

Pre-Print Manuscript:

Zhou, Y., and Bridgelall, R. (2020). A Review of Usage of Large-Scale Connected Vehicle Data. *Transportation Research Record*, DOI: 10.1177/0361198120940996.

1 **A Review of Usage of Large-Scale Connected Vehicle Data**

2

3 **Yun Zhou***

4 Doctoral Graduate Research Assistant

5 Department of Transportation, Logistics, and Finance, North Dakota State University

6 NDSU Department 2880

7 P.O. Box 6050, Fargo, ND 58108

8 Phone: (218) 289-7182; Email: yun.zhou@ndsu.edu

9 ORCID: 0000-0002-2782-0282

10

11 **Raj Bridgelall, Ph.D.**

12 Assistant Professor and Program Director

13 Department of Transportation, Logistics and Finance, North Dakota State University

14 NDSU Department 2880 P.O. Box 6050, Fargo, ND 58101-6050

15 Phone: (408) 607-3214; Email: raj@bridgelall.com

16 ORCID: 0000-0003-3743-6652

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

*Corresponding Author

35

Word count: 5,331 + 1,750 (7 Tables 250 each) = 7081 (7500 max)

36

Abstract: 208 words (250 max)

37

38

Funding Acknowledgements

39

The authors received no financial support for the research, authorship, and/or publication of this article.

40

41

42

Data Accessibility Statement

43

This paper (the *SPMD Dataset* subsection of the *Methodology* section) reviewed data that is

44

available on the USDOT public data portal (<https://data.transportation.gov/> or

45

<https://www.its.dot.gov/data/>) and articles available as referenced.

1 ABSTRACT

2 GPS loggers and cameras aboard connected vehicles can produce vast amounts of data. Analysts
3 can mine such data to decipher patterns in vehicle trajectories and driver-vehicle interactions. An
4 ability to process such large-scale data in real time can inform strategies to reduce crashes,
5 improve traffic flow, enhance system operational efficiencies, and reduce environmental
6 impacts. However, connected vehicle technologies are in the very early phases of deployment.
7 Hence, related datasets are extremely scarce, and the utility of such emerging datasets is largely
8 unknown. Subsequently, this paper provides a comprehensive review of studies that used large-
9 scale connected vehicle data from the United States Department of Transportation Connected
10 Vehicle Safety Pilot Model Deployment program. It is the first and only dataset available to the
11 public. The data contains real-world information about the operation of connected vehicles that
12 organizations are testing. The authors provide a summary of the available datasets, their
13 organization, the overall structure, and other characteristics of the data captured during pilot
14 deployments. Subsequently, the authors classify the data usage into three categories: driving
15 pattern identification, development of surrogate safety measures, and improvements in the
16 operation of signalized intersections. Finally, the authors identify some limitations experienced
17 with the existing dataset.

18

19 *Keywords:* Connected vehicle, Intelligent Transportation Systems, Smart cities, Signal phase and
20 timing, Surrogate safety measures, Risky driving pattern, Intersection safety.

21

22 INTRODUCTION

23 Connected vehicle (CV) technology enables real-time communications among users, vehicles,
24 and the multimodal infrastructure. Producers and transportation agencies assert that CV
25 technology will dramatically reduce the number of fatalities and serious injuries caused by
26 accidents on our roadways. The technology will achieve this by notifying and alerting drivers
27 about potentially dangerous driving situations. Examples include pedestrians or bicyclists
28 approaching an intersection, vehicles in blind areas beyond a curve, and oncoming cars swerving
29 into a lane to avoid an object or pothole on the road (1). CV technology can also smooth out
30 traffic flows, diminish congestion, and reduce travel time. Agencies can analyze CV data to
31 inform eco-friendly transportation planning (2).

32 GPS loggers and cameras aboard CVs can produce abundant travel data. CV data
33 exchanges include vehicle-to-vehicle (V2V) and vehicle-to-infrastructure (V2I). CVs
34 communicate basic safety messages (BSM) among vehicles and with roadside equipment (RSE)
35 (3-4). BSMs contain information such as vehicle location, speed, acceleration, and time. In
36 addition to BSMs, CVs also produce trajectory data, various driver-vehicle interaction data, and
37 contextual data. Analysts can mine such data to inform strategies that can reduce crashes,
38 improve traffic flow, enhance system operational efficiencies, and reduce environmental
39 impacts.

40 CV development is still in its infancy. Hence, there is very little information about
41 practical challenges and quantifiable benefits of real-world deployments. Therefore, the United
42 States Department of Transportation (USDOT) launched a one-year Connected Vehicle Safety
43 Pilot Model Deployment (SPMD) in August 2012 in Ann Arbor, Michigan to advance
44 knowledge about practical deployments. The deployment included 30 roadside equipment (RSE)
45 installations along approximately 75 lane-miles of roadway and approximately 3,000 equipped
46 vehicles (4). The data collected is now available to the public via the USDOT's public data

1 portal (<https://data.transportation.gov/> or <https://www.its.dot.gov/data/>) (4). As part of the
2 USDOT CV pilot programs, the agency awarded in 2015 a three-phase pilot to two cities and one
3 state—New York City, Tampa, and Wyoming (5). All three sites finished a 12-month period of
4 concept development (phase I) and a 24-month period of deployment design, build, and test
5 (Phase II). Each site is now entering the final deployment phase to operate and test deployed CV
6 systems for a minimum of 18 months. The sites will monitor key performance measures of the
7 CV pilot. Some datasets from this program are available at the USDOT’s public data portal.

8 Based on the current literature, analysts can use the CV data by itself or combine it with
9 other data sources to identify driving patterns, develop surrogate safety measures, evaluate
10 location-based intersection safety, and improve operations at signalized intersections. The
11 **contributions** of this paper are a comprehensive snapshot-in-time review of previous studies that
12 used connected vehicle data, and potential areas that researchers can advance using the newly
13 released large-scale CV data from the USDOT Pilot Deployment Program. This paper contains
14 all the information in one place to facilitate ongoing research about the potential value and utility
15 of emerging CV data.

16 The remainder of this paper includes: a methodology section which describes the
17 approach to the literature review and the strategy to classify the usage of the CV data; a results
18 section which summarizes the data usage and applications of the CV dataset, a conclusion
19 section which provides overall remarks about the usability and limitations of the CV dataset and
20 briefly discusses future work.

21 22 **METHODOLOGY**

23 The authors first researched the newly released large-scale CV datasets from the United States
24 Department of Transportation (USDOT) Pilot Deployment Program to understand its
25 organization, structure, and content. After conducting extensive literature searches using all the
26 traditional scientific databases of research output, the authors organized the data usage into three
27 categories of application development. The publication sources included the Transportation
28 Research Information Services (TRIS) Database and the International Transport Research
29 Documentation (ITRD) Databases.

30 31 **SPMD Dataset**

32 The SPMD data comprises of the following datasets: two driving datasets that include the Data
33 Acquisition System 1 (DAS1) and Data Acquisition System 2 (DAS2), a BSM dataset, an RSE
34 dataset, and three contextual dataset that include weather, network, and schedule. Each table in a
35 database contains comma-separated value (CSV) formatted information collected during a 24-
36 hour period (6).

37 DAS1 and DAS2 contain data from the DAS that the University of Michigan
38 Transportation Research Institute (UMTRI) and the Virginia Tech Transportation Institute
39 (VTTI) developed, respectively. The BSM dataset includes messages that a participating vehicle
40 transmitted and/or received, irrespective of the DAS installed. The RSE dataset contains data that
41 roadside units transmitted and/or received. The contextual datasets contain information about
42 conditions at the time of data collection. Contextual data include information about network
43 configuration and performance, weather, schedules (transit and special events), roadwork
44 activity, and traffic incidents.

45 Several studies indicated that the CV data was available at the Research data Exchange
46 (RDE, <https://www.its-rde.net/home>), but this link is no longer active. The CV data is now only

1 available at the USDOT's public data portal (<https://data.transportation.gov/> or
2 <https://www.its.dot.gov/data/>). The *DataWsu* and *DataFrontTargets* are considered as part of the
3 DAS1 dataset but they available as separate data files at the USDOT's public data portal. The
4 SPMD Sample Data Handbook provides a detailed introduction for each dataset and all the data
5 files under that dataset (6). Table 1 summarizes the datasets and a list of their accompanying files
6 (6). The DAS1 dataset was mentioned most as a data source in the reviewed studies. Table 2
7 further describes each file in the DAS1 dataset. The *DataWsu* and *DataFrontTargets* files are
8 most commonly used to capture the position and motion information of host vehicles. The
9 *DataWsu* file contains 27 fields, which is the most of all the DAS1 data files. Most of the data
10 logged in the *DataWsu* file comes from the onboard Wireless Safety Unit (WSU) that produces
11 GPS and inertial sensor data, and the Controller Area Network (CAN) that communicates vehicle
12 performance and status information. Table 3 describes the data elements in the *DataWsu* file.

13 The Mobileye system the Intel Corporation is a vision-based system that enables various
14 Advanced Driver Assistance System (ADAS) capabilities. The *DataFrontTargets* file contains
15 information from the installed Mobileye system that collects information from the scene ahead of
16 the vehicle. The system uses communicates measures and warnings based on a serious of
17 proprietary algorithms. Table 4 briefly describes the data elements of *DataFrontTargets* file.
18

19 **CV Pilot Deployment Program Dataset**

20 The program scale of the three-phase pilot sites is much larger than that of the SPMD program.
21 Furthermore, the estimated program duration of 54 months is more than four times that of the
22 SPMD program. At the time of this writing, there were only six sample RSE data files available
23 from the USDOT's public data portal. They are:

- 24 1. Wyoming
 - 25 a. BSM one-day sample file
 - 26 b. Two traveler information message (TIM) sample files
- 27 2. Tampa
 - 28 a. One signal phasing and timing (SPaT) sample file
 - 29 b. One BSM sample file
 - 30 c. One TIM sample file
- 31 3. New York
 - 32 a. None.

33 Because the data from the three CV Pilot Deployment sites are limited, this paper focuses
34 on studies that use the SPMD data or other CV data sources or probe data to determine the
35 research areas that can be advanced using emerging CV datasets from the three CV Pilot
36 Deployment sites.
37

1 **Table 1 Files Associated with the SPMD Dataset. Source: SPMD Sample Data Handbook (6).**

Driving Data		Message	Infrastructure	Contextual		
DAS1	DAS2	BSM	RSE	Weather	Network	Schedule
AudioTimes*	HV_Radar	BrakeByte1Events	BSM	Weather/ climatic data	Pointer to Resources	Pointer to* Resources
DataFrontTargets	HV_Primary	BrakeByte2Events	Geometry			
DataLane	DAS2_Trip_Summary*	BsmP1	Lane			
DataWsu		ExteriorLightsEvents	LaneNode			
DAS1_Trip_Summary*		PosAccurByte1Events	MAP			
		PosAccurByte2Events	Packet			
		PosAccurByte3Events	PCAPFile			
		PosAccurByte4Events	SPAT			
		SteerAngleEvents	SPATMovement			
		ThrottlePositionEvents	TIM			
		TransStateEvents	TIMRegion			
		WiperStatusFrontEvents	TIMRegionNode			
		BSM_Trip_Summary*				

2
3 *Not available at USDOT's public data portal (<https://data.transportation.gov/>).

4
5
6 **Table 2 Description of the Files in the DAS1 Dataset. Source: SPMD Sample Data Handbook (6).**

File Number	File	Description	Update Frequency
1	DataFrontTargets	Data collected by the Mobileye sensor. It captures information about the (vehicle) or object that is in front of the host vehicle.	10Hz
2	DataLane	Stores lane marking quality adjacent to the host vehicle and the distances between each side of the vehicle and each lane line.	10Hz
3	DataWsu	Logs of the GPS and CAN Bus data produced by an onboard device.	10Hz
4	DAS1_Trip_Summary*	A list of summary measures for each vehicle trip.	1 per trip

7
8 *Not available at USDOT's public data portal (<https://data.transportation.gov/>).

1 **Table 3 Data Elements of the DataWsu File. Source: SPMD Sample Data Handbook (6).**

Field Name	Type	Units	EnumId	Description
Device	Integer	none	-	A unique numeric ID assigned to each DAS. This ID also doubles as a vehicle's ID.
Trip	Integer	none	-	Count of ignition cycles—each cycle starts and ends when the ignition is in the on and off positions, respectively.
Time	Integer	centiseconds	-	Time (centiseconds) since the DAS started, which (generally) starts when the ignition is in the on position.
GpsValidWsu	Integer	none	1	Indicates whether or not a GPS data point is valid.
GpsTimeWsu	Integer	millisecond		Epoch GPS time received from the remote vehicle that has been targeted by the host vehicle's WSU.
LatitudeWsu	Float	deg	-	Latitude from WSU receiver .
LongitudeWsu	Float	deg	-	Longitude from WSU receiver.
AltitudeWsu	Real	m	-	Altitude from WSU receiver.
GpsHeadingWsu	Real	deg	-	Heading from WSU GPS receiver .
GpsSpeedWsu	Real	m/sec	-	Speed from WSU GPS receiver.
HdopWsu	Real	none	-	Horizontal dilution of precision.
PdopWsu	Real	none	-	Position dilution of precision.
FixQualityWsu	Integer	none	-	GPS Fix Quality.
GpsCoastingWsu	Integer	none	-	GPS Coasted.
ValidCanWsu	Integer	none	1	Valid Vehicle CAN Bus message to WSU.
YawRateWsu	Real	deg/sec	-	Yaw rate from vehicle CAN Bus via WSU.
SpeedWsu	Real	kph	-	Speed from vehicle CAN Bus via WSU.
TurnSngRWsu	Integer	none	11	Right turn signal from vehicle CAN Bus via WSU.
TurnSngLWsu	Integer	none	11	Left turn signal from vehicle CAN Bus via WSU.
BrakeAbsTcsWsu	Integer	none	-	Brake, ABS, and traction control from vehicle CAN Bus via WSU.
AxWsu	Real	m/sec ²	-	Longitudinal acceleration from vehicle CAN Bus via WSU.
PrndlWsu	Integer	none	403	Current transmission state (Park, Reverse, Neutral, Drive, Low) from vehicle CAN Bus via WSU.
VsaActiveWsu	Integer	none	-	Stability control active from vehicle CAN Bus via WSU.
HeadlampWsu	Integer	none	-	Headlamp state from vehicle CAN Bus via WSU.
WiperWsu	Integer	none	-	Wiper state from vehicle CAN Bus via WSU.
ThrottleWsu	Real	none	-	Throttle position from vehicle CAN Bus via WSU.
SteerWsu	Real	deg	-	Steering angle/position from vehicle CAN Bus via WSU.

1
2 **Table 4 Data Elements of the DataFrontTargets File. Source: SPMD Sample Data Handbook (6).**

Field Name	Type	Units	EnumId	Description
Device	Integer	none	-	A unique numeric ID for each DAS, which also doubles as a vehicle's ID.
Trip	Integer	none	-	Count of ignition cycles that begins and ends when the ignition is in the on and off position, respectively.
Time	Integer	centiseconds	-	Time (centiseconds) since DAS started, which (generally) starts when the ignition is in the on position.
TargetId	Integer	none	-	Numeric ID that the Mobileye sensor assigns to different objects being tracked, with a value of 1 assigned to the closest.
ObstacleId	Integer	none	-	ID of new obstacle that the Mobileye sensor assigns—the value is the last used free ID.
Range	Integer	m	-	Longitudinal position of an object (typically the closest) relative to a Mobileye defined reference point on the host vehicle.
RangeRate	Real	m/sec	-	Rate of change of the Range variable.
Transversal	Real	m	-	Mobileye assigned lateral position of the obstacle.
TargetType	Integer	none	409	Object classification (car, truck, pedestrian, etc.)
Status	Integer	none	410	Motion classification (stopped, moving, etc) of an identified obstacle/target.
CIPV	Integer	none	1	Indicates whether an obstacle in the vehicle's path is the closest.

3

1 RESULTS

2 Several studies reported on the use of CV data from SPMD. This section classifies them into
3 three categories: driving pattern identification, development of surrogate safety measures, and
4 improvement of signalized intersection operation.

6 Driving Pattern Identification

7 Driving pattern is one of the key factors that affect traffic safety. Vehicle speed, acceleration, and
8 deceleration are primary factors in the classification of driving patterns. Agencies consider that
9 driving above the speed limit is hazardous and risky. Speed restrictions can be a dynamic
10 function of road conditions and traffic situations. Researchers proposed several acceleration
11 thresholds to classify driving behavior as calm, normal, and aggressive (7-11). Risky driving
12 happens when longitudinal or lateral accelerations exceed certain thresholds. Researchers found
13 that risky driving patterns are highly correlated with the likelihood of crashes or near-crash
14 events (8-9) (11-13).

15 Three recent studies used the CV data from the SPMD program to propose
16 methodologies that use critical information from the instantaneous BSM exchanges between CVs
17 and roadside equipment to determine repeatable driving behaviors (14-16). Study (14)
18 investigated the longitudinal and lateral motion of the driving decision from BSMs and
19 established reasonable thresholds to identify potentially dangerous events such as hard
20 accelerations or braking, and quick lane changes. They used the DAS datasets and the 10-Hertz
21 motion data that contain speed along with longitudinal and lateral acceleration to visualize the
22 driving behavior. They investigated the relationships between speed and acceleration by
23 visualizing the distributions of acceleration in longitudinal and lateral directions. The results
24 validated that instantaneous driving decisions could provide valuable information to identify
25 extreme driving events such as sudden lane changes. The distributions of directional variations in
26 acceleration informed the thresholds of extreme acceleration in different directions. The study
27 presented valuable information on establishing context-relevant alerts, warnings, and control
28 assistance to nearby vehicles.

29 Study (15) conducted a time-series analysis to categorize driving patterns into different
30 regimes based on their volatility and average duration. The study explored correlations to
31 examine dwell times and switching times between regimes. The study used Google Earth to
32 visualize vehicle trajectories from the DAS dataset and to classify trips based on roadway type.
33 Aggregating the DAS datasets into one-second groups enabled a detailed econometric analysis of
34 instantaneous driving decisions. They analyzed the data using an expectation-maximization and
35 Dynamic Markov switching models of two-regimes and three-regimes. The results revealed that
36 acceleration and deceleration are two distinct regimes. The rate of deceleration was higher than
37 the rate of acceleration, and braking was more volatile during deceleration than during
38 acceleration.

39 Machine learning methods can identify the importance of motion-related variables in
40 classifying driving data into aggressive and normal driving patterns (17-18). The authors of (16)
41 applied machine learning to the CV data from the SPMD program to identify aggressive driving
42 patterns on horizontal curves. They used the random forest method of machine learning to
43 develop an aggressive driving detection model based on a time-to-lane crossing (TLC) metric,
44 under three scenarios. This detection model provided high classification accuracy in all three
45 scenarios, and it ranked the importance of the variable candidates in identifying aggressive
46 driving behaviors.

1 A common limitation for these three studies that used the SPMD dataset was the limited
2 number of observations. Study (14-15) used the one-day sample datasets but deleted some
3 observations that contained errors. Study (16) used a one-month dataset. Larger datasets require
4 more processing capacity to produce results in a reasonable amount of time, but they can
5 facilitate the removal of outlier situations, such as extremely bad weather conditions that could
6 bias the results.

7 **Surrogate Safety Measures Development**

8 Historical crash data can identify high-risk locations. However, because historical crash data
9 usually take a significant amount of time to collect, researchers developed surrogate safety
10 measures (SSMs). They are proactive solutions to assess safety risks by capturing near-crash
11 events when crash data is absent or limited. SSMs can quantify safety-related performance at a
12 road segment or evaluate the effectiveness of a safety treatment more efficiently (19). Data
13 collected from sensors can inform the development of SSMs to identify high-risk locations
14 accurately. The newly available CV data collected from CV devices and RSEs give researchers a
15 new opportunity to conduct SSM research. Some researchers attempted to develop SSMs from
16 the vehicle trajectory data of the SPMD program (20-21).

17 Study (20) proposed a framework to process CV data, calculate SSMs and their safety
18 indices, and analyze the correlation between crash records and the calculated safety indices.
19 They calculated three SSMs, which were time-to-collision (TTC), modified time-to-collision
20 (MTTC), and deceleration rate to avoid collision (DRAC), at the vehicle-level and their safety
21 indices at both the trip-level and the link-level. They used a negative binomial model to analyze
22 the relationship between crash records and the safety indices of three SSMs. They found that the
23 MTTC model provided the best overall performance. The study concluded that using SSMs with
24 motion-related CV data could improve overall safety evaluation.

25 Study (21) developed a new SSM called time-to-collision with disturbance (TTCD)
26 which can capture rear-end conflict risks in car-following scenarios. They used the CV data to
27 access the risk identified with the TTCD model by comparing that risk with the risk identified by
28 historical crash data. The result was that, among all accessed SSMs, TTCD can capture risk data
29 at the highest level of correlation with historical rear-end crash data. The associated high-risk
30 locations identified by TTCD were very similar to those identified by historical crash data. This
31 study suggested that researchers could use real-world CV data to identify high-risk locations as
32 crash predictors. This result validated the results from study (22), which presented a framework
33 and used simulated CV data to examine surrogate measures for evaluating the risk of secondary
34 crashes on highways.

35 Both studies (20-21) presented detailed procedures for cleaning and processing the CV
36 data from the SPMD program. They determined that CV data could help predict high-risk
37 locations in a very short period (one month and two months) before a substantial number of
38 crashes occur. Thus, CV data can help to develop proactive safety measures to improve road
39 safety management (21).

40 Intersections are one of the most dangerous locations on roadways based on their annual
41 crash history (23-24). Traditionally, the safety evaluation of an intersection depends on historical
42 crash frequency data and survey feedback from active users. Analysis of CV data can detect
43 high-risk intersections where historical crash data is limited. Three studies (25-27) combined
44 historical crash data and real-world CV data from the SPMD program to evaluate location-based
45 intersection safety. The purpose of this type of study is to conduct proactive safety measures on
46

1 specific intersections before the occurrence of crashes, and to seek solutions for improving their
2 safety. Studies (20-21) attempt to identify high-risk locations which are not limited to
3 intersections, whereas studies (25-27) focused on evaluating location-based intersection safety
4 due to the high crash frequency detected at intersections.

5 One study (25) provided an example of using CV data to assess location-based risk by
6 detecting extreme driving decisions. Researchers used the results of previous work that identified
7 extreme driving events to estimate crash risk as a function of instantaneous driving decisions at
8 specific locations. They introduced the concept of location-based volatility (LBV) to calculate
9 the coefficient of variation as a standardized measure of dispersion. The coefficient of variation
10 indicates the fluctuation of the longitudinal acceleration and deceleration at a specific location.
11 The study used rigorous fixed- and random-parameter Poisson regression models to investigate
12 the relationship between LBV and crash frequency. Results suggested that there is a statistically
13 significant relationship between LBV and crash frequencies for signalized intersections. One
14 limitation of the study was the limited number of variables available after the deletion of
15 inaccurate observations.

16 Another study (26) combined the CV BSM data with the crash and inventory data at
17 several intersections to investigate the relationship of different measures of volatility with crash
18 frequency. The study evaluated thirty-seven different measures of volatility. The researchers
19 used fixed- and random-parameter Poisson regression models in two levels of bulk passing and
20 individual passing to investigate the relationship between each measure of volatility and crash
21 frequency. Several methods investigated showed positive and statistically significant association
22 with crash frequency. They methods were the three measures at bulk passing level, the time-
23 varying stochastic volatility of speed, the percent laying beyond threshold-bonds of speed
24 created using mean plus two standard deviation at intersections, and the percent laying beyond
25 threshold-bonds of acceleration created using mean plus two standard deviation at intersections.

26 One research group (27) proposed a methodology to quantify driving volatility at each
27 intersection to assess intersection-based crash risk based on CV BSM data, crash data, and traffic
28 and intersection inventory data. They proposed to quantify driving volatility based on speed,
29 acceleration/deceleration, vehicular jerk, eight different volatility measures, coefficient of
30 variation, mean absolute deviation around mean, percentage of outliers, and time-dependent
31 dynamic volatility, at both aggregate intersection level and trip level. Subsequently, they used
32 Poisson and Poisson-lognormal regressions models to test the correlations between crash
33 frequency and intersection-based volatility with consideration of unobserved heterogeneity. They
34 used Full Bayesian estimation method and Markov Chain Monte-Carlo Gibbs Sampler
35 techniques to estimate the parameters. They calculated Moran's I statistics to investigate the
36 correlation between crash frequency and spatial factors. The results suggested that the crash
37 frequency significantly correlated with three measures: the two standard deviations threshold at
38 the intersection level, the coefficient of variation of speed at the pass level, and the mean
39 absolute deviance of vehicular jerk at the passing level.

40 The common limitation of the three studies was the limited CV data available during the
41 research period. The one-month or two-month CV data were relatively small sample sizes to
42 explain 5-year average crash rates. In addition, all studies only considered crash frequency but
43 not crash severity as a risk factor.

44

1 **Signalized Intersection Operation Improvement**

2 Signalized intersections are usually hot spots for traffic congestion, especially during rush hour.
3 They cause significant hours of delay and crash volume every year. Agencies often consider
4 adaptive signal control to accommodate varying demands. However, their significant cost to
5 install and maintain is a deterrent to deployment. A less-expensive alternative is to re-time the
6 signal by analyzing CV data from RSEs to estimated traffic volume (28-30). A downside of this
7 approach is the currently low penetration rate of CVs. Study (31) conducted a proof-of-concept
8 by using CV data in a low penetration rate environment to optimize signal coordination. They
9 could not use vehicle trajectories because the data was from a fixed location. Estimating traffic
10 volume from vehicle trajectories is an essential input for signal operation and algorithm design
11 optimization.

12 Study (32) proposed a method to estimate traffic volumes by using trajectory data from
13 CVs or trajectory data from navigation devices at locations with low CV penetration rates. The
14 study used the BSM data from the RSE in the SPMD program to capture trajectory variables
15 such as motion and position, and data from the signal phase and timing (SPaT) data to capture
16 the timing periods and signal status. They combined the data to produce a space-time trajectory
17 dataset. They modeled traffic arrivals as a time-dependent Poisson process. An expectation
18 maximization (EM) procedure provided an estimate of the arrival rate. They tested the estimation
19 procedure with CV trajectory data and vehicle trajectory data from navigation service users. The
20 results suggested that their proposed approach was effective and that agencies could use it to
21 improve signal control operation in environments with low CV penetration rates.

22 Monitoring queue status at signalized intersections could help optimize the available
23 capacity (33-35). For example, CV data can enable improved route selection through real-time
24 notifications of traffic status. Study (36) proposed an integrated macroscopic and microscopic
25 traffic flow model to estimate time-space queueing dynamics at signalized intersections using
26 RSE and SPaT data from the CV BSM repository. They identified the three regions of queue
27 formation region, queue region, and queue dissipation under the normal scenario and
28 oversaturated scenario based on vehicle deceleration, stop, and acceleration behaviors. This
29 integrated method estimated queue process in both queue length and queue time at signalized
30 intersections. Thus, this method could help to improve real-time traffic status estimation at the
31 signalized intersections equipped with connected vehicle technologies.

32 **Data Manipulation**

33 ***Error Checking***

34
35 Best practices in the use of large-scale data begin with quality evaluation and cleaning before
36 conducting any analysis or data mining tasks. However, relatively few studies reported
37 experience with error-checking and deep-cleaning SPMD datasets. The U.S. National Highway
38 Traffic Safety Administration published an independent evaluation of safety applications for
39 passenger vehicles in the SPMD program (37). They found several errors in the programming of
40 Volkswagen-Audi's forward-collision warning (FCW), intersection movement assist (IMA)
41 applications, and issues with a GPS on one vehicle that led to inaccuracies in some data records
42 in the SPMD datasets. Other researchers found errors in the DAS dataset such as speeds faster
43 than 200 mph and altitudes greater than 30,000 ft (14) (21). Another research found that a
44 significant portion (42%) of the "lateral acceleration" observations exceeded the maximum value
45

1 that the wireless communications device could record (25). Duplicated records and invalid
2 messages were found in DAS dataset (21).

3 Studies (15, 26) checked dataset and claimed that no error has occurred. Study (15)
4 indicated that they conducted error-checking process by linking microscopic trip data with a trip-
5 summary file to check the information consistency at trip-level and didn't state error-checking
6 for data values. Study (26) performed an error-checking process for its spatial data by mapping
7 data and resulted with a good match with the real map.

8 Among all the studies, study (21) presented a detailed data process procedure by a
9 detailed data preparation description and a data flow chart, which provides great guidelines for
10 both data cleaning and data process for later researches. Study (20) presented data process and
11 indicated the software applied for data process. Study (26) presented data integration and process
12 steps in a data flow chart. The authors of (38) demonstrated an automated method of cleaning the
13 data of a taxi probe dataset that utilizes known distributions of vehicle operations to detect
14 possible outliers for removal. Table 5 summarizes data shortcomings and provides
15 recommendations for data cleaning based on knowledge synthesized from the literature search.
16 Table 6 summarizes the data cleaning and processes, and the software tools used.
17

18 ***Data Mining Approaches***

19 Data mining techniques extract patterns from large-scale data that are interesting (39). Common
20 data mining approaches include statistical regression models and machine learning methods (39).
21 Statistical regression models estimate the numerical relationships between variables and can
22 predict new values, whereas machine learning methods recognize complex patterns and facilitate
23 decision-making based on data (39).

24 Statistical regression models are the most commonly used methods among all the studies
25 reviewed. As shown in Table 7, seven studies (14-15, 20, 25-27, 32) applied different regression
26 models to investigate the numerical relationships between factors. Study (21) calculated
27 Pearson's correlation coefficients as test statistics to quantify correlations between developed
28 measures and crash frequency. One study used the machine learning method of Random Forest
29 to identify aggressive/risky driving (16). Additional studies suggested potential application of
30 machine learning methods in CV data. Random forest and Support Vector Machines have been
31 applied to the driving style classification and transportation mode recognition problem (39-40).
32 Study (33) used CV data, without machine learning, to demonstrate the developed model.
33

34 ***Data Mining Challenges***

35 As summarized in Table 5, a common data cleaning recommendation is to detect and remove
36 erroneous records from the dataset, and to abandon data fields where there is a considerable
37 number of data records with errors. However, data elimination may reduce the data size, which
38 could result in lower model accuracy (15). Computation capacity may become another issue
39 when the dataset is very large. Studies (27) indicated that the workstation level computer (Dell
40 Precision T7600, 3.1 GHZ (32 CPUs) took a considerable amount of time to compute data
41 models for a data size of 230 million observations.

1 **Table 5 Data Shortcomings and Recommendation for Data Cleaning**

Type	Description	Examples	Recommendations for Data Cleaning
Outliers	Data values exceed maximum allowable value that can be recorded, or does not represent fact.	Value exceed maximum allowable value: Speeds > 200 mph [14]. Speeds > 415 mph and acceleration rate > 10 m/s/s [21]. Do not represent fact: Altitude > 30,000 ft [14].	Removed outliers [14][21]. Avoid to use all the values in the variable if outliers take a large portion of the data values (42% of data) [25].
Duplicated records	Duplicated records	Duplicated records in <i>DataWsu</i> data [21].	Check for duplicated data records and removed such records if they exist [21].
Invalid message	Invalid message	Invalid WSU or CAN Bus Message in <i>DataWsu</i> [21].	Filter records to remove invalid messages, e.g. filter with criterion " <i>GpsVaildWsu</i> = 1 and <i>VaildCanWsu</i> =1" [21].
Improperly recorded message	Activity recorded out of the scale that a sensor is designed for, thus recorded data values didn't fall into a normal data range.	Mobileye sensors may record the speed of vehicles in opposite directions if the road is narrow and doesn't have a median. Thus, the <i>Rangerate</i> values (speed of leading vehicle - speed of the following vehicle) in <i>DataFrontTargets</i> data may be negative and its absolute value is greater than the host vehicle's speed, indicating that the leading vehicle is backing up at a speed of <i>Rangerate-GpsSpeedWsu</i> [21].	Filter records to improperly recorded message, e.g. filter with criterion " <i>Rangerate-GpsSpeedWsu</i> >1" [21].

1 **Table 6 Summary of Data Processing**

Topics	Studies	Dataset	Data Sample Information	Data Cleaning and Processing	Software
Driving Pattern Identification	[14]	DAS	968,522 records of basic safety messages, from 155 trips made by 49 vehicles	Removed observations with errors. (Speed >200 mph, and altitude > 30,000 ft). Data visualization to show the extent of instantaneous driving volatility.	R, MATLAB, and Google Earth for data processing and visualization. Stata for modeling
	[15]	DAS	1,399,084 records of basic safety messages, from 184 trips made by 71 vehicles	Error-checked by linking microscopic trip data with a trip-summary file. Two datasets matched in terms of trip-level information. Data aggregation from 10 BSM per second to 1 BSM per second.	R, MATLAB, and Google Earth for data processing and visualization. Stata 14.1 for modeling
	[16]	BsmP1 (BSM)	1.5 billion rows of data, data size 204 GB	Data records on the eastbound of a horizontal curve were selected. East of (42.299469, -83.724666) were eliminated for study design.	R for data processing and extract information, Google Earth for extracting GPS coordinates.
Surrogate Safety Measures	[20]	<i>DataWsu</i> and <i>DataFrontTargets</i> (DAS1)	<i>DataWsu</i> file of 12 GB and <i>DataFrontTargets</i> file of 4.34 GB from nearly 100 vehicles.	Two datasets were read by Python to check data type and data organization. Then import to Hadoop for query using Apache Hive. Next, datasets were exported into small files and joined in PostgreSQL database. Fourthly, ArcGIS were used to ingrate link and intersection information. Fifthly, the data points around the intersections were removed by a 75-ft buffer zones created in PostgreSQL. Data processing framework figure in page 8	Python, Hadoop, Apache Hive, PostgreSQL, and ArcGIS,
	[21]	<i>DataWsu</i> and <i>DataFrontTargets</i> (DAS1)	62,589,725 BSMs from 90 vehicles.	<i>DataFrontTargets</i> file were filed for observations with vehicles in front of the host vehicle. Duplicates were removed from <i>DataWsu</i> file. Then <i>DataWsu</i> file were filtered remove invalid bus messages. Next, two cleaned datasets were merged, and cleaned to remove outliers. Finally, the datasets filtered out the vehicle movement in the opposite direction. Data process procedure figure in page 314	R for data manipulation and ArcGIS for spatial processing. R package ggmap for data visualization.
	[25]	Not stated	Not stated	Checked data values, 42% of data, 27,240,788 data point, had the lateral acceleration values exceeded the maximum value that the wireless communications device could record. Lateral acceleration variable was not used in this study.	Not stated

	[26]	BSM (RSE)	215,000,000 BSMS at selected intersections.	Data examination and error-checking process before data integration. Extra intersection data using geocodes to map the intersection location data from BSM, well matched with the real map. Appropriate geocoded polygons are used to filter BSM data for each selected intersection. A data integration and processing steps showed as a figure in page 295	Not stated	1 2 3
	[27]	Not stated	230 million BSMS	Appropriate geocodes are used to filter BSM data for each selected intersection. Zero speeds are removed from BSM data.	Stata's MATA language for modeling, and WinBUGS software for MCMC Gibbs sampling. noted that computations took long time at workstation level computer (Dell Precision T7600, 3.1 GHZ (32CPUs)).	
Signalized Intersection Operation Improvement	[32]	BSM data and SPAT (RSE)	Not stated	First select an interested movement and select GPSdata associated with the movement and time period based on direction of CV trajectories and prepare corresponding signal status data. Then, based on road geometry, calculate CVs' longitudinal position along the road from GPS positions, and generate time-space trajectories. Map GPS time into signal clock time and then aggregated trajectories to calculate the time dependent factor.	Not stated	
	[36]	SPAT and V2I driving records data (inferred to be BSM in RSE database)	2150 vehicle's daily trajectories	Not stated	Not stated	

1 **Table 7 Data-Driven Methods Applied to CV Data**

Topics	Studies	Technique and Purpose
Driving Pattern Identification	[14]	Negative binomial regression model: Correlation of extreme event frequency.
	[15]	Markov-switching dynamic regression model (time series analysis): Quantification and prediction of driving patterns.
	[16]	Random forest classification of risky driving behaviors.
Surrogate Safety Measures	[20]	Negative binomial regression model: Statistical relationship between the link developed safety surrogate measures and crash frequency.
	[21]	Pearson's correlation coefficients: correlation between developed safety surrogate measures and rear-end crashes.
	[25]	Fixed-and random-Poisson regression models: Quantification of the relationship between intersection-specific violations and crash frequency.
	[26]	Fixed-and random-Poisson regression models: Quantification of the relationship between intersection-specific violations and crash frequency.
	[27]	Hierarchical fixed- and random-parameter Poisson and Poisson log-normal models: Model crash function of intersection-specific volatilities and other factors. Full Bayesian estimation method and Markov Chain Monte-Carlo Gibbs sampler techniques: estimate parameter in Poisson models. Moran's I statistics: investigate correlation between crash frequency and spatial factors.
Signalized Intersection Operation Improvement	[32]	Time-dependent Poisson process: Model of traffic arrivals. Expectation Maximization: estimate parameter.
	[33]	Developed a mathematical model (without data mining) and used CV data to demonstrate the model.

2

3 **CONCLUSION**

4 This paper surveyed the literature and identified ten studies that used the only real-world
5 connected vehicle dataset currently available to the public—the large-scale CV datasets recently
6 released from the United States Department of Transportation (USDOT) Pilot Deployment
7 Program. This paper first provides a summary of the available datasets and describes their
8 organization, overall structure, and characteristics of the data captured during pilot deployments.
9 Secondly, the authors presented a summary of studies that used the data and classified the usage
10 into three categories involving application development, identifying driving patterns, developing
11 surrogate safety measures, and improving the operation of signalized intersections.

12 One common limitation indicated in some of studies is that only part of the dataset is
13 useful for analyses because of errors in the data collection processes and the large percentage of
14 erroneous attributes. All studies used one-day, one-month, or two-month CV data from the
15 USDOT program. A primary contribution of this review is a summary of current usage and
16 applications of the first and only dataset available to the public that contains real-world CV data.
17 This summary is in one place, thus providing a convenient reference to the research community.
18 Future work will extend the survey as more data becomes available.

19

20 **AUTHOR CONTRIBUTION STATEMENT:**

21 The authors confirm contribution to the paper as follows: study conception and
22 design: Yun Zhou and Raj Bridgelall; data collection: Yun Zhou and Raj Bridgelall; analysis and
23 interpretation of results: Yun Zhou and Raj Bridgelall; draft manuscript
24 preparation: Yun Zhou and Raj Bridgelall. All authors reviewed the results and approved
25 the final version of the manuscript. The authors do not have any conflicts of interest to declare.

26

REFERENCE

1. JPO, ITS. What Are Connected Vehicles and Why Do We Need Them? https://www.its.dot.gov/cv_basics/cv_basics_what.htm. Accessed on March 22, 2019
2. JPO, ITS. How Will Connected Vehicles Be Used? https://www.its.dot.gov/cv_basics/cv_basics_how_used.htm. Accessed on March 22, 2019
3. Greer, L., J., Fraser, D. Hicks, M. Mercer, and K. Thompson. Intelligent Transportation Systems Benefits, Costs, and Lessons Learned: 2018 Update Report. No. FHWA-JPO-18-641. United States Department of Transportation. Intelligent Transportation Systems Joint Program Office, 2018.
4. JPO, ITS. Safety Pilot Model Deployment Data. <https://datahub.transportation.gov/Automobiles/Safety-Pilot-Model-Deployment-Data/a7qq-9vfe>. Accessed on April 25, 2020
5. JPO, ITS. Connected Vehicle Pilot Deployment Program – Program Overall. https://www.its.dot.gov/pilots/pilots_overview.htm. Accessed on March 22, 2019
6. JPO, ITS. Safety Pilot Model Deployment Data, Safety Pilot Model Deployment Sample Data Handbook, [Safety Pilot Model Deployment Sample Data Handbook.docx](#). Accessed on March 22, 2019.
7. De Vlioger, I., D. De Keukeleere, and J.G. Kretzschmar. Environmental Effects of Driving Behaviour and Congestion Related to Passenger Cars. *Atmospheric Environment*. 2000. 34(27): 4649-55.
8. Simons-Morton, B.G., M.C. Ouimet, Z. Zhang, S.E. Klauer, S.E. Lee, J. Wang, P.S. Albert, and T.A. Dingus. Crash and Risky Driving Involvement Among Novice Adolescent Drivers and Their Parents. *American Journal of Public Health*. 2011.101(12) :2362-7.
9. Simons-Morton, B.G., Z. Zhang, J.C. Jackson, and P.S. Albert. Do Elevated Gravitational-Force Events While Driving Predict Crashes and Near Crashes? *American journal of epidemiology*. 2012. 175(10):1075-9.
10. Simons-Morton, B.G., K. Cheon, F. Guo, and P. Albert. Trajectories of Kinematic Risky Driving Among Novice Teenagers. *Accident Analysis & Prevention*. 2013. 51: 27-32.
11. Kim, E., and E. Choi. Estimates of Critical Values of Aggressive Acceleration from a Viewpoint of Fuel Consumption and Emissions. Transportation Research Board Annual Meeting, Washington, DC. 2013.
12. Paleti, R., N. Eluru and C.R. Bhat. Examining the Influence of Aggressive Driving Behavior on Driver Injury Severity in Traffic Crashes. *Accident Analysis & Prevention*. 2010. 42(6): 1839-1854.
13. Guo, F., and Y. Fang. Individual Driver Risk Assessment Using Naturalistic Driving Data. *Accident Analysis & Prevention*. 2013. 61:3-9.
14. Liu, J., and A.J. Khattak. Delivering Improved Alerts, Warnings, And Control Assistance Using Basic Safety Messages Transmitted Between Connected Vehicles. *Transportation research part C: emerging technologies*. 2016. 68: 83-100.
15. Khattak, A.J., and B. Wali. Analysis of Volatility in Driving Regimes Extracted from Basic Safety Messages Transmitted Between Connected Vehicles. *Transportation research part C: emerging technologies*, 2017. 84: 48-73.
16. Jahangiri, A., V.J. Berardi, and S.G., Machiani.. Application of Real Field Connected Vehicle Data for Aggressive Driving Identification on Horizontal Curves. *IEEE Transactions on Intelligent Transportation Systems*, 2018. 19(7): 2316-2324.

- 1 17. Wang, W. and J. Xi. A Rapid Pattern-Recognition Method for Driving Styles Using
2 Clustering-Based Support Vector Machines. In 2016 American Control Conference (ACC),
3 IEEE, 2016. 5270-5275.
- 4 18. Meiring, G. and H. Myburgh. A Review of Intelligent Driving Style Analysis Systems and
5 Related Artificial Intelligence Algorithms. *Sensors*, 2015. 15(12): 30653-30682.
- 6 19. Laureshyn A, C. Johnsson, T. De Ceunynck, Å. Svensson, M. de Goede, N. Saunier, and S.
7 Daniels. Review of current study methods for VRU safety. Appendix 6–Scoping review:
8 surrogate measures of safety in site-based road traffic observations. 2016.
- 9 20. He, Z., X. Qin, P. Liu and M.A. Sayed. Assessing Surrogate Safety Measures Using A Safety
10 Pilot Model Deployment Dataset. *Transportation Research Record*, 2018. 2672(38): 1-11.
- 11 21. Xie, K., D. Yang, K. Ozbay, and H. Yang. Use of Real-World Connected Vehicle Data in
12 Identifying High-Risk Locations Based on A New Surrogate Safety Measure. *Accident*
13 *Analysis & Prevention*. 2019. 125:311-319.
- 14 22. Yang, H., Z. Wang, and K. Xie. Impact of Connected Vehicles on Mitigating Secondary
15 Crash Risk. *International Journal of Transportation Science and Technology*. 2017. 3: 196-
16 207.
- 17 23. Abdel-Aty, M., and K. Haleem. Analyzing Angle Crashes at Unsignalized Intersections
18 Using Machine Learning Techniques. *Accident Analysis & Prevention*, 2011. 43(1): 461-470.
- 19 24. Persaud, B, T. Nguyen. Disaggregate Safety Performance Models for Signalized
20 Intersections on Ontario Provincial Roads. *Transportation Research Record*. 1998. 1635(1):
21 113-20.
- 22 25. Kamrani M., B. Wali, and A.J. Khattak. Can Data Generated by Connected Vehicles
23 Enhance Safety?: Proactive Approach to Intersection Safety Management. *Transportation*
24 *Research Record*. 2017. 2659(1): 80-90.
- 25 26. Kamrani M, R. Arvin, and A.J. Khattak. Extracting Useful Information From Basic Safety
26 Message Data: an Empirical Study of Driving Volatility Measures and Crash Frequency at
27 Intersections. *Transportation Research Record*. 2018. 2672 (38): 290-301
- 28 27. Wali B, A.J. Khattak, H. Bozdogan, and M. Kamrani. How is driving volatility related to
29 intersection safety? A Bayesian heterogeneity-based analysis of instrumented vehicles data.
30 *Transportation Research Part C: Emerging Technologies*. 2018.92:504-24.
- 31 28. Feng, Y., K.L. Head, S. Khoshmagham, and M. Zamanipour. A Real-Time Adaptive Signal
32 Control in a Connected Vehicle Environment. *Transportation Research Part C: Emerging*
33 *Technologies*. 2015. 55:460-73.
- 34 29. Guler, S.I., M. Menendez, and L. Meier. Using Connected Vehicle Technology to Improve
35 the Efficiency of Intersections. *Transportation Research Part C: Emerging Technologies*.
36 2014. 46:121-31.
- 37 30. Lee, J., B.B. Park, K. Malakorn, and J.J. So. Sustainability Assessments of Cooperative
38 Vehicle Intersection Control at an Urban Corridor. *Transportation Research Part C:*
39 *Emerging Technologies*. 2013. 32:193-206.
- 40 31. Day, C.M., and D.M. Bullock. Detector-Free Signal Offset Optimization with Limited
41 Connected Vehicle Market Penetration: Proof-of-Concept Study. *Transportation Research*
42 *Record*. 2016. 2558(1):54-65.
- 43 32. Zheng, J., and H.X. Liu. Estimating Traffic Volumes for Signalized Intersections Using
44 Connected Vehicle Data. *Transportation Research Part C: Emerging Technologies*. 2017.
45 79:347-62.

- 1 33. Newell, G.F. Approximation Methods for Queues with Application to the Fixed-Cycle 39
2 Traffic Light. *Siam Review*, 1965. 7(2): 223-240.
- 3 34. Mirchandani, P.B., and N. Zou. Queuing Models for Analysis of Traffic Adaptive Signal
4 Control. *IEEE Transactions on Intelligent Transportation Systems*, 2007. 8(1): 50-59.
- 5 35. Chang, T.H. and J.T. Lin. Optimal Signal Timing for an Oversaturated Intersection.
6 *Transportation Research Part B: Methodological*, 2000. 34(6): 471-491.
- 7 36. Zhao S, and K. Zhang. Observing Space-time Queueing Dynamics at a Signalized
8 Intersection using Connected Vehicles as Mobile Sensors. InProc., 96th Annual Meeting of
9 the Transportation Research Board. 2017.
- 10 37. Nodine, E., S. Stevens, A. Lam, C. Jackson, and W.G. Najm. Independent Evaluation of
11 Light-Vehicle Safety Applications Based on Vehicle-To-Vehicle Communications Used in
12 the 2012-2013 Safety Pilot Model Deployment. No. DOT HS 812 222. United States.
13 National Highway Traffic Safety Administration, 2015.
- 14 38. Bridgelall, Raj, P. Lu, D.D. Tolliver, and T. Xu. Mining Connected Vehicle Data for
15 Beneficial Patterns in Dubai Taxi Operations. *Journal of Advanced Transportation*, 2018
16 (2018).
- 17 39. Han, J, M. Kamber, and J. Pei. In *Data mining: concepts and techniques*. Elsevier, 2011.
- 18 40. Wang, W., and J. Xi. A rapid pattern-recognition method for driving styles using clustering-
19 based support vector machines. In 2016 American Control Conference (ACC), IEEE.
20 2016:5270-5275.,
- 21 41. Jahangiri, A., and H. A. Rakha. Applying machine learning techniques to transportation
22 mode recognition using mobile phone sensor data. *IEEE transactions on intelligent*
23 *transportation systems*. 2015.16(5): 2406-2417.