

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

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

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

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

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

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

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

Abstract

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

12 1. Introduction

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

Traditional crash prediction studies collect retrospective reports of crashes in a given road section for a specified time period, match these crash cases with non-crash controls (typically 1 to 4 matching), and then build statistical models, such as logistic regression and neural networks, to study risk factors associated with higher risk of crashes (Blower et al., 2010; Meuleners et al., 2017; Sharwood et al., 2013). This case-control study design is efficient and less time-consuming in trucking safety field since crashes are very rare. However, case-control studies, by nature, are limited in study design. Firstly, it is impossible to estimate and compare the rate of crashes since the number of

*Corresponding Author

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)

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
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**

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, trucks in Company A were equipped
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, in this study, we selected
62 496 regional truck drivers who move freights in a region and surrounding states. 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, instead, a unique identification number was provided for each driver to
65 link the three data sources. The study protocol was reviewed and approved by the Institutional Review Board of
66 Saint 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 identity number, and speed at that second. In total, 13,187,289
71 rows of ping data were generated by the 496 truck drivers.

72 Apart from ping data, Company A also collected real-time SCEs data for all their trucks. In contrast to
73 irregularly collected ping data, SCEs were recorded whenever pre-determined kinematic thresholds were triggered.
74 There were 12,458 critical events occurred to these 496 truck drivers during the study period. Four types of critical
75 events were recorded in this critical events data. The number of SCEs.

76 *2.2. Weather*

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

82 *2.3. Other available sources*

83 Traffic and road geometry can be collected from Google map API and OpenStreet 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.

84 **3. Data preparation**

85 *3.1. Data aggregation*

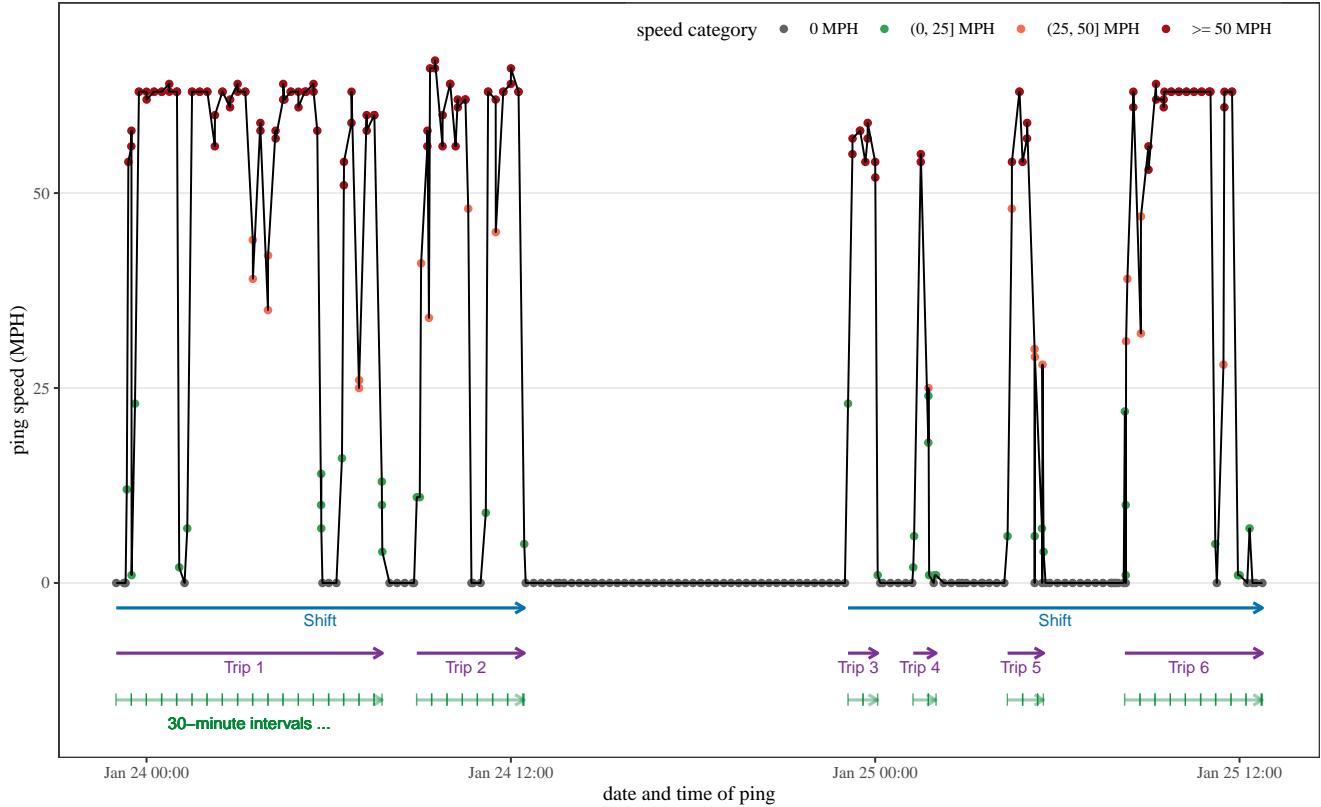


Figure 1: Data aggregation process from pings to shifts, trips, and 30-minute intervals.

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

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

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

98 3.2. Cumulative driving time as a measure of fatigue

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

101 Driver's fatigue is difficult to measure in real life. In this study, we attempt to use three proxies to measure the
102 fatigue of the truck drivers: cumulative driving time in a shift, the rest time before a shift, and the rest time before
103 a trip.

104 4. Methodology

105 4.1. Statistical models

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

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

Here we model the probability of a critical event occurred using a Bayesian hierarchical Bernoulli regression. We categorized the number of safety events during a trip into a binary variable Y 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. Since each trip i has a different travel time t_i , we derived the Bernoulli distribution parameter p_i using the probability density function of the Poisson distribution, with the parameter λ_i equaled a linear combination of β_i and x_i .

$$\begin{aligned} P_i &= P(\text{at least one event in trip } i) \\ &= 1 - P(\text{no event in trip } i) \\ &= 1 - \frac{e^{-t_i\lambda_i}(t_i\lambda_i)^0}{0!} \\ &= 1 - \exp(-t_i\lambda_i) \\ &= 1 - \exp(-t_i e^{\beta_0 + \beta_i x_i}) \end{aligned} \tag{1}$$

Transform that into a linear function of β_i , x_i and t_i

$$\begin{aligned} 1 - P_i &= \text{EXP}(-t_i e^{\beta_0 + \beta_i x_i}) \\ \log(1 - P_i) &= -t_i e^{\beta_0 + \beta_i x_i} \\ \log \frac{1}{1 - P_i} &= e^{\beta_0 + \beta_i x_i + \log(t_i)} \\ \log \left(\log \frac{1}{1 - P_i} \right) &= \beta_0 + \beta_i x_i + \log(t_i) \end{aligned} \tag{2}$$

Then, the random effects logistic model is

$$\begin{aligned} Y_i &\sim \text{Bern}(P_i) \\ \log \left(\log \frac{1}{1 - P_i} \right) &= \beta_{0,d(i)} + \beta_{1,d(i)} \cdot \text{CT}_i + \xi \cdot \mathbf{W} + \nu \cdot \mathbf{D}_i + \log(t_i) \end{aligned} \tag{3}$$

118 Here the trip is indexed by i , Y_i is the binary outcome variable of whether at least one critical event occurred in
 119 trip i ; $d(i)$ is the driver for trip i , $\beta_{0,d(i)}$ is the random intercept for driver $d(i)$; $\beta_{1,d(i)}$ is the random slope for the
 120 cumulative time (CT_i) of driving in the shift (the sum of driving time for all previous trips) for driver $d(i)$; \mathbf{W} is a
 121 vector of external environment fixed effects, including precipitation intensity and probability, visibility, and whether
 122 it was sunrise or sunset time; \mathbf{D}_i are driver level fixed effects, including age group and business unit; t_i is the travel
 123 time for the trip i .

We assume that the drivers are random effects, and we assume exchangeable priors of the form

$$\beta_{0,d(1)}, \beta_{0,d(2)}, \dots, \beta_{0,d(n)} \sim \text{i.i.d.}N(\mu_0, \sigma_0^2)$$

and

$$\beta_{1,d(1)}, \beta_{1,d(2)}, \dots, \beta_{1,d(n)} \sim \text{i.i.d.}N(\mu_1, \sigma_1^2)$$

The parameters μ_0 , σ_0 , μ_1 , and σ_1 are hyperparameters with priors. Since we do not have much prior knowledge on the hyperparameters, we assigned diffuse priors for these hyperparameters.

$$\begin{aligned} \mu_0 &\sim N(0, 10^2) \\ \mu_1 &\sim N(0, 10^2) \\ \sigma_0 &\sim \text{GAMMA}(1, 1) \\ \sigma_1 &\sim \text{GAMMA}(1, 1) \end{aligned} \tag{4}$$

124 Since μ_0 and μ_1 can be any real number, so we assigned two normal distributions with mean 0 and standard
 125 deviation of 10 as the priors for these two hyperparameters. In comparison, σ_0 and σ_1 must be strictly positive, so

¹²⁶ we assigned GAMMA(1, 1) with wide distribution on positive real numbers as their priors.

¹²⁷ 5. Results

¹²⁸ 5.1. Sample description

¹²⁹ 5.2. Statistical models

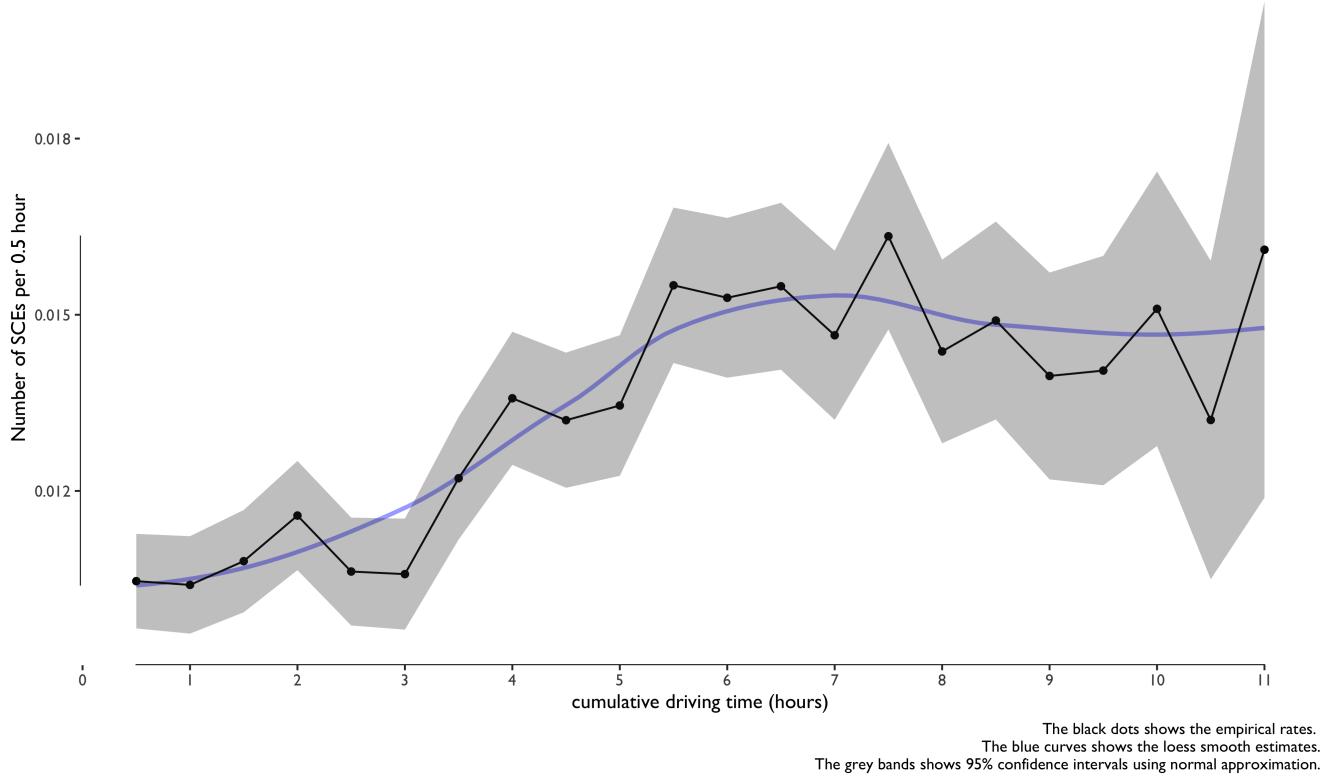


Figure 3: The rate of safety critical events and cumulative driving time

¹³⁰ 6. Discussion

¹³¹ 7. Conclusions

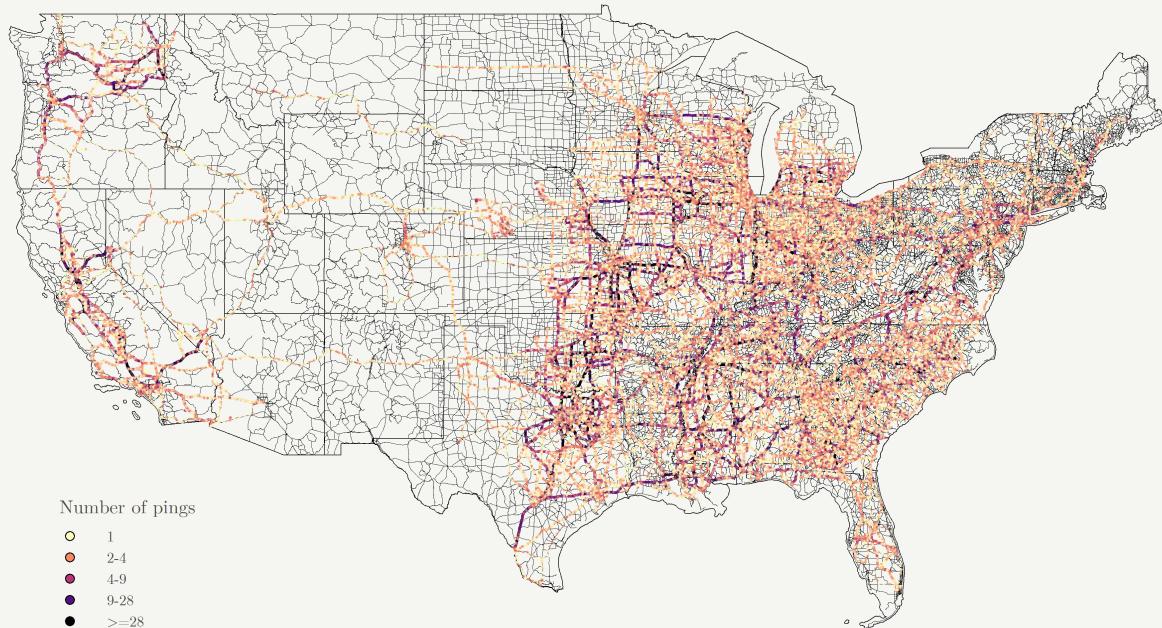
¹³² Subsampling MCMC (Dang et al., 2019; Quiroz et al., 2019, 2018, 2016).

¹³³ Acknowledgement

¹³⁴ This work was supported in part by the National Science Foundation (CMMI-1635927 and CMMI-1634992), the
¹³⁵ Ohio Supercomputer Center (PMIU0138 and PMIU0162), the American Society of Safety Professionals (ASSP)
¹³⁶ Foundation, the University of Cincinnati Education and Research Center Pilot Research Project Training Program,
¹³⁷ and the Transportation Informatics Tier I University Transportation Center (TransInfo). We also thank the DarkSky
¹³⁸ company for providing us five million free calls to their historic weather API.

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

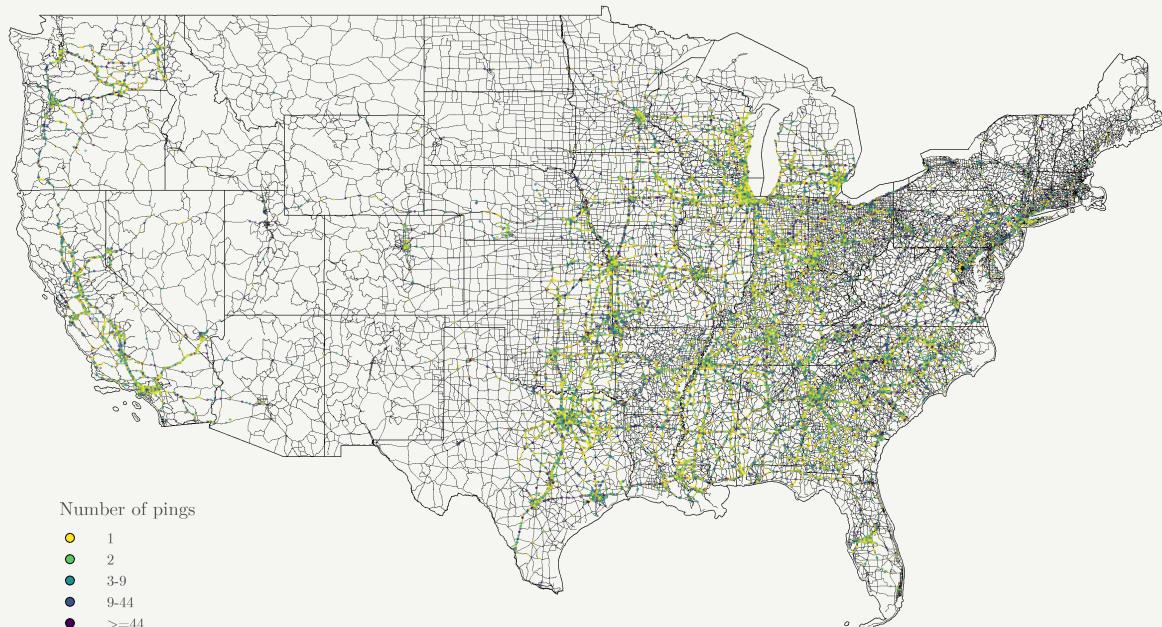
The drivers were employees in large commercial truck company in the United States



(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



(b) Inactive pings

Figure 2: Geographical point patterns of moving and stopped pings generated by the 496 sample drivers.

139 **References**

- 140 Blower, D., Green, P.E., Matteson, A., 2010. Condition of trucks and truck crash involvement: Evidence from
141 the large truck crash causation study. *Transportation Research Record* 2194, 21–28.
- 142 Cantor, D.E., Corsi, T.M., Grimm, C.M., Özpolat, K., 2010. A driver focused truck crash prediction model.
143 *Transportation Research Part E: Logistics and Transportation Review* 46, 683–692.
- 144 Chen, C., Zhang, G., Tian, Z., Bogus, S.M., Yang, Y., 2015. Hierarchical bayesian random intercept model-based
145 cross-level interaction decomposition for truck driver injury severity investigations. *Accident Analysis & Prevention*
146 85, 186–198.
- 147 Dang, K.-D., Quiroz, M., Kohn, R., Tran, M.-N., Villani, M., 2019. Hamiltonian Monte Carlo with energy
148 conserving subsampling. *Journal of Machine Learning Research* 20, 1–31.
- 149 Dingus, T.A., Hanowski, R.J., Klauer, S.G., 2011. Estimating crash risk. *Ergonomics in Design* 19, 8–12.
- 150 Dong, C., Dong, Q., Huang, B., Hu, W., Nambisan, S.S., 2017. Estimating factors contributing to frequency and
151 severity of large truck-involved crashes. *Journal of Transportation Engineering, Part A: Systems* 143, 04017032.
- 152 Guo, F., 2019. Statistical methods for naturalistic driving studies. *Annual Review of Statistics and Its Application*
153 6, 309–328.
- 154 Guo, F., Klauer, S.G., Hankey, J.M., Dingus, T.A., 2010. Near crashes as crash surrogate for naturalistic driving
155 studies. *Transportation Research Record* 2147, 66–74.
- 156 Han, C., Huang, H., Lee, J., Wang, J., 2018. Investigating varying effect of road-level factors on crash frequency
157 across regions: A bayesian hierarchical random parameter modeling approach. *Analytic methods in accident research*
158 20, 81–91.
- 159 Hickman, J.S., Hanowski, R.J., Bocanegra, J., 2018. A synthetic approach to compare the large truck crash
160 causation study and naturalistic driving data. *Accident Analysis & Prevention* 112, 11–14.
- 161 Mannering, F.L., Bhat, C.R., 2014. Analytic methods in accident research: Methodological frontier and future
162 directions. *Analytic methods in accident research* 1, 1–22.
- 163 Meuleners, L., Fraser, M.L., Govorko, M.H., Stevenson, M.R., 2017. Determinants of the occupational environment
164 and heavy vehicle crashes in western australia: A case–control study. *Accident Analysis & Prevention* 99, 452–458.
- 165 Naik, B., Tung, L.-W., Zhao, S., Khattak, A.J., 2016. Weather impacts on single-vehicle truck crash injury
166 severity. *Journal of Safety Research* 58, 57–65.
- 167 National Highway Traffic Safety Administration, 2017. A Compilation of Motor Vehicle Crash Data from the
168 Fatality Analysis Reporting System and the General Estimates System.
- 169 Pantangi, S.S., Fountas, G., Sarwar, M.T., Anastasopoulos, P.C., Blatt, A., Majka, K., Pierowicz, J., Mohan,
170 S.B., 2019. A preliminary investigation of the effectiveness of high visibility enforcement programs using naturalistic
171 driving study data: A grouped random parameters approach. *Analytic Methods in Accident Research* 21, 1–12.

- 172 Quiroz, M., Kohn, R., Villani, M., Tran, M.-N., 2019. Speeding up MCMC by efficient data subsampling. *Journal*
173 of the American Statistical Association
- 174 Quiroz, M., Tran, M.-N., Villani, M., Kohn, R., 2018. Speeding up MCMC by delayed acceptance and data
175 subsampling. *Journal of Computational and Graphical Statistics* 27, 12–22.
- 176 Quiroz, M., Tran, M.-N., Villani, M., Kohn, R., Dang, K.-D., 2016. The block-Poisson estimator for optimally
177 tuned exact subsampling MCMC. arXiv preprint arXiv:1603.08232.
- 178 Sharwood, L.N., Elkington, J., Meuleners, L., Ivers, R., Boufous, S., Stevenson, M., 2013. Use of caffeinated
179 substances and risk of crashes in long distance drivers of commercial vehicles: Case-control study. *BMJ* 346, f1140.
- 180 Stern, H.S., Blower, D., Cohen, M.L., Czeisler, C.A., Dinges, D.F., Greenhouse, J.B., Guo, F., Hanowski, R.J.,
181 Hartenbaum, N.P., Krueger, G.P., others, 2019. Data and methods for studying commercial motor vehicle driver
182 fatigue, highway safety and long-term driver health. *Accident Analysis & Prevention* 126, 37–42.
- 183 The Dark Sky Company, LLC, 2019. Dark Sky API — Overview.
- 184 Uddin, M., Huynh, N., 2017. Truck-involved crashes injury severity analysis for different lighting conditions on
185 rural and urban roadways. *Accident Analysis & Prevention* 108, 44–55.
- 186 WHO, 2018. The top 10 causes of death.
- 187 Zhu, X., Srinivasan, S., 2011. A comprehensive analysis of factors influencing the injury severity of large-truck
188 crashes. *Accident Analysis & Prevention* 43, 49–57.