

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**

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 (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 vehicle-related fatalities (Hickman et al., 2018). Truck drivers are often on the road for long routes
18 under on-time demands, complex traffic and weather conditions, with little to no supervision and contact with fellow
19 workers. Therefore, trucking safety is an important research topic and a number of studies have been published to
20 predict and reduce crash risk associated with trucks (Cantor et al., 2010; Chen et al., 2015; Dong et al., 2017).

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 (Blower et al., 2010; Meuleners et al., 2017; Sharwood et al., 2013). This
25 case-control study design is efficient and less time-consuming in the field of trucking safety since crashes are very
26 rare. However, case-control studies, by nature, are limited in study design since a) it is impossible to estimate and
27 compare the rate of crashes since the number of non-crashes is unknown, b) retrospective reports are often subject
28 to recall and report bias: the drivers may not accurately recall the exact conditions at the time of the event, c) the

*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)

29 drivers may intentionally conceal some critical facts to escape from legal punishment (Dingus et al., 2011; Stern et
30 al., 2019).

31 Naturalistic driving studies (NDSs) have been emerging in the past decade thanks to the advancement of
32 technology. An NDS continuously collects driving data (including latitude, longitude, and speed) under real-world
33 conditions using on-board unobtrusive equipment (Guo, 2019). In contrast to retrospective reports, an NDS resembles
34 a cohort study: a pre-determined set of drivers are prospectively followed for a certain amount of time. Therefore,
35 NDS has several advantages. First, NDS collects both crashes and non-crashes, so it is more useful in comparing the
36 rates of events. Second, since vehicle crashes are extremely rare, it may take a huge amount of driving time to have
37 sufficient sample of crashes. Instead, NDS focus safety-critical events (SCEs), which is defined as events that avoid
38 crashes by last-second evasive maneuver (Dingus et al., 2011). SCEs can be 1000 times as high as real crashes and
39 are argued to be good surrogates of crashes (Dingus et al., 2011; Guo et al., 2010; Johnsson et al., 2018; Mahmud et
40 al., 2017). Third, NDS data are collected using programmed instruments or sensors, so they are less likely to be
41 subject to human error, recall bias, or misinformation. Lastly, NDS collects data every a few seconds to minutes,
42 and this large-scale high-resolution data provide a promising opportunity to quantifying driving risk (Guo, 2019).

43 However, many issues arise given the characteristics of NDSs. First, the sheer volume of NDS data creates a
44 challenge to data management and aggregation (Mannering and Bhat, 2014). For example, a NDS data set can have
45 billions rows of real-time speeds and locations, and it is important to have scalable and high-performance tools to
46 aggregate these data into units that fit into the framework of statistical modeling. Second, routinely collected NDS
47 data only have vehicle driving data. Crucial environmental variables such as weather and traffic need to be accessed
48 from other data sources and merged back to the driving data. Third, even with these data sources, management,
49 and aggregation issues solved, there is a lack of consensus on choosing the statistical models that are both sufficiently
50 complex to account for the characteristics of NDS and computationally feasible to fit the large-scale data. With
51 increasing companies collecting NDS data on a regular basis, a scalable and generalizable analyzing framework can
52 serve as a pattern for researchers to better understand NDS data and gain insights into trucking and transportation
53 safety.

54 This paper aims to propose and showcase a generalizable data analytic framework (data collecting, aggregating,
55 fusing, and statistical modeling) that accounts for the features of NDS data. To achieve this aim, we have answers
56 the following questions:

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

60 The remainder of this paper is organized as follows. Section 2 provides a brief literature review on previously
61 published studies that use NDS data sets. Then, Section 3 presents our NDS data and other third-party data

sources. Section 4 demonstrates how we aggregate the ping data into shifts, trips, and 30-minute intervals, and merge different data sources. Section 5 details the driver hierarchical logistic and negative binomial model. Section 6 presents the statistical results and interpretation and conclusions and implications are discussed in Section 7.

2. Literature review

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

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

- (A) The number of sample drivers are small (around 100 drivers) except for Wali et al. (2019). The studies may not have sufficient statistical power due to the small sample size, and the generalizability may be limited.
- (B) The data sources come from only NDS data sets, which increases the workload and difficulty of data collection. In secondary data analysis, exclusively replying on one data source may limit our power to answer the question. With various organizations collecting data, we can exploit the power of third-party data providers, integrate different sources of data, and as a consequence, improve the prediction accuracy of statistical models.
- (C) Although the listed papers occasionally used hierarchical models, relatively few actually used driver-centric hierarchical models. NDS data sets are naturally generated by a driver-centric process: recruited drivers are followed for a certain amount of time, and all relevant data are collected in this process.
- (D) No consistent framework for cleaning, aggregating, and statistical modeling has been proposed in these papers. NDS data sets collects large-scale high-resolutional data, which rely on a context-specific, statistically sensible, and computationally affordable to analyze and empower policy making.

Our study serves as a complement to the existing literature and a model for future NDS studies. Firstly, we combined routinely collected driving data of 497 commercial truck drivers and a third-party weather data provider. Then, a contextual sensible data aggregation framework is proposed to reduce the original driving data to shifts, trips, and 30-minute intervals. Lastly, we propose to use driver-centric mixed-effect statistical models to analyze the 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
Soccolich	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	
Ghasemzadeh	2018	141 general drivers	lane-keeping behavior	driver characteristics, weather and traffic conditions	logistic regression and multivariate adaptive regression splines	Traffic, age and experience, and speed limits are significant factors
Zhu	2018	42 general drivers	car-following	NA	Gazis-Herman-Rothery, Gipps, intelligent driver, full velocity difference, and 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

92 **3. Data sources**

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

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

104 *3.1. Ping and SCEs data*

105 The DAGs ping irregularly (typically every a couple of seconds to minutes) as the truck goes on road. Each ping
106 collects several key variables, including the date and time (year, month, day, hour, minute, and second), latitude
107 and longitude (specific to five decimal places), driver identification number (ID), and speed at that second. In total,
108 13,187,289 rows of ping data were generated by the 497 truck drivers, with 8,029,087 (60.89%) of them were active
109 pings (speed of the ping is not zero).

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

115 *3.2. Weather*

116 Weather is one of the most studied risk factors associated with trucking safety (Naik et al., 2016; Uddin and
117 Huynh, 2017; Zhu and Srinivasan, 2011). In this study, we obtained historic weather data from the DarkSky
118 Application Programming Interface (API), which allows us to query historic real-time and hour-by-hour nationwide
119 historic weather conditions according to latitude, longitude, date, and time (The Dark Sky Company, LLC, 2019).
120 The primary weather variables included visibility, precipitation probability¹, precipitation intensity, wind speed,

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.

121 and others. To reduce the cost of querying all 13 million ping data from the DarkSky API, we rounded the GPS
122 coordinates to the second decimal places, which are worth up to 1.1 kilometers, and we also round the time to the
123 nearest hour. Then the weather variables were queried from the DarkSky API using the approximated latitudes,
124 longitudes, date and hour.

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

131 4. Data preparation

132 4.1. Shifts, trips, and 30-minute intervals

133 To convert this 13 million row real-time ping data into analyzable units, we aggregate them into *shifts*, *trips*, and
134 *30-minute intervals*, which are inspired by real world truck transporting practice and the hours-of-service policy by
135 Federal Motor Carrier Safety Administration (2013). Shifts are on-duty periods with no breaks longer than eight
136 hours (there can be short breaks less than 8 hours). Trips are continuous driving periods with no breaks less than
137 half an hour. These trips are further divided into 30-minute fixed intervals. This is because trips can vary from
138 several minutes to several hours, which are not a good analyzable unit for statistical modeling. The details of the
139 aggregation process is as follows:

- 140 • *Shifts*: for each of the sample truck drivers, if the ping data showed that the truck was not moving for more
141 than eight hours, the ping data were separated into two different shifts on the left and right side of this long
142 break. There could be several short breaks (less than eight hours) within each shift.
- 143 • *Trips*: for each shift, if the ping data showed that the truck was not moving for more than half an hour, the
144 ping data were separated into different trips. These ping data were then aggregated into different trips. The
145 drivers are assumed to be fully driving within each trip since there are not breaks longer than 30 minutes
146 within each trip. The trips are nested within shifts.
- 147 • *30-minute intervals*: each trip is further decomposed into 30-minute fixed intervals according to the start and
148 end time of the trip. The last interval of the trip is typically less than 30 minutes. The 30-minute intervals are
149 nested within trips.

150 Figure 1 visually present the data aggregation process of ping → shifts → trips → 30-minute intervals, as well as
151 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

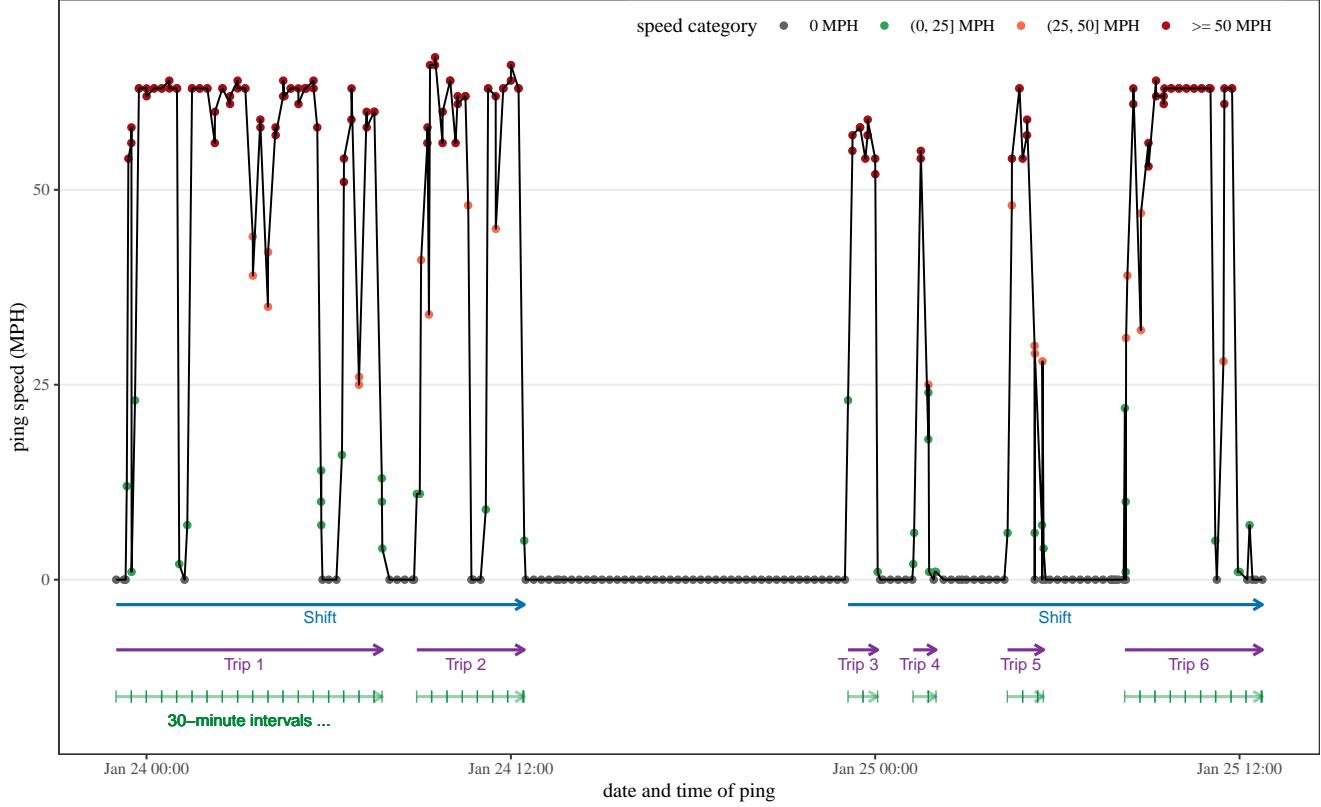


Figure 1: Data aggregation process from pings to shifts, trips, and 30-minute intervals.

152 the current speed. Grey dots indicate stopping pings with the current speed of zero. The arrows in the lower part
 153 represent the aggregated shifts (blue), trips (purple), and 30-minute intervals (green). The long blue arrows (shifts)
 154 are separated and defined by long grey dots (more than eight hours) in the middle of the figure. Similarly, the
 155 shorter purple arrows are separated and defined by shorter grey dots (greater than half an hour but less than eight
 156 hours). The shortest green line segments (30-minute intervals) are defined by the start and end time of the purple
 157 arrows, and these 30-minute intervals are much more homogeneous in length than shifts and trips.

158 4.2. Data fusion

159 Driver demographic variables were merged to the 30-minute intervals using driver unique IDs. Weather variables
 160 were firstly merged to the original ping data using unique latitude, longitude, and time combinations, and then
 161 aggregated to 30-minute intervals by taking the average of the weather variables. The SCEs were merged to the
 162 30-minute intervals by matching driver IDs and if the time of the SCEs falls in between the start and end time of
 163 the intervals. The data aggregation and fusion process is conducted using the R package `data.table` to leverage its
 164 high-speed in-memory aggregation, grouping, and joining performance for large data (Dowle and Srinivasan, 2019),
 165 and the code is shared in the Supplementary materials.

166 4.3. Cumulative driving time as a measure of fatigue

167 Fatigue is the most important predictor of truck crashes (Cavuoto et al., 2016; Maman et al., 2017; Stern et
168 al., 2019). However, driver fatigue is difficult to measure in real life (Hartley et al., 1994). In this study, we use
169 cumulative driving time within each shift for each driver as a proxy measure of the fatigue of truck drivers(McCauley
170 et al., 2013). It is calculated by adding up the 30-minute interval times in each shift for each driver, and the rest
171 time between trips and shifts were not included.

172 5. Statistical models

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

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

184 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
185 for the cumulative driving time (CT_i) in the shift (the sum of driving time for all previous intervals within that
186 shift) for driver $d(i)$. These random intercepts and random slopes are assumed to have a hyper-distribution with

187 hyperparameters $\mu_0, \sigma_0, \mu_1, \sigma_1$. x_2, \dots, x_k are other fixed-effect variables including driver demographics (age, gender,
 188 and race), weather (visibility, precipitation intensity and probability), interval specific variables (mean and standard
 189 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:

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

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

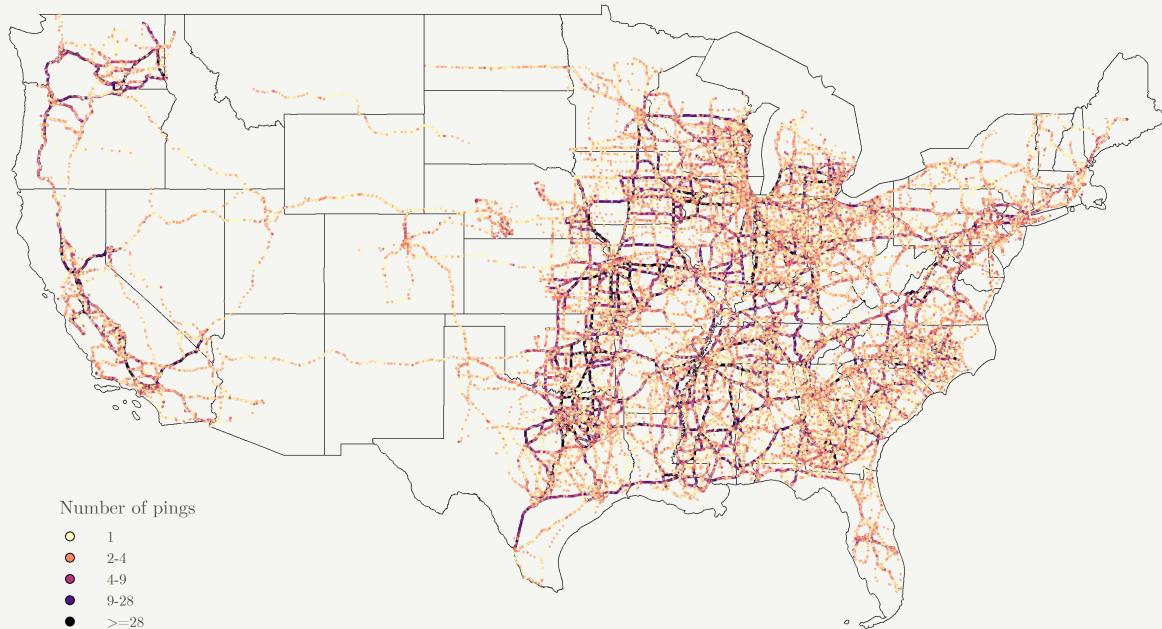
195 To compare with models without driver-level random effects, we also estimated logistic and NB regression models
 196 without any random effects. Log likelihood, the Akaike Information Criterion (AIC), the Bayesian Information
 197 Criterion (BIC), and c -statistic were reported to assess model fit. The hierarchical logistic and NB models were
 198 estimated using the `lme4` R package (Bates et al., 2015), and model fit statistics were generated using `finalfit` R
 199 package (Harrison et al., 2019). All the analyses were conducted in statistical computing environment R 3.6.2 (R
 200 Core Team, 2019). The data and associated R code can be accessed in the supplementary materials.

201 6. Results and discussion

202 6.1. Geographic distribution of sample pings

203 Figure 2 demonstrates the geographical point patterns of the actively moving pings (Figure 2a) and stopped
 204 pings (2b) generated by the 497 sample drivers. In both of the two figures, the grey thinner lines are major highways
 205 in the U.S., the black thicker lines are state borders, and darker color represents higher ping density at that location.
 206 The two plots shows that the majority of the transporting tasks was in the middle and east parts, with a few in the
 207 west (California and Seattle), while very few points were in the Midwest. The coverage of locations all around the
 208 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 drivers, 2015-2016
The drivers were employees in large commercial truck company in the United States



(a) Active pings

Geographical distribution of the stopped pings generated by the 497 drivers, 2015-2016
The drivers were employees in large commercial truck company in the United States



(b) Inactive pings

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

209 6.2. Statistical models

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

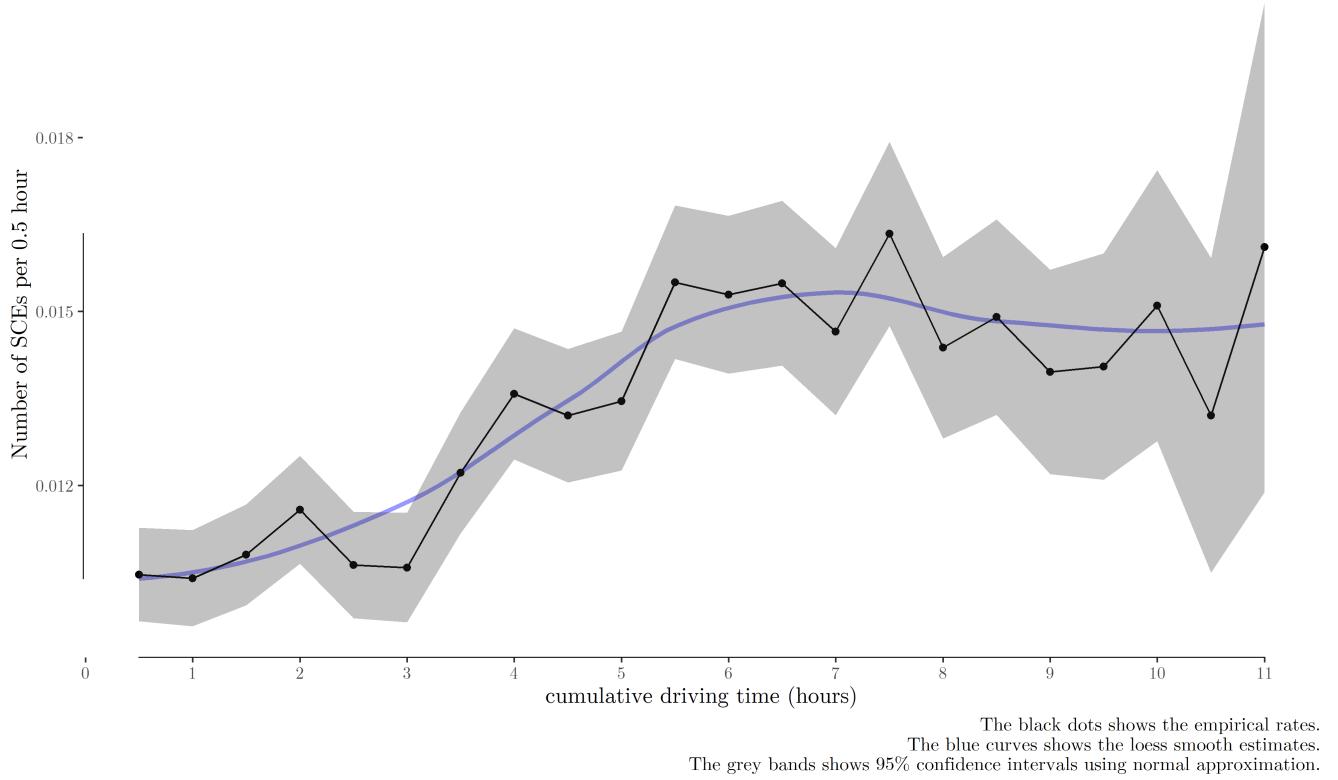


Figure 3: The rate of safety critical events and cumulative driving time

216 Table 2 presents the results of the four statistical models: (1) logistic regression without random effects, (2)
 217 NB regression without random effects, (3) hierarchical logistic regression with driver-level random intercepts and
 218 random slopes for cumulative driving time, and (4) hierarchical NB regression with driver-level random intercepts
 219 and random slopes. Compared to model (1) and (2), in which most predictors are significant, the predictors in
 220 model (3) and (4) are less significant. This reduction in the significance of predictors is because the variation of the
 221 outcome variable in model (3) and (4) is explained by the driver-level random effects, instead of other fixed-effect
 222 predictors. In all four models, the estimated parameters for cumulative driving time were not significant and the
 223 values were close to zero, indicating that cumulative driving time was not associated with the risk of SCEs among
 224 the sample drivers. The estimated values of the hyperparameters (σ_0 and σ_1) were not small, which suggests that

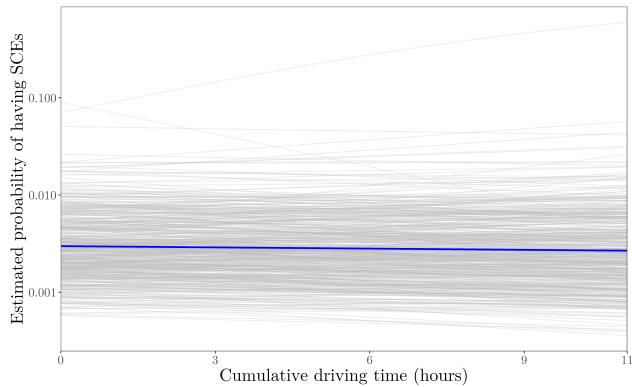


Figure 4: Hierarchical Logistics model

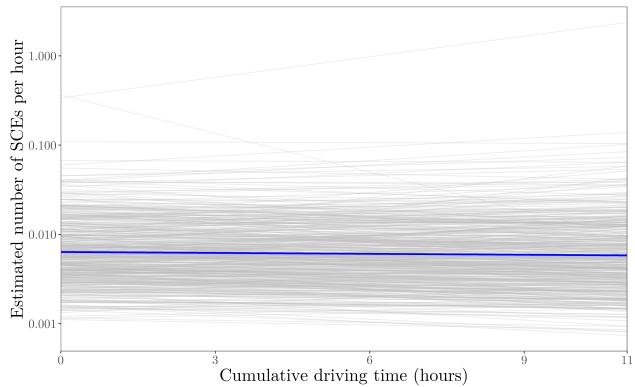


Figure 5: Hierarchical negative binomial model

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.

225 there were fair amount of variability across drivers.

226 To better understand the relationship between cumulative driving time and the risk of SCEs, as well as driver-to-
227 driver variability, we visualized the estimated risk of SCEs and cumulative driving hours for each driver (the grey
228 lines) and the overall trend (the bold blue lines), as shown in Figure 6. It worths noting that the y -axis in the two
229 plots are on the log 10 scale to avoid an overwhelm of grey lines on the lower part of the plots. Both of the two
230 figures suggest that there seems to be no association between cumulative driving time and SCEs among the sample
231 drivers, although there is fair amount of variability in both the intercept and slope across drivers.

232 Possible explantion why there is no relationship, and possibly the hour-of-service rule.

233 6.3. Model evaluation

234 Table 3 presents model fit statistics in the four models. Higher log likelihood values and c-statistics indicate
235 better model fit, while lower AIC and BIC values suggest better model fit. All four model fit statistics suggest that
236 the hierarchical logistic regression model has the best fit among the four models. Adding driver-level random effects
237 substantially improved the model fit statistics, with c -statistics increased by 0.17 and 0.169 for the logistic and NB
238 regression models.

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

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

246 7. Conclusions and implications

247 This paper provides a preliminary analysis of the association between cumulative driving time and the risk
 248 of SCEs among 497 commercial truck drivers, using a driver-centric analysis framework. To accomplish this, we
 249 pulled weather data from a third-party data provider, merged four different sources of data (ping, SCEs, driver
 250 demographics, and weather), and aggregated the fused data into shifts, trips, and 30-minute intervals. Hierarchical
 251 logistic and negative binomial models were used to explore the relationship between cumulative driving time and
 252 SCEs.

253 Although this case study is based on NDS data generated from large commercial truck drivers, we argue that the
 254 driver-centric data collection, aggregation, fusing, and statistical modeling framework is generalizable to other types
 255 of drivers since the original ping and SCEs data are similar among different types of drivers.

256 **Hour-of-service rule.**

257 Acknowledgement

258 The research work presented in this study was supported in part by the National Science Foundation (CMMI-
 259 1635927 and CMMI-1634992), the Ohio Supercomputer Center (PMIU0138 and PMIU0162), the American Society of
 260 Safety Professionals (ASSP) Foundation, the University of Cincinnati Education and Research Center Pilot Research
 261 Project Training Program, and the Transportation Informatics Tier I University Transportation Center (TransInfo).
 262 We also thank the DarkSky company for providing us five million free calls to their historic weather API.

263 References

- 264 Bates, D., Mächler, M., Bolker, B., Walker, S., 2015. Fitting linear mixed-effects models using lme4. *Journal of*
 265 *Statistical Software* 67, 1–48. <https://doi.org/10.18637/jss.v067.i01>
- 266 Blower, D., Green, P.E., Matteson, A., 2010. Condition of trucks and truck crash involvement: Evidence from
 267 the large truck crash causation study. *Transportation Research Record* 2194, 21–28.
- 268 Cantor, D.E., Corsi, T.M., Grimm, C.M., Özpolat, K., 2010. A driver focused truck crash prediction model.
 269 *Transportation Research Part E: Logistics and Transportation Review* 46, 683–692.

- 270 Cavuoto, L., Megahed, F., others, 2016. Understanding fatigue and the implications for worker safety, in: ASSE
271 Professional Development Conference and Exposition. American Society of Safety Engineers.
- 272 Chen, C., Zhang, G., Tian, Z., Bogus, S.M., Yang, Y., 2015. Hierarchical bayesian random intercept model-based
273 cross-level interaction decomposition for truck driver injury severity investigations. Accident Analysis & Prevention
274 85, 186–198.
- 275 Dingus, T.A., Hanowski, R.J., Klauer, S.G., 2011. Estimating crash risk. Ergonomics in Design 19, 8–12.
- 276 Dong, C., Dong, Q., Huang, B., Hu, W., Nambisan, S.S., 2017. Estimating factors contributing to frequency and
277 severity of large truck-involved crashes. Journal of Transportation Engineering, Part A: Systems 143, 04017032.
- 278 Dowle, M., Srinivasan, A., 2019. Data.table: Extension of ‘data.frame’.
- 279 Federal Motor Carrier Safety Administration, 2013. Summary of hours of service regulations.
- 280 Guo, F., 2019. Statistical methods for naturalistic driving studies. Annual Review of Statistics and Its Application
281 6, 309–328.
- 282 Guo, F., Klauer, S.G., Hankey, J.M., Dingus, T.A., 2010. Near crashes as crash surrogate for naturalistic driving
283 studies. Transportation Research Record 2147, 66–74.
- 284 Han, C., Huang, H., Lee, J., Wang, J., 2018. Investigating varying effect of road-level factors on crash frequency
285 across regions: A bayesian hierarchical random parameter modeling approach. Analytic methods in accident research
286 20, 81–91.
- 287 Harrison, E., Drake, T., Ots, R., 2019. Finalfit: Quickly create elegant regression results tables and plots when
288 modelling.
- 289 Hartley, L.R., Arnold, P.K., Smythe, G., Hansen, J., 1994. Indicators of fatigue in truck drivers. Applied
290 Ergonomics 25, 143–156.
- 291 Hickman, J.S., Hanowski, R.J., Bocanegra, J., 2018. A synthetic approach to compare the large truck crash
292 causation study and naturalistic driving data. Accident Analysis & Prevention 112, 11–14.
- 293 Johnsson, C., Laureshyn, A., De Ceunynck, T., 2018. In search of surrogate safety indicators for vulnerable road
294 users: A review of surrogate safety indicators. Transport reviews 38, 765–785.
- 295 Mahmud, S.S., Ferreira, L., Hoque, M.S., Tavassoli, A., 2017. Application of proximal surrogate indicators for
296 safety evaluation: A review of recent developments and research needs. IATSS research 41, 153–163.
- 297 Maman, Z.S., Yazdi, M.A.A., Cavuoto, L.A., Megahed, F.M., 2017. A data-driven approach to modeling physical
298 fatigue in the workplace using wearable sensors. Applied ergonomics 65, 515–529.
- 299 Mannering, F.L., Bhat, C.R., 2014. Analytic methods in accident research: Methodological frontier and future
300 directions. Analytic methods in accident research 1, 1–22.
- 301 McCauley, P., Kalachev, L.V., Mollicone, D.J., Banks, S., Dinges, D.F., Van Dongen, H.P., 2013. Dynamic
302 circadian modulation in a biomathematical model for the effects of sleep and sleep loss on waking neurobehavioral

- 303 performance. *Sleep* 36, 1987–1997.
- 304 Meuleners, L., Fraser, M.L., Govorko, M.H., Stevenson, M.R., 2017. Determinants of the occupational environment
305 and heavy vehicle crashes in western australia: A case-control study. *Accident Analysis & Prevention* 99, 452–458.
- 306 Naik, B., Tung, L.-W., Zhao, S., Khattak, A.J., 2016. Weather impacts on single-vehicle truck crash injury
307 severity. *Journal of Safety Research* 58, 57–65.
- 308 National Highway Traffic Safety Administration, 2017. A Compilation of Motor Vehicle Crash Data from the
309 Fatality Analysis Reporting System and the General Estimates System.
- 310 Pantangi, S.S., Fountas, G., Sarwar, M.T., Anastasopoulos, P.C., Blatt, A., Majka, K., Pierowicz, J., Mohan,
311 S.B., 2019. A preliminary investigation of the effectiveness of high visibility enforcement programs using naturalistic
312 driving study data: A grouped random parameters approach. *Analytic Methods in Accident Research* 21, 1–12.
- 313 R Core Team, 2019. R: A language and environment for statistical computing. R Foundation for Statistical
314 Computing, Vienna, Austria.
- 315 Sharwood, L.N., Elkington, J., Meuleners, L., Ivers, R., Boufous, S., Stevenson, M., 2013. Use of caffeinated
316 substances and risk of crashes in long distance drivers of commercial vehicles: Case-control study. *BMJ* 346, f1140.
- 317 Stern, H.S., Blower, D., Cohen, M.L., Czeisler, C.A., Dinges, D.F., Greenhouse, J.B., Guo, F., Hanowski, R.J.,
318 Hartenbaum, N.P., Krueger, G.P., others, 2019. Data and methods for studying commercial motor vehicle driver
319 fatigue, highway safety and long-term driver health. *Accident Analysis & Prevention* 126, 37–42.
- 320 The Dark Sky Company, LLC, 2019. Dark Sky API — Overview.
- 321 Uddin, M., Huynh, N., 2017. Truck-involved crashes injury severity analysis for different lighting conditions on
322 rural and urban roadways. *Accident Analysis & Prevention* 108, 44–55.
- 323 Wali, B., Khattak, A.J., Karnowski, T., 2019. Exploring microscopic driving volatility in naturalistic driving
324 environment prior to involvement in safety critical events—concept of event-based driving volatility. *Accident
325 Analysis & Prevention* 132, 105277.
- 326 WHO, 2018. The top 10 causes of death.
- 327 Zhu, X., Srinivasan, S., 2011. A comprehensive analysis of factors influencing the injury severity of large-truck
328 crashes. *Accident Analysis & Prevention* 43, 49–57.