logistic and negative binomial model. Section 6
62 presents the statistical results and interpretation and conclusions and implications are discussed in Section 7.

63 **2. Literature review**

64 Although NDS data only emerge in the recent decade and are relatively new, there are an increasing number
65 of data analytic studies published using this data. In this section, instead of exhaustively reviewing all published
66 papers, we introduced a few recent papers that build statistical models using NDS data sets (either trucks or more
67 general vehicles). The data, methods, and results of these papers are briefly outlined and compared, we then identify
68 and summarize the research gaps.

69 Table 1 presents eight data analytic studies that use NDS data. These studies extract the outcomes and features
70 (such as driving time, sleep patterns, and traffic) potentially associated with driving risk from NDS data sets, then
71 the relationship between the outcome variables and predictors is explored using statistical models. From the listed
72 papers, we could observe the following issues in previously published studies:

- 73 (A) The number of sample drivers are small (around 100 drivers) except for Wali et al. [36]. The studies may not
74 have sufficient statistical power due to the small sample size, and the generalizability may be limited.
- 75 (B) The data sources come from only NDS data sets, which increases the workload and difficulty of data collection.
76 In secondary data analysis, exclusively replying on one data source may limit our power to answer the question.
77 With various organizations collecting data, we can exploit the power of third-party data providers, integrate
78 different sources of data, and as a consequence, improve the prediction accuracy of statistical models.
- 79 (C) Although the listed papers occasionally used hierarchical models, relatively few actually used driver-centric
80 hierarchical models. NDS data sets are naturally generated by a driver-centric process: recruited drivers are
81 followed for a certain amount of time, and all relevant data are collected in this process.
- 82 (D) No consistent framework for cleaning, aggregating, and statistical modeling has been proposed in these papers.
83 NDS data sets collects large-scale high-resolutinal data, which rely on a context-specific, statistically sensible,
84 and computationally affordable to analyze and empower policy making.

85 Our study serves as a complement to the existing literature and a model for future NDS studies. Firstly, we
86 combined routinely collected driving data of 497 commercial truck drivers and a third-party weather data provider.
87 Then, a contextual sensible data aggregation framework is proposed to reduce the original driving data to shifts,
88 trips, and 30-minute intervals. Lastly, we propose to use driver-centric mixed-effect statistical models to analyze the
89 aggregated data.

Table 1: A review of sample size, outcomes, predictors, statistical models, and results in previous NDS data analytic studies

Authors	Year	Sample	Outcomes	Predictors	Statistical model	Results
Soccilich	2013	97 truck drivers	SCEs	driving time	mixed-effect negative binomial model	there is an increase in risk in the 11th driving hour
Chen	2016	96 truck drivers	SCEs	sleep patterns	negative binomial regression	less sleep in the early stage of non-work periods associated with higher risk at least two nighttime periods
Sparrow	2016	106 truck drivers	fatigue	Sleep/wake patterns	mixed-effects ANOVA	Traffic, age and experience, and speed limits are significant factors
Ghasemzadeh	2018	141 general drivers	lane-keeping behavior	driver characteristics, weather and traffic conditions	logistic regression and multivariate adaptive regression splines	Gazis-Herman-Rothery, Gipps, intelligent driver, full velocity difference, and Wiedemann models
Zhu	2018	42 general drivers	car-following	NA	Wiedemann models	intelligent driver model had the best performance
Mollicone	2019	106 truck drivers	hard-braking events	official duty logs, sleep patterns	Poisson regression	frequency of hard-braking events positively associated with predicted fatigue
Pantangi	2019	54 general drivers	speeding, tailgating	driver-, trip-, vehicle-, weather- characteristics	grouped random parameter probit model	the high-visibility enforcement has mixed effects
Wali	2019	3300 general drivers	acceleration, vehicular jerk	traffic and roadway factors	fixed- and random-parameter discrete choice model	intentional volatility is associated with crash and near-crash events

90 **3. Data sources**

91 The data were collected by a leading freight shipping trucking company (we will name it as Company A for
92 confidentiality reasons) in the United States. From April 2015 to March 2016, the company equipped all their trucks
93 with in-vehicle data acquisition systems (DAGs) that collect real-time *ping* and *SCEs* data. Details of these two data
94 sources will be introduced in Subsection 3.1. The study protocol was reviewed and approved by the Institutional
95 Review Board of Saint Louis University.

96 For demonstration purposes, we sampled 497 regional truck drivers who move freights in a region and surrounding
97 states in this study. Apart from these vehicle driving data, demographic variables including age, gender, and race
98 were also provided to the research team. The drivers were anonymized to ensure confidentiality, while a unique
99 identification number was provided for each driver to link the three data sources. The average age of the sample
100 drivers was 45.83 (standard deviation: 12.03), with 36 female drivers (7.2%). There were 247 whites (49.7%), 206
101 blacks (41.4%), and 44 other races (8.9%).

102 *3.1. Ping and SCEs data*

103 The DAGs ping irregularly (typically every a couple of seconds to minutes) as the truck goes on road. Each ping
104 collects several key variables, including the date and time (year, month, day, hour, minute, and second), latitude
105 and longitude (specific to five decimal places), driver identification number (ID), and speed at that second. In total,
106 13,187,289 rows of ping data were generated by the 497 truck drivers, with 8,029,087 (60.89%) of them were active
107 pings (speed of the ping is not zero).

108 Apart from ping data, Company A also collected real-time SCEs data for all their trucks. In contrast to irregularly
109 collected ping data, SCEs were recorded whenever pre-determined kinematic thresholds were triggered. There were
110 9,032 critical events occurred to these 497 truck drivers during the study period. Four types of critical events were
111 recorded in this critical events data, including 3,944 headway (43.67%), 3,588 hard brakes (39.72%), 869 collision
112 mitigation (9.62%), 631 rolling stability (6.99%).

113 *3.2. Weather*

114 Weather is one of the most studied risk factors associated with trucking safety [26, 35, 38]. In this study, we
115 obtained historic weather data from the DarkSky Application Programming Interface (API), which allows us to
116 query historic real-time and hour-by-hour nationwide historic weather conditions according to latitude, longitude,
117 date, and time [34]. The primary weather variables included visibility, precipitation probability¹, precipitation
118 intensity, wind speed, and others. To reduce the cost of querying all 13 million ping data from the DarkSky API,

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.

119 we rounded the GPS coordinates to the second decimal places, which are worth up to 1.1 kilometers, and we also
120 round the time to the nearest hour. Then the weather variables were queried from the DarkSky API using the
121 approximated latitudes, longitudes, date and hour.

122 Traffic and road geometry can be collected from Google map API and OpenStreetMap API. However, querying
123 historic traffic data for all our sample pings from Google map will create costs higher than the budget of the research
124 team. The OpenStreetMap API is open-sourced and free platform that provides road geometry data (including
125 speed limit and the number of lanes), but the missing rate ($> 50\%$) is too high to be of practical use for sample pings
126 in this study. Therefore, we did not use traffic data or road geometry data in this study. We shared our R code to
127 extract weather (the DarkSky API) and road geometry data (the OpenStreetMap) in the Supplementary materials.

128 4. Data preparation

129 4.1. Shifts, trips, and 30-minute intervals

130 To convert this 13 million row real-time ping data into analyzable units, we aggregate them into *shifts*, *trips*, and
131 *30-minute intervals*, which are inspired by real world truck transporting practice and the hours-of-service policy
132 by Federal Motor Carrier Safety Administration [11]. Shifts are on-duty periods with no breaks longer than eight
133 hours (there can be short breaks less than 8 hours). Trips are continuous driving periods with no breaks less than
134 half an hour. These trips are further divided into 30-minute fixed intervals. This is because trips can vary from
135 several minutes to several hours, which are not a good analyzable unit for statistical modeling. The details of the
136 aggregation process is as follows:

- 137 • *Shifts*: for each of the sample truck drivers, if the ping data showed that the truck was not moving for more
138 than eight hours, the ping data were separated into two different shifts on the left and right side of this long
139 break. There could be several short breaks (less than eight hours) within each shift.
- 140 • *Trips*: for each shift, if the ping data showed that the truck was not moving for more than half an hour, the
141 ping data were separated into different trips. These ping data were then aggregated into different trips. The
142 drivers are assumed to be fully driving within each trip since there are not breaks longer than 30 minutes
143 within each trip. The trips are nested within shifts.
- 144 • *30-minute intervals*: each trip is further decomposed into 30-minute fixed intervals according to the start and
145 end time of the trip. The last interval of the trip is typically less than 30 minutes. The 30-minute intervals are
146 nested within trips.

147 Figure 1 visually present the data aggregation process of ping → shifts → trips → 30-minute intervals, as well as
148 the nested structure. The y-axis is speed and x-axis is time. Each dot is a ping, and the color of that ping indicate
149 the current speed. Grey dots indicate stopping pings with the current speed of zero. The arrows in the lower part

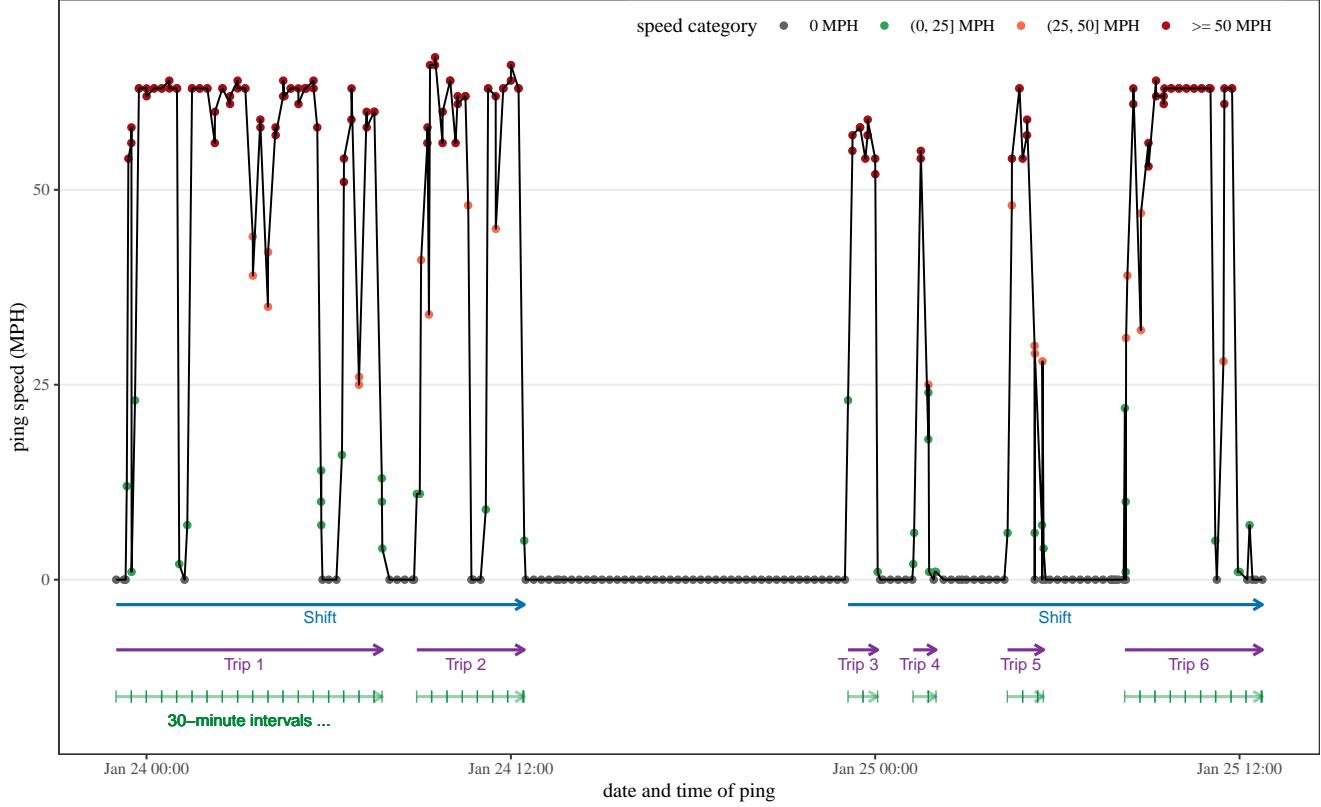


Figure 1: Data aggregation process from pings to shifts, trips, and 30-minute intervals.

represent the aggregated shifts (blue), trips (purple), and 30-minute intervals (green). The long blue arrows (shifts) are separated and defined by long grey dots (more than eight hours) in the middle of the figure. Similarly, the shorter purple arrows are separated and defined by shorter grey dots (greater than half an hour but less than eight hours). The shortest green line segments (30-minute intervals) are defined by the start and end time of the purple arrows, and these 30-minute intervals are much more homogeneous in length than shifts and trips.

4.2. Data fusion

Driver demographic variables were merged to the 30-minute intervals using driver unique IDs. Weather variables were firstly merged to the original ping data using unique latitude, longitude, and time combinations, and then aggregated to 30-minute intervals by taking the average of the weather variables. The SCEs were merged to the 30-minute intervals by matching driver IDs and if the time of the SCEs falls in between the start and end time of the intervals. The data aggregation and fusion process is empowered by the R package `data.table` to leverage its high-speed and multi-thread in-memory modification, aggregation, grouping, and joining performance for large data sets [10]. The code is shared in the Supplementary materials.

163 *4.3. Cumulative driving time as a measure of fatigue*

164 Fatigue is the most important predictor of truck crashes [4, 20, 33]. However, driver fatigue is difficult to measure
165 in real life [16]. In this study, we use cumulative driving time within each shift for each driver as a proxy measure of
166 the fatigue of truck drivers[23]. It is calculated by adding up the 30-minute interval times in each shift for each
167 driver, and the rest time between trips and shifts were not included.

168 **5. Statistical models**

169 Traditional statistical models assume that observations are independent from each other given their predictor
170 variables. However, natural data are almost never independent given the predictor variables. In the example of truck
171 driver's safety events, if we assume the external traffic, weather and driver's socioeconomic status are fixed, truck
172 drivers may exhibit similar driving patterns in multiple trips, and then drivers hired by the same company may
173 share similar culture and safety atmospheres. Therefore, traffic accidents are naturally nested within drivers and
174 drivers are nested within companies. Traditional statistical models that assume independence between observations
175 are not appropriate in this case since objects tend to be similar within a group. Hierarchical models, also known as
176 multilevel model, random-effects model or mixed model, have been developed to allow for the nested nature of data.
177 Instead of assuming independence given predictor variables, hierarchical models assume conditional independence.
178 Hierarchical models are advocated to be the default method since they can produce more precise prediction and
179 more robust results than traditional models. [14, 28]

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:

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

180 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
181 for the cumulative driving time (CT i) in the shift (the sum of driving time for all previous intervals within that
182 shift) for driver $d(i)$. These random intercepts and random slopes are assumed to have a hyper-distribution with
183 hyperparameters $\mu_0, \sigma_0, \mu_1, \sigma_1$. x_2, \dots, x_k are other fixed-effect variables including driver demographics (age, gender,

184 and race), weather (visibility, precipitation intensity and probability), interval specific variables (mean and standard
 185 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:

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

186 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
 187 parameter. Since there is no good solution to estimate the θ parameter here, it was set as a fixed value estimated
 188 from a Poisson regression using maximum likelihood estimation. Other parameters are similar and explained in the
 189 previous hierarchical logistic regression model, and we put a $*$ on the parameter to note the difference between the
 190 parameters of the two models.

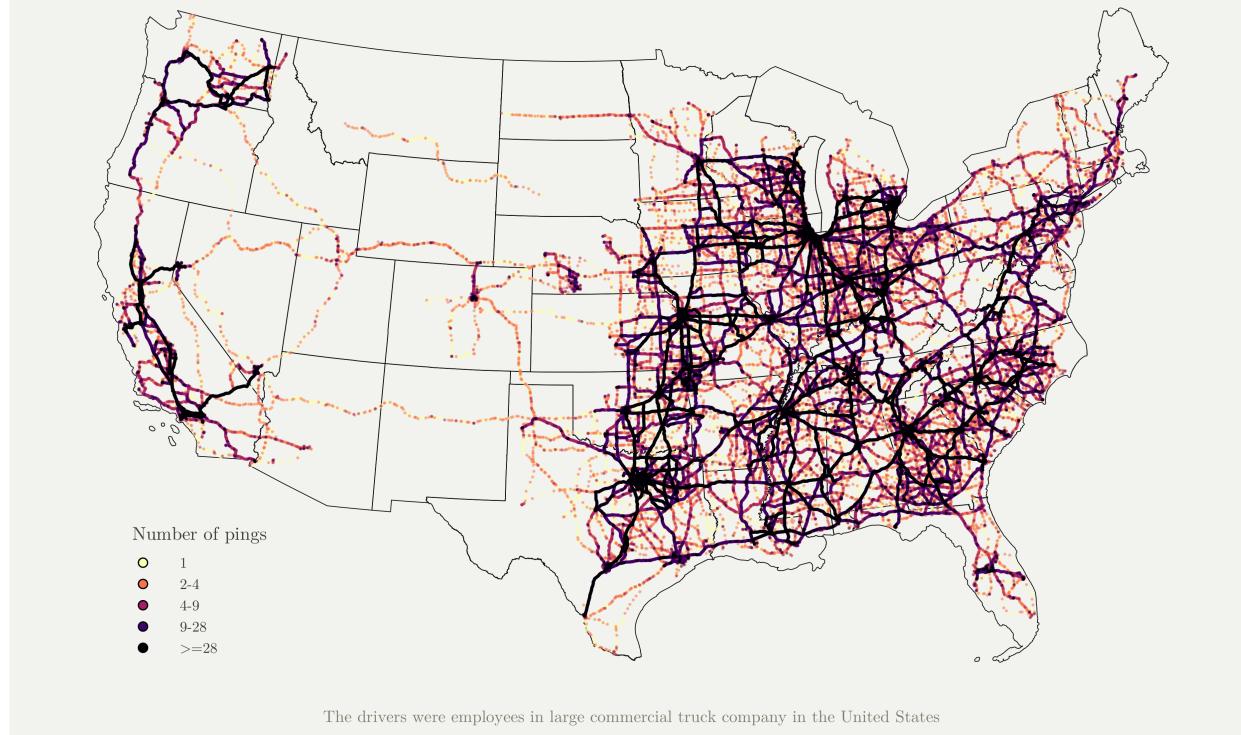
191 To compare with models without driver-level random effects, we also estimated logistic and NB regression models
 192 without any random effects. Log likelihood, the Akaike Information Criterion (AIC), the Bayesian Information
 193 Criterion (BIC), and c -statistic were reported to assess model fit. The hierarchical logistic and NB models were
 194 estimated using the `lme4` R package [1], and model fit statistics were generated using `finalfit` R package [15]. All
 195 the analyses were conducted in statistical computing environment R 3.6.2 [29]. The data and associated R code can
 196 be accessed in the supplementary materials.

197 6. Results and discussion

198 6.1. Geographic distribution of sample pings

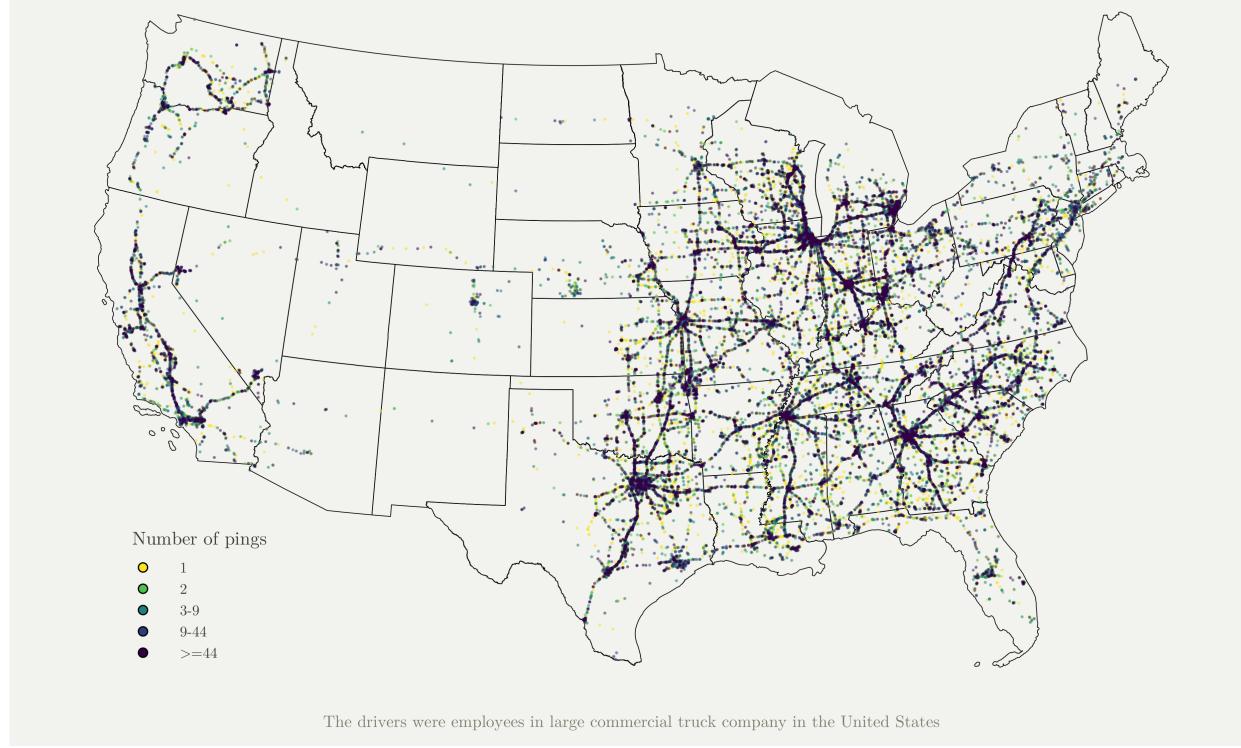
199 Figure 2 demonstrates the geographical point patterns of the actively moving pings (Figure 2a) and stopped
 200 pings (2b) generated by the 497 sample drivers. In both of the two figures, the grey thinner lines are major highways
 201 in the U.S., the black thicker lines are state borders, and darker color represents higher ping density at that location.
 202 The two plots shows that the majority of the transporting tasks was in the middle and east parts, with a few in the
 203 west (California and Seattle), while very few points were in the Midwest. The coverage of locations all around the
 204 U.S. makes the sample in this study generally representative of the regional driving tasks in this country.

Geographical distribution of the **moving pings** generated by the 497 truck drivers, 2015-2016



(a) Active pings

Geographical distribution of the **stopped pings** generated by the 497 truck drivers, 2015-2016



(b) Inactive pings

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

205 *6.2. Exploratory analysis*

206 Figure 3 presents the univariate relationship between cumulative driving time and the rate of SCEs (the number
 207 of SCEs per 0.5 hour). The black points are the rates calculated from the aggregated data, surrounded by 95%
 208 confidence interval grey bands, and the blue curve is the Locally Weighted Scatterplot Smoothing (LOESS) estimates
 209 of the black points. It shows that the rate of SCEs increases as cumulative driving time goes from zero to six hours,
 210 while the trend levels off after six hours of cumulative driving. It worths attention that the magnitude of change in
 211 the y -axis is very small, and this is the raw curve estimate, without adjusting for other variables.

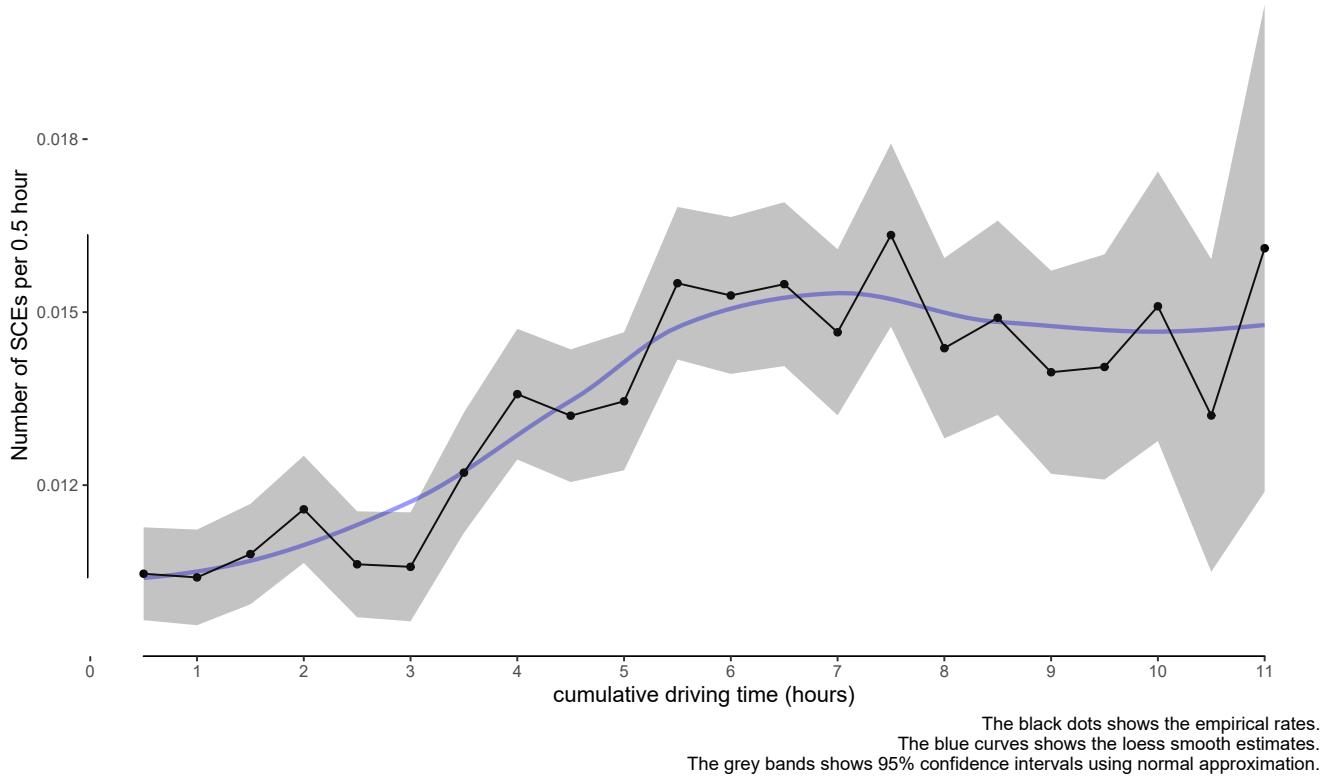


Figure 3: The rate of safety critical events and cumulative driving time

212 *6.3. Statistical models*

213 Table 2 presents the results of the four statistical models: (1) logistic regression without random effects, (2)
 214 NB regression without random effects, (3) hierarchical logistic regression with driver-level random intercepts and
 215 random slopes for cumulative driving time, and (4) hierarchical NB regression with driver-level random intercepts
 216 and random slopes. Compared to model (1) and (2), in which most predictors are significant, the predictors in
 217 model (3) and (4) are less significant. This reduction in the significance of predictors is because the variation of the
 218 outcome variable in model (3) and (4) is explained by the driver-level random effects, instead of other fixed-effect
 219 predictors. In all four models, the estimated parameters for cumulative driving time were not significant and the

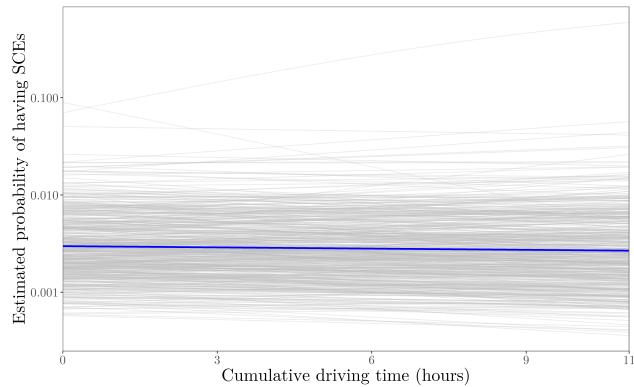


Figure 4: Hierarchical Logistics model

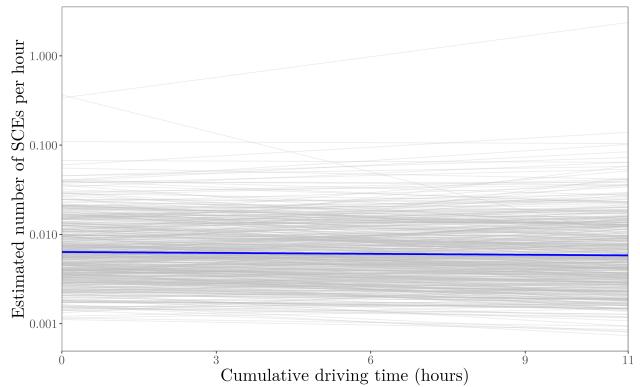


Figure 5: Hierarchical negative binomial model

Figure 6: Simulated relationship between cumulative driving time and probability (logistics model)/rate (negative binomial model) of SCEs the 497 sample drivers. The y -axes are on the log 10 scale.

values were close to zero, indicating that cumulative driving time was not associated with the risk of SCEs among the sample drivers. The estimated values of the hyperparameters (σ_0 and σ_1) were not small, which suggests that there were fair amount of variability across drivers.

To better understand the relationship between cumulative driving time and the risk of SCEs, as well as driver-to-driver variability, we visualized the estimated risk of SCEs and cumulative driving hours for each driver (the grey lines) and the overall trend (the bold blue lines), as shown in Figure 6. It worths noting that the y -axis in the two plots are on the log 10 scale not on a linear scale, which is to avoid an overwhelm of grey lines on the lower part of the plots. Both of the two figures suggest that there seems to be no association between cumulative driving time and SCEs among the sample drivers, although there is fair amount of variability in both the intercept and slope across drivers.

It is known that fatigue, driving time, and work schedule are major risk factors for trucking safety [5, 7, 25, 31, 32]. However, this study found no significant relationship between cumulative driving time and the risk of SCEs. This null relationship could be explained by several reasons. Firstly, the sample 497 regional drivers moved freights within a region that included the surrounding states, and they were on the road for around five days and took breaks on a weekly basis. Their schedule are busier than local drivers but less intensive than over-the-road drivers. Secondly, other crucial variables such as traffic and road geometry were not available in this study, and missing these variables may have nullified the relationship. Thirdly, the data quality and integrity of the third-party weather data provider cannot be validated. The weather stations did not cover every corner of the places traveled and the data of the uncovered places were inferred using computational algorithms.

6.4. Model evaluation

Table 3 presents model fit statistics in the four models. Higher log likelihood values and c-statistics indicate better model fit, while lower AIC and BIC values suggest better model fit. All four model fit statistics suggest that

Table 2: Estimated results for the standard and hierarchical logistic and NB models

	Logistic (1)	NB (2)	Hierarchical logistic (3)	Hierarchical NB (4)
Intercept (μ_0)	-4.691*** (0.094)	-6.985*** (0.084)	-5.812*** (0.235)	-8.459*** (0.237)
Cumulative driving (μ_1)	-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)	0.094 (0.105)	0.096 (0.109)
Race: other	0.235*** (0.042)	0.141*** (0.046)	0.370** (0.179)	0.348* (0.186)
Gender: female	-0.288*** (0.050)	-0.347*** (0.053)	-0.085 (0.184)	-0.086 (0.191)
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
θ		0.036*** (0.001)		0.145
sd: Intercept (σ_0)			0.956	1.01
sd: cumulative driving (σ_1)			0.078	0.084
cor: μ_0 & μ_1			-0.222	-0.262

Note:

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

Table 3: Model fit statistics for the standard and hierarchical logistic and NB models

Fit statistics	Logistic	NB	Hierarchical logistic	Hierarchical NB
Log likelihood	-46,304	-49,627	-43,042	-45,961
AIC	92,634	99,280	86,117	91,954
BIC	92,788	99,434	86,306	92,144
c-statistic	0.590	0.571	0.760	0.740

242 the hierarchical logistic regression model has the best fit among the four models. Adding driver-level random effects
 243 substantially improved the model fit statistics, with *c*-statistics increased by 0.17 and 0.169 for the logistic and NB
 244 regression models.

245 Although the model fit can be improved substantially by adding driver-level random effects, we should acknowledge
 246 that the models are generally underfitting. The models without driver-level random effects have *c*-statistics of 0.59
 247 and 0.71, which are only slightly higher than a random classification model that has the *c*-statistics of 0.5. Even for
 248 the best fit model, the *c*-statistic is only 0.76, which is good but not strong enough (a model with the *c*-statistics
 249 of 0.8 is usually viewed as a strong model). The model fit statistics could be further improved by adding other
 250 important predictors such as traffic and road geometry, which are current not available or accessible to the research
 251 team.

252 7. Conclusions and implications

253 7.1. Summary and impact

254 This paper provides a preliminary analysis of the association between cumulative driving time and the risk
 255 of SCEs among 497 commercial truck drivers, using a driver-centric analysis framework. To accomplish this, we
 256 pulled weather data from a third-party data provider, merged four different sources of data (ping, SCEs, driver
 257 demographics, and weather), and aggregated the fused data into shifts, trips, and 30-minute intervals. Exploratory
 258 analysis indicated that the rate of SCEs has a slightly increasing and non-linear (concave down) trend as cumulative
 259 driving time increases. However, hierarchical logistic and negative binomial models with driver-specific random
 260 intercepts and slopes showed no significant association between cumulative driving time and th risk of SCEs.

261 Although this case study is based on NDS data generated from large commercial truck drivers, we argue that the
 262 data collection, aggregation, fusing, and driver-centric statistical modeling framework is generalizable to other types
 263 of drivers. The sensors to collect ping and SCEs data are similar for different vehicles, and the third-party data
 264 sources can be merged onto ping data using space (GPS locations) and time information.

265 The advantages of this driver-centric analysis framework.

266 *7.2. Limitations and future work*

267 Although the sample size in this study is relatively large (~500 truck drivers), there are several aspects that can
268 be explored in future studies. Future studies can similar analysis framework but different data sets to check the
269 consistency of conclusions. More importantly, over-the-road truck drivers deserve special attention as they have
270 more long-distance tasks and are subject to fatigue and falling asleep issues at wheel [22]. Besides, as more insurance
271 companies are routinely collecting ping and critical events data similar to ours, researchers are encourage to conduct
272 analysis on passenger-carrying vehicle using this proposed driver-centric analysis framework. A large sample size can
273 yield more statistical power and robust results, and as a consequence, provide high-quality evidence to optimize
274 scheduling, empower policy making, and improve transportation safety.

275 **Supplementary Materials**

276 To make this study replicable and promote future studies using similar data sets, we provided our R code
277 associated with data cleaning, aggregation, merging, statistical modeling and validation, and data visualization, as
278 well as a portion of data on a GitHub repository, which can be accessed at [GitHub repo](#). An associated website is
279 created to briefly introduce our workflow and preliminary results.

280 **Acknowledgement**

281 The research work presented in this study was supported in part by the National Science Foundation (CMMI-
282 1635927 and CMMI-1634992), the Ohio Supercomputer Center (PMIU0138 and PMIU0162), the American Society of
283 Safety Professionals (ASSP) Foundation, the University of Cincinnati Education and Research Center Pilot Research
284 Project Training Program, and the Transportation Informatics Tier I University Transportation Center (TransInfo).
285 We also thank the DarkSky company for providing us five million free calls to their historic weather API.

286 **References**

- 287 [1] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. Fitting linear mixed-effects models using lme4.
288 *Journal of Statistical Software*, 67(1):1–48, 2015. doi: 10.18637/jss.v067.i01.
- 289 [2] Daniel Blower, Paul E Green, and Anne Matteson. Condition of trucks and truck crash involvement: Evidence
290 from the large truck crash causation study. *Transportation Research Record*, 2194(1):21–28, 2010.
- 291 [3] David E Cantor, Thomas M Corsi, Curtis M Grimm, and Koray Özpolat. A driver focused truck crash prediction
292 model. *Transportation Research Part E: Logistics and Transportation Review*, 46(5):683–692, 2010.
- 293 [4] Lora Cavuoto, Fadel Megahed, et al. Understanding fatigue and the implications for worker safety. In *ASSE*
294 *Professional Development Conference and Exposition*. American Society of Safety Engineers, 2016.

- 295 [5] Chen Chen and Yuanchang Xie. Modeling the safety impacts of driving hours and rest breaks on truck drivers
296 considering time-dependent covariates. *Journal of safety research*, 51:57–63, 2014.
- 297 [6] Cong Chen, Guohui Zhang, Zong Tian, Susan M Bogus, and Yin Yang. Hierarchical bayesian random intercept
298 model-based cross-level interaction decomposition for truck driver injury severity investigations. *Accident
299 Analysis & Prevention*, 85:186–198, 2015.
- 300 [7] Guang Xiang Chen, Youjia Fang, Feng Guo, and Richard J Hanowski. The influence of daily sleep patterns of
301 commercial truck drivers on driving performance. *Accident analysis & prevention*, 91:55–63, 2016.
- 302 [8] Thomas A Dingus, Richard J Hanowski, and Sheila G Klauer. Estimating crash risk. *Ergonomics in Design*, 19
303 (4):8–12, 2011.
- 304 [9] Chunjiao Dong, Qiao Dong, Baoshan Huang, Wei Hu, and Shashi S Nambisan. Estimating factors contributing
305 to frequency and severity of large truck-involved crashes. *Journal of Transportation Engineering, Part A:
306 Systems*, 143(8):04017032, 2017.
- 307 [10] Matt Dowle and Arun Srinivasan. *data.table: Extension of ‘data.frame’*, 2019. URL <https://CRAN.R-project.org/package=data.table>. R package version 1.12.6.
- 308 [11] Federal Motor Carrier Safety Administration. Summary of hours of service regulations. <https://cms8.fmcsa.dot.gov/regulations/hours-service/summary-hours-service-regulations>, 2013. [Online; accessed 06-February-2020].
- 311 [12] Feng Guo. Statistical methods for naturalistic driving studies. *Annual Review of Statistics and Its Application*,
312 6:309–328, 2019.
- 313 [13] Feng Guo, Sheila G Klauer, Jonathan M Hankey, and Thomas A Dingus. Near crashes as crash surrogate for
314 naturalistic driving studies. *Transportation Research Record*, 2147(1):66–74, 2010.
- 315 [14] Chunyang Han, Helai Huang, Jaeyoung Lee, and Jie Wang. Investigating varying effect of road-level factors on
316 crash frequency across regions: a bayesian hierarchical random parameter modeling approach. *Analytic methods
317 in accident research*, 20:81–91, 2018.
- 318 [15] Ewen Harrison, Tom Drake, and Riinu Ots. *finalfit: Quickly Create Elegant Regression Results Tables and Plots
319 when Modelling*, 2019. URL <https://CRAN.R-project.org/package=finalfit>. R package version 0.9.7.
- 320 [16] Laurence R Hartley, Pauline K Arnold, G Smythe, and J Hansen. Indicators of fatigue in truck drivers. *Applied
321 Ergonomics*, 25(3):143–156, 1994.
- 322 [17] Jeffrey S Hickman, Richard J Hanowski, and Joseph Bocanegra. A synthetic approach to compare the large
323 truck crash causation study and naturalistic driving data. *Accident Analysis & Prevention*, 112:11–14, 2018.

- 324 [18] Carl Johnsson, Aliaksei Laureshyn, and Tim De Ceunynck. In search of surrogate safety indicators for vulnerable
325 road users: a review of surrogate safety indicators. *Transport reviews*, 38(6):765–785, 2018.
- 326 [19] SM Sohel Mahmud, Luis Ferreira, Md Shamsul Hoque, and Ahmad Tavassoli. Application of proximal surrogate
327 indicators for safety evaluation: A review of recent developments and research needs. *IATSS research*, 41(4):
328 153–163, 2017.
- 329 [20] Zahra Sedighi Maman, Mohammad Ali Alamdar Yazdi, Lora A Cavuoto, and Fadel M Megahed. A data-driven
330 approach to modeling physical fatigue in the workplace using wearable sensors. *Applied ergonomics*, 65:515–529,
331 2017.
- 332 [21] Fred L Mannering and Chandra R Bhat. Analytic methods in accident research: Methodological frontier and
333 future directions. *Analytic methods in accident research*, 1:1–22, 2014.
- 334 [22] Anne T McCartt, John W Rohrbaugh, Mark C Hammer, and Sandra Z Fuller. Factors associated with falling
335 asleep at the wheel among long-distance truck drivers. *Accident Analysis & Prevention*, 32(4):493–504, 2000.
- 336 [23] Peter McCauley, Leonid V Kalachev, Daniel J Mollicone, Siobhan Banks, David F Dinges, and Hans PA
337 Van Dongen. Dynamic circadian modulation in a biomathematical model for the effects of sleep and sleep loss
338 on waking neurobehavioral performance. *Sleep*, 36(12):1987–1997, 2013.
- 339 [24] Lynn Meuleners, Michelle L Fraser, Matthew H Govorko, and Mark R Stevenson. Determinants of the
340 occupational environment and heavy vehicle crashes in western australia: A case–control study. *Accident
341 Analysis & Prevention*, 99:452–458, 2017.
- 342 [25] Daniel Mollicone, Kevin Kan, Chris Mott, Rachel Bartels, Steve Bruneau, Matthew van Wollen, Amy R Sparrow,
343 and Hans PA Van Dongen. Predicting performance and safety based on driver fatigue. *Accident Analysis &
344 Prevention*, 126:142–145, 2019.
- 345 [26] Bhaven Naik, Li-Wei Tung, Shanshan Zhao, and Aemal J Khattak. Weather impacts on single-vehicle truck
346 crash injury severity. *Journal of Safety Research*, 58:57–65, 2016.
- 347 [27] National Highway Traffic Safety Administration. A Compilation of Motor Vehicle Crash Data from the Fatality
348 Analysis Reporting System and the General Estimates System. [https://crashstats.nhtsa.dot.gov/Api/Public/
349 ViewPublication/812384](https://crashstats.nhtsa.dot.gov/Api/Public/ViewPublication/812384), 2017. [Online; accessed 23-November-2019].
- 350 [28] Sarvani Sonduru Pantangi, Grigoris Fountas, Md Tawfiq Sarwar, Panagiotis Ch Anastopoulos, Alan Blatt,
351 Kevin Majka, John Pierowicz, and Satish B Mohan. A preliminary investigation of the effectiveness of high
352 visibility enforcement programs using naturalistic driving study data: A grouped random parameters approach.
353 *Analytic Methods in Accident Research*, 21:1–12, 2019.

- 354 [29] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
355 Computing, Vienna, Austria, 2019. URL <https://www.R-project.org/>.
- 356 [30] Lisa N Sharwood, Jane Elkington, Lynn Meuleners, Rebecca Ivers, Soufiane Boufous, and Mark Stevenson. Use
357 of caffeinated substances and risk of crashes in long distance drivers of commercial vehicles: case-control study.
358 *BMJ*, 346:f1140, 2013.
- 359 [31] Susan A Soccilich, Myra Blanco, Richard J Hanowski, Rebecca L Olson, Justin F Morgan, Feng Guo, and
360 Shih-Ching Wu. An analysis of driving and working hour on commercial motor vehicle driver safety using
361 naturalistic data collection. *Accident Analysis & Prevention*, 58:249–258, 2013.
- 362 [32] Amy R Sparrow, Daniel J Mollicone, Kevin Kan, Rachel Bartels, Brieann C Satterfield, Samantha M Riedy,
363 Aaron Unice, and Hans PA Van Dongen. Naturalistic field study of the restart break in us commercial motor
364 vehicle drivers: truck driving, sleep, and fatigue. *Accident Analysis & Prevention*, 93:55–64, 2016.
- 365 [33] Hal S Stern, Daniel Blower, Michael L Cohen, Charles A Czeisler, David F Dinges, Joel B Greenhouse, Feng
366 Guo, Richard J Hanowski, Natalie P Hartenbaum, Gerald P Krueger, et al. Data and methods for studying
367 commercial motor vehicle driver fatigue, highway safety and long-term driver health. *Accident Analysis &*
368 *Prevention*, 126:37–42, 2019.
- 369 [34] The Dark Sky Company, LLC. Dark Sky API — Overview. <https://darksky.net/dev/docs>, 2019. [Online;
370 accessed 20-February-2019].
- 371 [35] Majbah Uddin and Nathan Huynh. Truck-involved crashes injury severity analysis for different lighting conditions
372 on rural and urban roadways. *Accident Analysis & Prevention*, 108:44–55, 2017.
- 373 [36] Behram Wali, Asad J Khattak, and Thomas Karnowski. Exploring microscopic driving volatility in naturalistic
374 driving environment prior to involvement in safety critical events—concept of event-based driving volatility.
375 *Accident Analysis & Prevention*, 132:105277, 2019.
- 376 [37] WHO. The top 10 causes of death. <http://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death>, 2018. [Online; accessed 20-February-2019].
- 377 [38] Xiaoyu Zhu and Sivaramakrishnan Srinivasan. A comprehensive analysis of factors influencing the injury severity
378 of large-truck crashes. *Accident Analysis & Prevention*, 43(1):49–57, 2011.