

Received August 30, 2018, accepted September 27, 2018, date of publication October 4, 2018, date of current version October 29, 2018.

Digital Object Identifier 10.1109/ACCESS.2018.2873569

# Congestion Prediction With Big Data for Real-Time Highway Traffic

FAN-HSUN TSENG<sup>1</sup>, (Member, IEEE), JEN-HAO HSUEH<sup>2</sup>,  
CHIA-WEI TSENG<sup>2</sup>, (Student Member, IEEE),  
YAO-TSUNG YANG<sup>2</sup>, (Student Member, IEEE),  
Han-Chieh Chao<sup>3,4</sup>, (Senior Member, IEEE),  
AND LI-DER CHOU<sup>2</sup>, (Member, IEEE)

<sup>1</sup>Department of Technology Application and Human Resource Development, National Taiwan Normal University, Taipei 106, Taiwan

<sup>2</sup>Department of Computer Science and Information Engineering, National Central University, Taoyuan 320, Taiwan

<sup>3</sup>Department of Electrical Engineering, National Dong Hwa University, Hualien 974, Taiwan

<sup>4</sup>Department of Computer Science and Information Engineering and the Department of Electronic Engineering, National Ilan University, I-Lan 260, Taiwan

Corresponding author: Li-Der Chou (cld@csie.ncu.edu.tw)

This work was supported in part by the Young Scholar Fellowship Program by the Ministry of Science and Technology (MOST), Taiwan, under Grant MOST107-2636-E-003-001, and in part by the MOST under Grant MOST102-2221-E-008-039-MY3, Grant MOST103-2221-E-008-090-MY3, Grant MOST104-2221-E-008-039-MY3, Grant MOST105-2221-E-008-071-MY3, and Grant MOST107-2221-E-259-005-MY3.

**ABSTRACT** By collecting and analyzing a vast quantity and different categories of information, traffic flow and road congestion can be predicted and avoided in intelligent transportation system. However, how to tackle with these big data is vital but challenging. Most of the existing literatures utilized batch method to process a bunch of road data that cannot achieve real-time traffic prediction. In this paper, we use the spouts and bolts in Apache Storm to implement a real-time traffic prediction model by analyzing enormous streaming data, such as road density, traffic events, and rainfall volume. The proposed SVM-based real-time highway traffic congestion prediction (SRHTCP) model collects the road data from the Taiwan Area National Freeway Bureau, the traffic events reported by road users from the Police Broadcasting Service in Taiwan, and the weather data from the Central Weather Bureau in Taiwan. We use fuzzy theory to evaluate the traffic level of road section in real time with considering road speed, road density, road traffic volume, and the rainfall of road sections. In addition, the SRHTCP model predicts the road speed of next time period by exploring streaming traffic and weather data. Results showed that the proposed SRHTCP model improves 25.6% prediction accuracy than the prediction method based on weighted exponential moving average method under the measurement of mean absolute relative error.

**INDEX TERMS** Big data, fuzzy theory, intelligent transportation system, real-time streaming data, support vector machine.

## I. INTRODUCTION

According to the study of International Data Corporation, the usage of digital data worldwide is about 1.8 zettabytes in 2011. The study further predicts that data amount will be 44 times more than nowadays in 2020, about 35.2 zettabytes. These digital data are generated by a variety of different ways. The online auction company eBay achieves online transactions in millions every day. There are more than 88 million users and more than millions of merchandise queries so that eBay's database increases more than 50 terabytes data every day. To analyze user behavior, the online system of eBay deals

with more than 50 petabytes data and executes more than 5 thousands items of business analysis per day. The enormous amount of data information is regarded as big data [1].

The well-known technology research Gartner points out that big data should be provided with high capacity, high growth capacity and high variability of characteristics. In 2001, Doug [2] pointed out that there are three data growth directions, i.e., volume (data size), velocity (data transfer speed) and variety (diversity of information). Big data feature. In [3], Chandarana and Vijayalakshmi pointed out that the characteristics of 5V of big data.

- 39 • *Volume*: A great quantity of data has been created  
40 quickly since everyone has one or more than one mobile  
41 devices. The data size is large and growing quickly  
42 reached terabytes even petabytes.
- 43 • *Velocity*: The data will be generated rapidly. In big data,  
44 the data generation and deletion, data flow, data change,  
45 and data processing are fast beyond our imagination.
- 46 • *Variety*: Based on data structure, big data can be classi-  
47 fied into two categories, i.e., structured and unstructured  
48 data. There are various diversities of data types and  
49 forms, e.g., data from different mobile phones, differ-  
50 ent social media data like Facebook or Twitter, data  
51 from sensor networks, vehicle-to-vehicle and device-to-  
52 device networks.
- 53 • *Veracity*: The uncertainty and reliability of big data is a  
54 big issue. The unstructured data usually has problems  
55 in imprecision of data. For instance, text usually has  
56 more than one meaning. Therefore, big data should be  
57 analyzed preciously.
- 58 • *Value*: The value of big data is a vital issue. How to  
59 explore the real value of big data by processing enor-  
60 mous and varied data is an important research topic.

61 The growing speed of automotive industry is as fast as  
62 big data. The Organisation Internationale des Constructeurs  
63 d'Automobiles (OICA) [4] pointed out that there are 70 mil-  
64 lion cars produced worldwide every year in recent years.  
65 It even reaches 90 million cars in both 2013 and 2014.  
66 According to Ministry of Transportation and Communica-  
67 tions (MOTC) Republic of China (R.O.C.) [5] statistics, vehi-  
68 cle usage in highways attains to the amount of more than  
69 530 million every year during 2010 to 2013. Since 2014,  
70 the National Freeway Bureau in Taiwan launches its new  
71 "Pay as You Go" toll system. As a result, it can be observed  
72 that highways play an important role in Taiwan no matter in  
73 cities, towns or rural areas. However, the growing highway  
74 usage raises the probability of traffic jam, e.g. periodical  
75 congestion sections and sudden traffic accidents. Various  
76 literatures investigate traffic jam issues but most of them  
77 focus on historical data analysis. On the other hand, existing  
78 literatures dichotomize traffic status into traffic jam or not as  
79 their prediction results. It is unable to describe driver's feel-  
80 ing of congestion event precisely. Unlike existing literatures,  
81 the contribution of this work is as follows.

- 82 • The paper proposes an SVM-based Real-time Highway  
83 Traffic Congestion Prediction (SRHTCP) model that  
84 instantaneously forecasts the car speed of next time  
85 period and analyzes traffic jams in highways.
- 86 • The traffic analysis and prediction is accomplished by  
87 collecting different data formats and sources. To deal  
88 with traffic data, more than 150 thousand data in Taiwan  
89 Area National Freeway Bureau (TANFB) collected from  
90 3617 vehicle detectors (VDs) along highways in Taiwan,  
91 9.6 thousand real-time weather data from the Central  
92 Weather Bureau of Taiwan (CWBT) and social media  
93 from the Police Broadcasting Service of Taiwan (PBST)  
94 are processed.

- 95 • The paper uses Apache Storm [6] to process real-time  
96 streaming data, and utilizes fuzzy theory to analyze  
97 driver's feeling of traffic jam. The proposed SRHTCP  
98 model is superior to other methods in terms of prediction  
99 accuracy.

100 The rest of this paper is organized as follows. Section II  
101 discusses and compares related works. In Section III,  
102 the designed system based on Apache Storm framework is  
103 introduced. Section IV illustrates experiment and prediction  
104 results. Finally, Section V concludes this work.

## 105 II. RELATED WORKS

106 Based on the concept of big data, this work collects and  
107 analyzes various data types and formats to predict real-time  
108 highway traffic. Several technologies are used in the paper,  
109 i.e., Apache Storm, traffic theory, fuzzy theory and support  
110 vector machine.

### 111 A. APACHE STORM

112 Apache Storm [7] is an open source and distributed real-  
113 time computing system. It can easily and quickly process the  
114 undivided streaming data and has following features.

- 115 • *Widely used*: Apache Storm is able to process mes-  
116 sage and update database. Moreover, it can continuously  
117 query and report streaming data thus proceed with a  
118 great number of requests.
- 119 • *Scalability*: Apache Storm is provided with excellent  
120 scalability that enormous machines can be adjusted and  
121 configured at the same time. It achieves one million  
122 data processing per second by using 10 cluster nodes  
123 equipped with Apache Storm.
- 124 • *Availability*: A real-time system must ensure that each  
125 data are proceeded successfully. Apache Storm traces  
126 every data through its message ID to guarantee data  
127 availability.
- 128 • *Robustness*: The goal of Apache Storm is easy manage-  
129 ment. The administrator easily obtains user experience  
130 and monitors machines conveniently.
- 131 • *Fault tolerance*: While a machine is out of order, Apache  
132 Storm can reboot the broken machine without influenc-  
133 ing other on-line machines. It guarantees that a task  
134 can be executed infinitely unless the task is terminated  
135 manually.
- 136 • *Variety of programming languages*: Apache Storm is  
137 robust and flexible so that developers can program it  
138 with multiple different programming languages.

139 Apache Storm is composed of Nimbus, Supervisor and  
140 Zookeeper. Nimbus is the brain of Apache Storm cluster and  
141 runs on the master node. It is responsible for sending tasks to  
142 other nodes and monitoring operation status of the cluster.  
143 Note that there is only one Nimbus in the whole cluster.  
144 Supervisor tackles task reception and monitoring. It runs on  
145 all working nodes and turns on or off task process based on  
146 its received task. Zookeeper plays the communicator role in

147 Apache Storm cluster. All nodes send heartbeat message to  
148 Zookeeper.

149 The particular design of Apache Storm achieves real-time  
150 streaming data processing. The task submission program in  
151 Apache Storm is called *Topology*. A topology is composed of  
152 several spouts and bolts. The smallest processing unit is called  
153 tuple. These three components are introduced as follows.

- 154 • Tuple: The smallest unit to compose a stream. In the  
155 topology of Apache Storm, all data are transmitted with  
156 stream format.
- 157 • Spout: An interface between topology and streaming  
158 data. When data are settled and mapped, data spout to  
159 topology and forwarded to bolt for processing.
- 160 • Bolt: A bolt is an element for data processing in topol-  
161 ogy. Data calculations are conducted in bolts. A bolt  
162 forwards data to another bolt after its data processing.

163 The relation between spouts and bolts in a topology is  
164 illustrated with Fig. 1. When a spout receives streaming  
165 data, it divides the data into several tuple and forwards to  
166 corresponding bolts. Note that different bolts tackle with  
167 different data processing then output data or send to next  
168 bolt. Therefore, we need to design different topologies to  
169 accomplish different data processing.

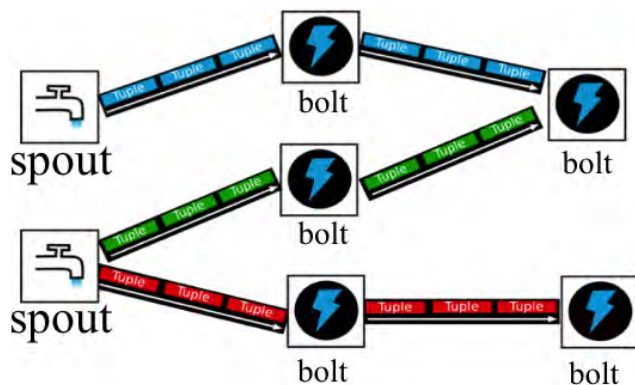


FIGURE 1. The relation between spouts and bolts in topology.

## 170 B. OPEN DATA

171 In Taiwan, government provides almost 30 thousand open  
172 data sets in the open data platform [8], e.g., weather, traffic,  
173 infrastructure, education, construction, election. In the paper,  
174 traffic data, weather data and social media data are integrated  
175 and analyzed.

- 176 • Traffic open data: The paper retrieved the traffic data  
177 of highways in Taiwan from the TANFB [9]. Open  
178 data in the platform includes speed of sections, VDs,  
179 closed-circuit televisions, changeable message signs and  
180 automatic vehicle identification. We use the information  
181 provided by VDs in this paper.
- 182 • Weather open data: The paper adopted weather data  
183 from the CWBT [10]. Open data in the platform includes  
184 weather forecasts, observation, earthquake, tsunami,

185 climate and weather. We use rainfall information from  
186 weather stations in this paper.

- 187 • Social media open data: The paper collects social media  
188 data from the PBST [11] so that the prediction tallies  
189 with real-time traffic status. Drivers are capable  
190 of informing traffic reports and events to PBST. Four  
191 kinds of traffic report are recorded in PBST, i.e., traffic  
192 barriers, road construction, traffic congestion and other  
193 events. In this paper, we use real-time traffic reports in  
194 PBST.

## 195 C. TRAFFIC THEORY

196 In traffic theory, traffic stream models are used to represent  
197 the relation between volume, speed and density. In [12],  
198 Greenshields proposed a parabolic function and named it  
199 as Greenshields model. The Greenshields model describes  
200 that the higher density and volume result in the lower car  
201 speed, and vice versa. It means that the probability of traffic  
202 congestion is higher while car density is rising. In [13],  
203 Nakayama *et al.* verified the relation between car density  
204 and speed with some experiments. The experiment is imple-  
205 mented in a ring road without any obstacle. The experiment  
206 is designed to examine that whether car volume influences car  
207 speed or leads to traffic accident or not. The results showed  
208 that traffic jam happens when there are 22 cars in a 230 meters  
209 ring road. Rainfall is an important factor that impacts on traf-  
210 fic flow. In [14], the most congested M42 highway in United  
211 Kingdom was investigated for analyzing drivers' behaviors.  
212 The researchers found that a sudden slowdown or change  
213 lines gives rise to the reason of traffic jam.

## 214 D. FUZZY THEORY

215 In 1965, Zadeh firstly proposed fuzzy theory [15]. The fuzzy  
216 theory was proposed to solve uncertainties in the real world  
217 by using computers. It uses Fuzzy Control Logic in its mech-  
218 anism, includes fuzzification, fuzzy database, fuzzy inference  
219 and defuzzification. In the fuzzification, input parameters  
220 are converted into the membership level of fuzzy sets by  
221 the membership function. Then, various "if-then" judgment  
222 equations are pre-defined in fuzzy database. Fuzzy inference  
223 estimates the membership level of input parameters based  
224 on the defined fuzzy database. Lastly, defuzzification step  
225 converts inputs into numerical results so that computers are  
226 able to determine final result. In the paper, we utilize fuzzy  
227 theory to evaluate traffic jam level.

## 228 E. SUPPORT VECTOR MACHINE

229 Support Vector Machine (SVM) [16] is a well-known classi-  
230 fication technique. It can be applied to solve various research  
231 issues, e.g., virtual machines classification in clouds [17],  
232 anomaly detection [17], data diagnosis [18] and sensor fault  
233 classification [19]. In classification process, the main con-  
234 cept of SVM is to construct an optimal hyper-plane to be  
235 the boundary for making decisions. Three characteristics of  
236 SVM are as follows. (i) It establishes the maximum boundary  
237 so that the coordinates of cluster nodes have the maximum

possible distance from the decision boundary. (ii) It uses kernel function to establish a separated linear hyper-plane, thus cluster nodes in higher dimension can be divided into multiple clusters easily. (iii) It adopts a small part of data as training data so that the prediction can be more accurate after training process. In the paper, we utilize SVM to forecast the car speed of highways in Taiwan.

**F. COMPARISON OF RELATED WORKS**

During past years in Intelligent Transportation System (ITS), various researchers have utilized different methods to monitor traffic events [21], [22] and to predict traffic congestion [23]–[27]. In [21], Milojevic and Rakocevic proposed a vehicle-to-vehicle congestion detection algorithm based on the IEEE 802.11p standard. The proposed algorithm permits vehicles to be self-aware so that vehicles are able to monitor speed and cooperate with each other. In [22], Cheng *et al.* proposed a new automatic incident detection method for urban expressways based on geometric conditions and detector locations. However, these two articles only focus on traffic detection rather the congestion prediction in this paper.

In [23], Ji *et al.* utilized Kalman filter with the Global Positioning System (GPS) location reported by drivers to forecast travel time dynamically. The results showed that the prediction accuracy of improved model is superior to the original method with Kalman filter. In [24], Feng *et al.* also utilized Kalman filter but to predict vehicle’s future location. Results showed that the proposed method is superior to a prediction method based on neural network. However, both two methods did not consider weather data and real-time traffic events. In addition, the authors investigated travel time prediction and vehicle’s location prediction rather than traffic congestion prediction in this paper. Moreover, the researchers in [23] did not consider highway.

In [25], Wang used Grey prediction to detect traffic incidents. The researcher used actual examples to compare the difference between prediction and reality. Results showed that the proposed method achieves acceptable false-alarm rate to determine whether incident happened. However, the paper is different to our work that traffic incidents are collected from PBST to forecast traffic jams in next time period. In [26], Kuo applied the Kalman filter to Support Vector Machine for achieving travel time prediction. The air pollution index was considered in the prediction model to fit in with real traffic status. However, the author did not consider traffic incidents and real-time weather data. In [27], Li *et al.* proposed a bipolar traffic density awareness routing protocol for vehicle ad hoc networks. The average inter-vehicle space of vehicle networks was predicted, but both traffic events and weather data were unmentioned. In addition, the category and scale of collected data in these three papers is less and much smaller than our work. Moreover, they did not collect real-time traffic and weather data which is achieved in our work. Most of existing literatures used a batch of traffic data

to predict car speed, however it cannot achieve instantaneous traffic prediction.

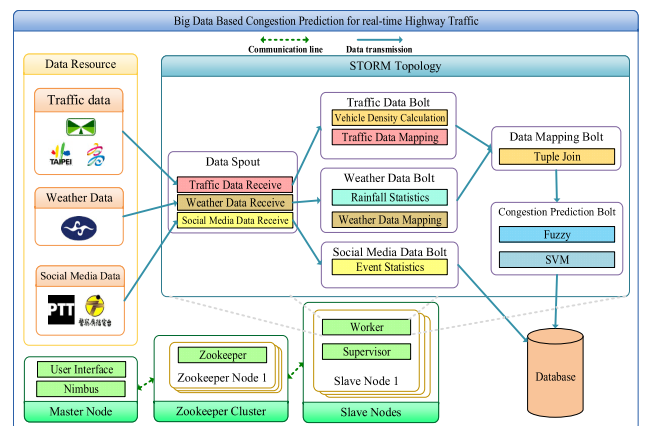
Some researchers have utilized Apache Storm to process streaming data. In [28], Wang *et al.* used Apache Storm their main processing system. The simulation-based result showed that the system can be applied to practical situation. In [29], Chardonnes *et al.* utilized Storm to realize real-time integration and detection on Twitter [30] message keyword statistics and Bitly [31] short Uniform Resource Locator (URL) location statistics. In [32], Bifet and Morales proposed an open source platform called Scalable Advanced Massive Online Analysis (SAMOA). The proposed SAMOA integrates Storm for mining big data streams. In this paper, we utilize the Apache Storm platform to collect traffic data, weather data and social media data instantaneously, and then predict the car speed of freeways in real-time.

**III. SYSTEM DESIGN**

Based on the concept of big data, the paper proposes a mechanism that uses Apache Storm system to collect and analyze freeway traffic information, weather information and social network information. The proposed SRHTCP model utilizes SVM to forecast the speed of next road section, and instantaneously evaluates freeway condition by using fuzzy theory.

**A. SYSTEM ARCHITECTURE AND WORKFLOW**

The system architecture of this work is captured in Fig. 2. The proposed cluster architecture is based on the open source distributed real-time computation system *Apache Storm* (or simply Storm throughout this paper). With the master node, it contains Nimbus and user interface, Zookeeper for data exchange. The slave node is composed of supervisor and worker. The proposed congestion prediction method for real-time freeway traffic is implemented in the Storm topology framework. It is composed of six components, i.e., data spout, traffic data bolt, weather data bolt, social media data bolt, data mapping bolt and congestion prediction bolt. The traffic data bolt calculates vehicle density on roads then sends



**FIGURE 2. System architecture.**

329 the calculation results to the traffic data mapping module. The  
 330 weather data bolt gathers rainfall statistics then examines and  
 331 sends to weather data mapping module. The social media data  
 332 bolt collects four types of road events then sends to database.  
 333 The computing cluster of Storm is composed of master  
 334 node, slave node and Zookeeper. The master node is regarded  
 335 as the brain of Apache Storm that executes the Storm Nimbus  
 336 procedure. It submits and distributes all tasks by the Nimbus,  
 337 and provides administrator with user interface to manage  
 338 Storm system. The slave node executes Storm Supervisor pro-  
 339 cedure that distributes tasks at any time. The worker module  
 340 is launched once the slave node receives task. The Zookeeper  
 341 plays the communicator role in Storm cluster and records the  
 342 system status.

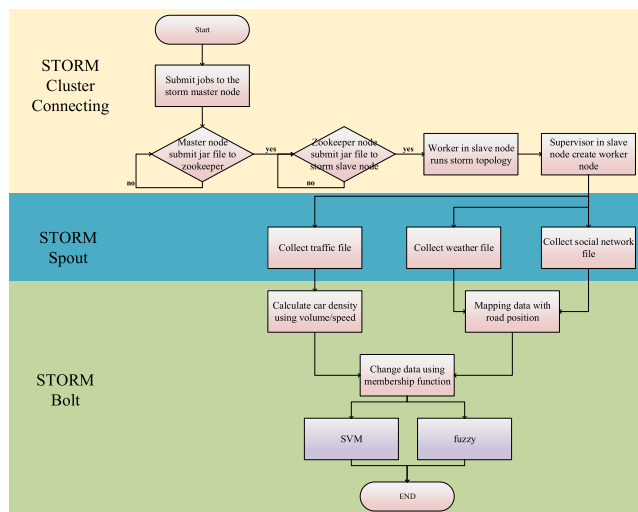


FIGURE 3. System workflow.

343 The system workflow is shown as Fig. 3. The constructed  
 344 Storm cluster is responsible for job execution by creating  
 345 master and slave node. In slave node, worker nodes are  
 346 created to run Storm topology. After that, the Storm spout  
 347 collects traffic, weather and social network data from the  
 348 TANFB, CWBT and PBST. In the Storm bolt layer, Fuzzy  
 349 theory is utilized to analyze real-time traffic congestion level,  
 350 and SVM is utilized to forecast car speed of lane in next time  
 351 period.

**B. DATA SPOUT AND BOLTS**

352 The Apache Hadoop applies batch method to separately pro-  
 353 cess enormous data, however it cannot immediately process  
 354 data. As a result, Apache Storm was proposed to deal with  
 355 real-time streaming data through the spout and bolt. The  
 356 designed data spout and bolts are explained as follows.  
 357

**1) DATA SPOUT**

358 The data spout for data collection is composed of traffic  
 359 data receive module, weather data receive module and social  
 360 media data receive module. The traffic data receive module  
 361 is responsible to retrieve data from Taiwan Area National  
 362

Freeway Bureau to Storm system. Firstly, the module calls the  
 363 open() function in spout to set the traffic data path for loading  
 364 the XML file about traffic status. Then, the nextTuple() func-  
 365 tion reads data from the XML file one by one. A piece of data  
 366 includes device ID, lane number, vehicle speed in lane and the  
 367 number of vehicles with different types. Note that four types  
 368 of vehicles are used in the paper, i.e., motorcycle, compact,  
 369 van and trailer.  
 370

The weather data receive module collects weather data  
 371 from CWBT to Storm system. The module calls the open()  
 372 function in spout to retrieve the weather of a location from  
 373 the XML file about weather information, then reads data by  
 374 using nextTuple() function from the XML file one by one.  
 375 Note that the module retrieves the rainfall value of selected  
 376 location in every ten minutes.  
 377

The social media data receive module collects media data  
 378 from social networks for reporting traffic conditions. The  
 379 retrieve method is same with the other two modules but from  
 380 the JSON file. The collected media data includes the informa-  
 381 tion of event such as location, freeway number, direction and  
 382 type. For instance, a piece of data reported from social media  
 383 data is recorded as freeway number 1, 235 km at southbound,  
 384 and car accident.  
 385

**2) TRAFFIC DATA BOLT**

386 The traffic data bolt for data processing contains vehicle den-  
 387 sity calculation module and vehicle speed calculation mod-  
 388 ule. The vehicle density calculation module receives tuple  
 389 data from the data spout by using execute() function. The  
 390 traffic volume is formulated as  
 391

$$q = ku, \tag{1}$$

392 where  $q$  is traffic volume as well as the number of vehicles  
 393 per hour,  $k$  is road density that represents the average number  
 394 of vehicles in one kilometer, and  $u$  is car speed in kilometer  
 395 per hour. The obtained road density is attached behind the  
 396 tuple data of traffic data receive module and sent to traffic  
 397 data mapping module with new tuple format.  
 398

The traffic data mapping module receives the tuple data  
 399 from vehicle density calculation module. It loops up the  
 400 value in density field from the tuple data, and retrieves the  
 401 corresponding latitude and longitude then adds to the tuple  
 402 data. The updated tuple data is sent to data mining bolt by  
 403 using the emit() function.  
 404

**3) WEATHER DATA BOLT**

405 The weather data bolt contains rainfall statistics module and  
 406 weather data mapping module. The rainfall statistics module  
 407 receives weather tuple data from data spout by using execute()  
 408 function. It examines the correctness of the weather tuple  
 409 data. If the rainfall value is less than zero, the weather station  
 410 observes none of rain so that the negative value should be  
 411 revised to zero. The updated rainfall value in weather tuple  
 412 data is sent to weather data mapping module.  
 413

The weather data mapping module receives weather tuple  
 414 data from the rainfall statistics module and retrieves latitude  
 415

416 and longitude value. Then it searches vehicle detectors with  
 417 1 km distance from the weather station, and adds vehicle  
 418 detector ID to the weather tuple data. The updated weather  
 419 tuple data is transmitted to the data mapping bolt by using  
 420 the emit() function.

421 4) OTHER BOLTS

422 In this subsection three bolts are introduced, i.e., social media  
 423 data bolt, data mapping bolt and congestion prediction bolt.  
 424 In social media data bolt, the event statistics module uses execute()  
 425 function to receive tuple data from data spout for data  
 426 selection. The tuple data includes four types of road events,  
 427 i.e., traffic congestion, traffic barrier, road construction and  
 428 other events. The event statistics module uses emit() function to  
 429 transfer road type reports to database.

430 In data mapping bolt, the tuple join module receives vehicle  
 431 tuple and weather tuple data from the traffic data bolt and  
 432 weather data bolt respectively. It integrates vehicle detector  
 433 tuple and weather tuple with the same vehicle detector IDs.  
 434 The integrated tuple data records vehicle detector's status,  
 435 lane ID, average speed of lane in kilometer per hour, vehicle  
 436 type, and the volume of lane. There are three different statuses  
 437 of a vehicle detector, i.e., functional, time-out and malfunction.  
 438 The lane volume represents the number of vehicles in a  
 439 lane in one minute.

440 The fuzzy module and SVM module are in congestion  
 441 prediction bolt. The fuzzy module converts input data into  
 442 different membership functions of fuzzy set. The fuzzy inference  
 443 is executed and completed based on fuzzy database.  
 444 The output data is the membership between traffic congestion  
 445 level and other parameters. Last, the fuzzy module quantifies  
 446 output results through defuzzification and acquires the real-  
 447 time traffic congestion level. On the other hand, the SVM  
 448 module predicts the speed of next time period based on the  
 449 received tuple data from tuple join module. The historical  
 450 data includes real-time speed, the speed of 5 and 10 minutes  
 451 before.

452 C. DATA PROCESSING AND PARAMETER SELECTION

453 With the vehicle data, Taiwan Area National Freeway Bureau  
 454 provides data from vehicle detectors every 5 minutes. The  
 455 used parameters and the flowchart of vehicle data processing  
 456 are listed in Table 1 and captured in Fig. 4. In the beginning,  
 457 system tries to access data in Taiwan Area National Freeway  
 458 Bureau. If the system cannot retrieve data in time, it will  
 459 wait for a time period to access in next time period. After  
 460 obtaining the vehicle information by using XML parser to  
 461 read the data, the XML parser segments data into several data  
 462 sets based on vehicle detector. The parameter  $VD\_status_x$   
 463 represents the status of vehicle detector  $x$ . If  $VD\_status_x = 1$ ,  
 464 the connection between vehicle detector  $x$  and Taiwan Area  
 465 National Freeway Bureau is disconnected. If  $VD\_status_x = 2$ ,  
 466 the vehicle detector is malfunction. The output data  
 467 are negative in both statuses. In order not to affect the  
 468 fuzzy and SVM modules, the parameters  $Speed_{x, lane}$  and  
 469  $Density_{x, lane}$  of this vehicle detector are set to 1. If the

TABLE 1. Parameter of vehicle data.

Parameter	Description
$X$	Total number of vehicle detector
$VD\_URL$	Traffic file's URL
$ID_x$	Vehicle detector ID
$VD\_status_x$	The status of vehicle detector $x$
$lane$	Total number of lanes
$Speed_{x, lane}$	Average speed for each lanes and $ID_x$
$Volume_{x, lane}$	Car volume of each lane for $ID_x$
$Density_{x, lane}$	Density of each lane for $ID_x$
$Data\_time$	Data collected time

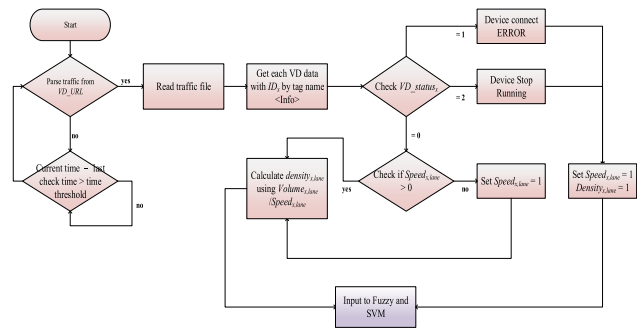


FIGURE 4. Flowchart of vehicle data processing.

parameter  $VD\_status_x = 0$ , the vehicle detector is functional  
 so that the obtained  $Density_{x, lane}$  and  $Speed_{x, lane}$  are sent to  
 fuzzy module and SVM module respectively. The parameter  
 $Density_{x, lane}$  is calculated as

$$Density_{x, lane} = \frac{Volume_{x, lane, car\_all}}{Speed_{x, lane}}, \quad (2)$$

where  $volume_{x, lane, car\_all}$  is the obtained car volume of a line  
 from vehicle detector  $x$ .

With the weather data, the Central Weather Bureau in  
 Taiwan provides information from rainfall detector every ten  
 minutes. The used parameters and the flowchart of weather  
 data processing are listed in Table 2 and captured in Fig. 5.  
 First, the system tries to access data from Central Weather  
 Bureau, and waits for next time period if it fails to access the  
 data. After obtaining the weather data by using XML parser,  
 the XML parser segments data into several data sets based on  
 rainfall detector. If the parameter  $Value\_10_x < 0$ , the connection  
 between rainfall detector  $x$  and Central Weather Bureau  
 is incorrect. To avoid influencing fuzzy module, the rainfall

TABLE 2. Parameter of weather data.

Parameter	Description
$ST\_URL$	Weather file's URL.
$X$	Total number of rainfall detector
$ID_x$	Rainfall detector's ID
$ST\_lat_x$	Latitude of rainfall detector $x$
$ST\_lon_x$	Longitude of rainfall detector $x$
$Value\_10_x$	Rainfall value in 10 minute at rainfall detector $x$
$Data\_time$	Data collected time

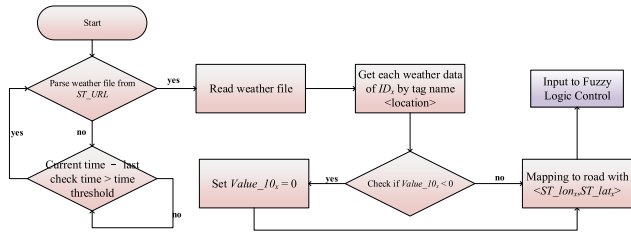


FIGURE 5. Flowchart of weather data processing.

data  $Value_{10}_x$  is set to zero. Then the system sends corresponding rainfall information to adjacent vehicle detectors range from the rainfall detector in 1 km.

The social media data is retrieved from the Police Broadcasting Service in Taiwan every ten minutes. The used parameters and the flowchart of social media data processing are listed in Table 3 and captured in Fig. 6. At the beginning, the system tries to access social data and will retry to access data in next time period if it fails to obtain the social media data of current time period. The parameter  $Event\_type_x$  is utilized to check the type of traffic event  $x$ , such as traffic barrier, car accident, block, road construction. If the traffic event belongs to traffic congestion and traffic barrier, the reported event is mapped to the vehicle detectors. Therefore the reported event can be looked up to those vehicle detectors in the range of this event. Based on the comparison result, the system examines lane status and analyzes several reported information, e.g., average speed, the number of report events from different sections, and the location of report events.

TABLE 3. Parameter of social media data.

Parameter	Description
$EV\_URL$	Social media file's URL
$X$	Total number of social media event
$Event\_type_x$	Type of report traffic event $x$
$Event\_location_x$	Location of reported traffic event $x$
$Event\_des_x$	Description of reported traffic event $x$
$Event\_date_x$	Date of reported traffic event $x$
$Event\_time_x$	Time of reported traffic event $x$
$Event\_source_x$	Source of reported traffic event $x$

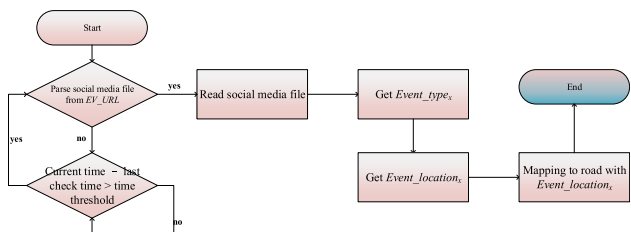


FIGURE 6. Flowchart of social media data processing.

#### D. CALCULATION OF CONGESTION LEVEL BY FUZZY

Vehicle data and weather data are sent to fuzzy module for road status analysis. The three stages of fuzzy theory are

TABLE 4. Parameter of fuzzy theory.

Parameter	Description
$P_s$	Car speed as the input of fuzzy logic control
$P_d$	Car density as the input of fuzzy logic control
$P_r$	Rainfall volume as the input of fuzzy logic control
$P_v$	Difference between historical and real-time car volume as the input of fuzzy logic control
$a$	The lower limit of membership function $P_s, P_d$
$b$	The middle limit of membership function $P_s, P_d$
$c$	The higher limit of membership function $P_s, P_d$
$d$	The lower limit of membership function $P_r, P_v$
$e$	The higher limit of membership function $P_r, P_v$

fuzzification, fuzzy inference and defuzzification. The used parameters in the proposed fuzzy module are listed in Table 4. In fuzzification stage, the membership function includes car speed, road density, rainfall and car volume. The fuzzy subsets of car speed with respect to low, normal and high are formulated as

$$f_{low}(P_s) = \begin{cases} 1, & \text{if } 0 \leq P_s < a \\ \frac{P_s - b}{b - a} + 1, & \text{if } a \leq P_s < b \\ 0, & \text{if } P_s \geq b, \end{cases} \quad (3)$$

$$f_{normal}(P_s) = \begin{cases} \frac{P_s - b}{b - a} + 1, & \text{if } a \leq P_s < b \\ \frac{a - P_s}{b - c} + 1, & \text{if } b \leq P_s < c \\ 0, & \text{otherwise,} \end{cases} \quad (4)$$

$$f_{high}(P_s) = \begin{cases} 0, & \text{if } P_s < b \\ \frac{P_s - c}{c - b} + 1, & \text{if } b \leq P_s < c \\ 1, & \text{if } b \leq P_s. \end{cases} \quad (5)$$

The fuzzy subsets of road density with respect to low, normal and high are formulated as

$$f_{low}(P_d) = \begin{cases} 1, & \text{if } 0 \leq P_d < a \\ \frac{P_d - b}{b - a} + 1, & \text{if } a \leq P_d < b \\ 0, & \text{if } P_d \geq b, \end{cases} \quad (6)$$

$$f_{normal}(P_d) = \begin{cases} \frac{a - P_d}{b - c} + 1, & \text{if } a \leq P_d < b \\ \frac{P_d - b}{b - a} + 1, & \text{if } b \leq P_d < c \\ 0, & \text{otherwise,} \end{cases} \quad (7)$$

$$f_{high}(P_d) = \begin{cases} 0, & \text{if } P_d < b \\ \frac{P_d - c}{c - b} + 1, & \text{if } b \leq P_d < c \\ 1, & \text{if } b \leq P_d. \end{cases} \quad (8)$$

The fuzzy subsets of rainfall volume with respect to low and high are formulated as

$$f_{low}(P_r) = \begin{cases} 1, & \text{if } 0 \leq P_r < d \\ \frac{P_r - e}{e - d} + 1, & \text{if } d \leq P_r < e \\ 0, & \text{if } P_r \geq e, \end{cases} \quad (9)$$

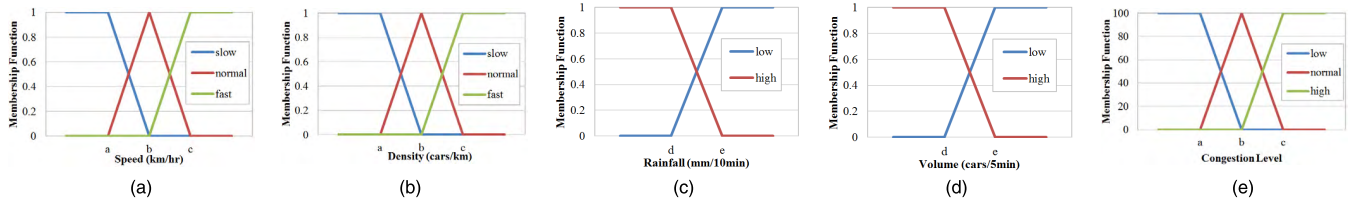


FIGURE 7. The membership function of proposed model. (a) Speed. (b) Density. (c) Rainfall. (d) Volume. (e) Congestion level.

$$f_{high}(P_r) = \begin{cases} 0, & \text{if } P_r < d \\ \frac{P_r - e}{e - d} + 1, & \text{if } d \leq P_r < e \\ 1, & \text{if } e \leq P_r. \end{cases} \quad (10)$$

The fuzzy subsets of car volume between historical and real-time with respect to low and high are formulated as

$$f_{low}(P_v) = \begin{cases} 1, & \text{if } 0 \leq P_v < d \\ \frac{P_v - e}{e - d} + 1, & \text{if } d \leq P_v < e \\ 0, & \text{if } P_v \geq e, \end{cases} \quad (11)$$

$$f_{high}(P_v) = \begin{cases} 0, & \text{if } P_v < d \\ \frac{P_v - e}{e - d} + 1, & \text{if } d \leq P_v < e \\ 1, & \text{if } e \leq P_v. \end{cases} \quad (12)$$

The obtained membership functions of car speed, road density, rainfall and car volume are shown as Fig. 7. In fuzzy inference stage, the fuzzy inference is conducted based on the fuzzy rules in fuzzy database. The fuzzy rule is defined as

$$R^{(l)}: \text{if } x \text{ is } A_i^l \text{ then } y \text{ is } B^l, \quad (13)$$

where  $R^{(l)}$  is the  $l$ -th rule,  $x$  and  $y$  are the input and output of fuzzy module. The parameters  $A_i^l$  and  $B^l$  are fuzzy sets. The proposed fuzzy inference is based on minimum inference mechanism, which is defined as

$$B^l(y) = \max_{1 \leq l \leq m} [A_1^l(x_1) \cdot A_2^l(x_2) \cdot B^l(y)]. \quad (14)$$

In the last stage, the proposed defuzzification is based on the center of area method and defined as

$$y = \frac{\int_Y yB(y)dy}{\int_Y B(y)dy}. \quad (15)$$

According the three stages, the congestion level of a road section can be quantified to 0 to 100.

**E. PREDICTION OF ROAD SPEED BY SVM**

The data format should be defined thus the car speed data can be applied to SRHTCP’s prediction process. The data format of SRHTCP is listed as Table 5. The label stands for the category of car speed per hour. First, the label is used to classify car speed. The index represents the dimension of SRHTCP’s training set as well as the data’s feature value. The value stands for the realistic value of the dimension.

In the work, the car speeds in previous three time periods are used to be the dimension of training data in SRHTCP. The

TABLE 5. Data format of SRHTCP.

[Label]	[Index1]:[Value1]	[Index2]:[Value2]	[Index3]:[Value3]
[Label]	[Index1]:[Value1]	[Index2]:[Value2]	[Index3]:[Value3]
[Label]	[Index1]:[Value1]	[Index2]:[Value2]	[Index3]:[Value3]
...	...	...	...

parameter  $S_t$  represents the car speed in time  $t$ , and  $S_{t-1}$ ,  $S_{t-2}$  and  $S_{t-3}$  are the car speed of time  $t - 1$ ,  $t - 2$  and  $t - 3$  respectively. The  $S_{t-1}$ ,  $S_{t-2}$  and  $S_{t-3}$  are used to conduct the car speed per hour of  $S_t$ . Therefore the travel time of  $S_t$  can be calculated by three previous time periods  $S_{t-1}$ ,  $S_{t-2}$  and  $S_{t-3}$ . The flowchart of data processing is captured in Fig. 8.

First of all, parsing input variable initializes the value of  $feature\_amount_g$ . It means that how many features are needed to describe each label in the calculation of SRHTCP. In addition, the number of support vector  $h$  is initialized to zero. Then, features are retrieved from each data file so that SRHTCP examines where the number of feature is larger than the  $feature\_amount_g$  author defined or not. If the number of features is less than the defined  $feature\_amount_g$ , it means the number of features is insufficient so that the label’s feature number should be increased continuously. When feature number exceeds the  $feature\_amount_g$ , SRHTCP starts to retrieve next feature until there is no more support vectors.

After retrieving features, the proposed system classifies training data into several groups based on the parameters user defined. The system utilizes one arbitrary parameter to train and predict the accuracy of this training data. Then it uses another parameter to execute data training and others are used to predict accuracy. After the training process, one of the parameters is selected to be the optimal solution and utilized in support vectors for data training. After training process, the trained model predicts data sets and obtains final result.

**IV. EXPERIMENTAL RESULT AND ANALYSIS**

In this section, congestion prediction with big data for free-way traffic is implemented. Based on the concept of big data, Apache Storm is used to implement platform that collects



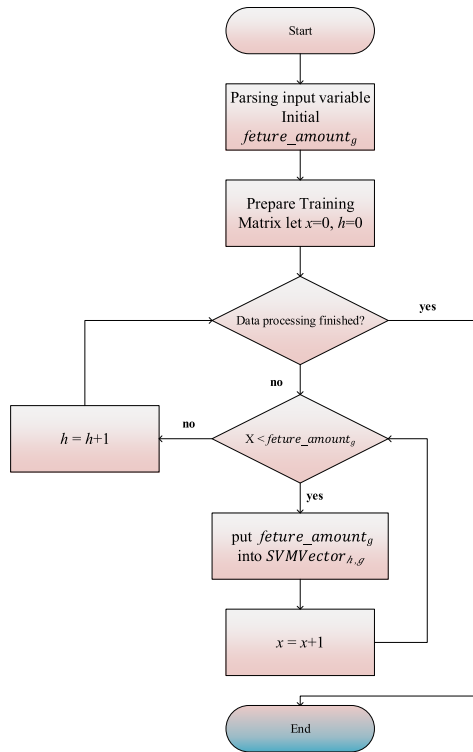


FIGURE 8. The flowchart of data processing.

TABLE 6. Specification of XenServer.

Hardware/Software	Specification
CPU	AMD Opteron(tm) Processor 6172
Memory	16382 Mb
Network interface card	Intel® 82576 Gigabit Ethernet Controller
Network interface	1000 Base TX
Operating system	Xen server 6.2

592 traffic, weather and social data. The congestion level of  
 593 freeway is analyzed by the proposed fuzzy model and the  
 594 car speed of next time period is predicted by the proposed  
 595 SRHTCP model.

596 **A. EXPERIMENT SET UP AND PARAMETERS**

597 The experiments in this paper are implemented by the VMs in  
 598 XenServer. The hardware and software of used environment  
 599 in XenServer are captured in Table 7. In the constructed  
 600 Apache Storm platform, there is a master node, a ZooKeeper  
 601 node and a slave node. Some of used hardware and software  
 602 resources in the master node, ZooKeeper node and slave  
 603 node are the same with used specification in XenServer, i.e.,  
 604 CPU, network interface card and network interface. In master,  
 605 ZooKeeper and slave node, the used memory is 4096 Mb.  
 606 In master node, the operating system (OS) version is Ubuntu  
 607 14.04.2 LTS and the OS kernel is 3.13.0-55-generic. The  
 608 operating system of ZooKeeper node is XenServer 6.2. Both  
 609 master node and slave node are equipped with Apache Storm  
 610 in version 0.9.0.1.

**B. MEASUREMENT OF PREDICTION ACCURACY**

611 In the paper, the prediction accuracy is measured by the  
 612 Mean Absolute Relative Error (MARE) and Mean Square  
 613 Error (MSE) methods. The MARE is calculated by  
 614

$$615 \text{MARE} = \frac{1}{n} \sum_{i=1}^n \frac{|a_i - b_i|}{a_i}, \quad (16)$$

616 where  $a_i$  is the prediction result and  $b_i$  is the observation data.  
 617 The MSE is calculated by

$$618 \text{MSE} = \frac{1}{n} \sum_{i=1}^n (a_i - b_i)^2, \quad (17)$$

619 where  $a_i$  and  $b_i$  is the value of prediction and observation data  
 620 respectively.

621 The MARE and MSE are familiar methods to evaluate the  
 622 difference between prediction and reality [33]. The MARE is  
 623 a percentage and the MSE is a value. The MARE percentage  
 624 represents the difference between prediction result and obser-  
 625 vation data. A lower MARE percentage stands for the higher  
 626 prediction accuracy, and vice versa. However, the MARE  
 627 is hard to explore the difference between prediction and  
 628 observation when the error difference is small. The MSE is  
 629 more useful to enlarge error difference so that the prediction  
 630 error will be more obvious. In other words, the prediction  
 631 accuracy and inaccuracy are easier to be explored. A lower  
 632 MSE implies that the higher prediction accuracy, and vice  
 633 versa. Both methods are used to evaluate prediction accuracy  
 634 in the paper.

635 **C. EXPERIMENT 1: REAL-TIME TRAFFIC ANALYSIS BASED  
 636 ON TRAFFIC REPORTS FROM DRIVERS**

637 In this experiment, the real-time traffic reports from Police  
 638 Broadcasting Service in Taiwan are collected and matched  
 639 with the information retrieved by those vehicle detectors  
 640 where in the same region. As a result, the matched informa-  
 641 tion is more valuable than the information from vehicle detec-  
 642 tors without matching with traffic reports. The information  
 643 reported by drivers or passengers can reflect the congestion  
 644 level at that time. We investigate the number of traffic reports  
 645 under different car speed, which is captured in Fig. 9. It can  
 646 be observed that drivers report traffic jam more times when  
 647 the car speed equals to 50 to 60 and 90 to 100. It implies that  
 648 drivers not only report traffic jam when congestion happened,  
 649 sometimes drivers but also report traffic jam when there is a  
 650 slight congestion.

651 In addition, we also investigate the car speed of different  
 652 counties and cities in Taiwan, which is captured in Fig. 10.  
 653 The first four counties and cities are in north Taiwan, the last  
 654 three counties are in south Taiwan, and others are in the  
 655 middle of Taiwan. It can be observed that there are more  
 656 traffic report times in north Taiwan and fewer reports in  
 657 middle and south Taiwan. This result is mainly attributed to  
 658 the fact that drivers in north Taiwan are used to listen the  
 659 Police Broadcasting Service radio. Thus there are more traffic  
 660 reports in north than other counties and cities in Taiwan.

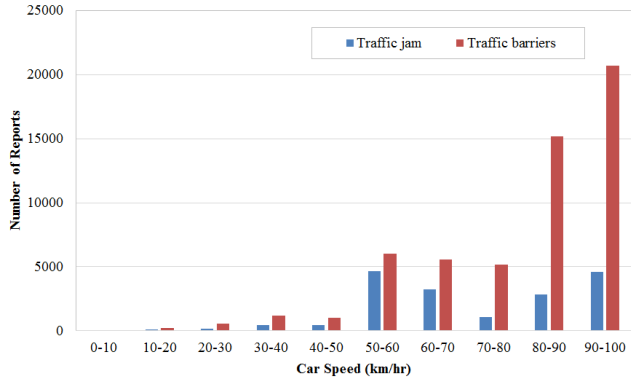


FIGURE 9. Number of reports under different car speed.

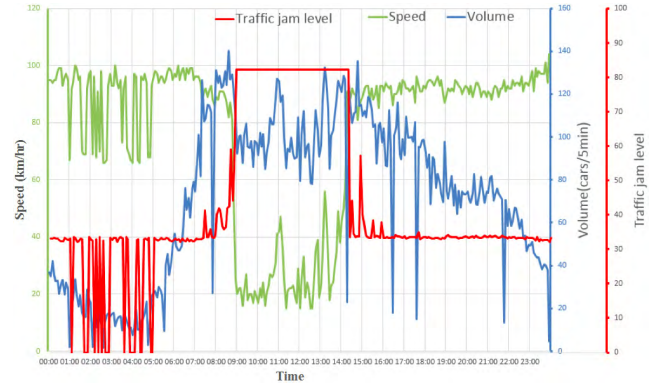


FIGURE 11. Traffic jam level in a non-rainy day.

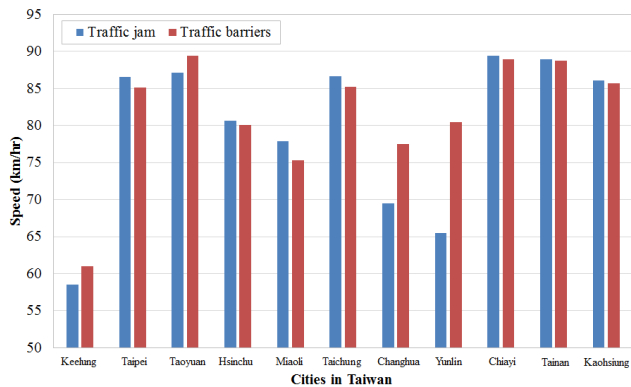


FIGURE 10. Number of reports in different counties and cities.

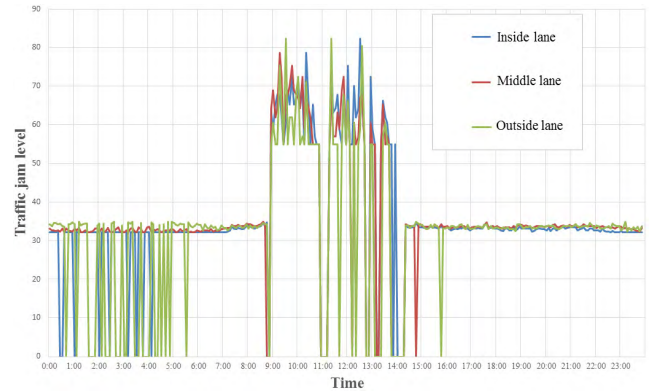


FIGURE 12. Traffic jam level of different lanes.

661 **D. EXPERIMENT 2: REAL-TIME TRAFFIC ANALYSIS**  
 662 **OF A NON-RAINY DAY**

663 In this experiment, real-time traffic in non-rainy day is  
 664 analyzed to examine whether weather influences traffic or  
 665 not. The traffic jam level of a non-rainy day is captured  
 666 in Fig. 11. The observation site is located in Hsinchu inter-  
 667 change, and the date is May 1, 2015. It can be observed that  
 668 from 9 am to 2 pm, the higher vehicle volume and lower  
 669 car speed leads to the higher traffic jam level. This result is  
 670 mainly attributed to the fact that Hsinchu is an industrial area  
 671 and from 9 am to 11 am is the commute time. In addition,  
 672 from 11 am to 2 pm is the lunchtime, Hsinchu industrial area  
 673 is lack of good restaurant so that engineers are used to drive  
 674 their cars for lunch.

675 In Fig. 12, the traffic jam level of different lanes is captured.  
 676 It can be observed that the traffic jam level from 9 am to  
 677 2 pm is higher than other time periods. The commute time  
 678 and lunchtime result in the higher congestion level. On the  
 679 other hand, it can be observed that the congestion level of  
 680 inner lane and middle lane is steadier than that of outside  
 681 lane. This is attributed to the fact that cars in outside lane  
 682 have higher probability of leaving freeway and blocked by traffic  
 683 lights.

684 **E. EXPERIMENT 3: REAL-TIME TRAFFIC ANALYSIS**  
 685 **OF A RAINY DAY**

686 In this experiment, real-time traffic of a rainy day is analyzed.  
 687 In Fig. 13, the traffic data is captured by the vehicle detectors  
 688 from Taipei to Sanchong interchange, and the observation  
 689 data is June 14, 2015. It can be observed that the rain reaches  
 690 about 3 mm every ten minutes at 3:30 pm and slightly  
 691 decreases to 1.5 mm every ten minutes at 4 pm and 5 pm.  
 692 In general, rainy days affect driver's sight and vision so that  
 693 drivers drive slowly. The car speed is obviously lower at 3 pm  
 694 than other time periods due to the heavy rain. In addition,  
 695 it can be observed that the vehicle volume at 3 pm is lower  
 696 than other time periods around 80 cars every five minutes but  
 697 the car speed is still low. This result is mainly attributed to the  
 698 fact that heavy rain leads to slower driving.

699 In Fig. 14, we also investigate the traffic jam level of the  
 700 same rainy day. The traffic jam level in Fig. 13 is evalu-  
 701 ated and obtained by the proposed fuzzy model. It can be  
 702 observed that the traffic jam level at 3 pm is higher than time,  
 703 the same phenomenon happened at 4:30 pm. Both results  
 704 mainly attributed to the fact that heavy rain results in drivers  
 705 drive slowly, thus the slower car speed leads to higher traffic  
 706 jam level. The fuzzy inference of proposed model is validated  
 707 by the results.

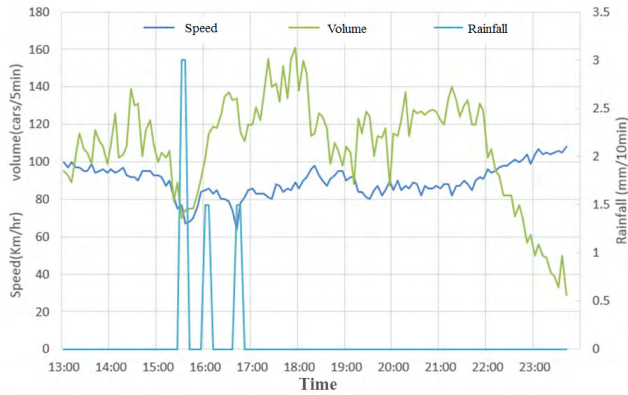


FIGURE 13. Traffic analysis of a rainy day.

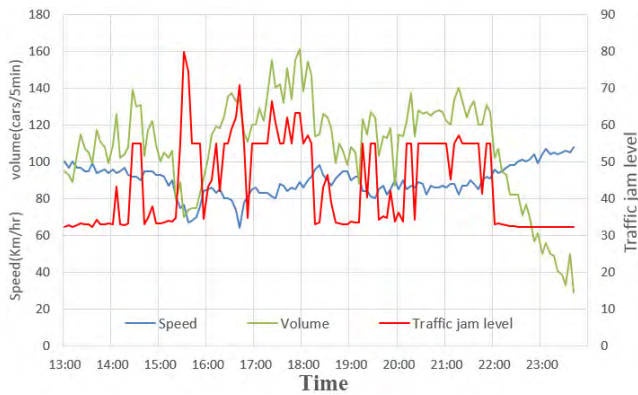


FIGURE 14. Traffic jam level of a rainy day.

**F. EXPERIMENT 4: CAR SPEED PREDICTION BY SRHTCP AND EWMA**

In this experiment, real-time car speed is predicted by the proposed SRHTCP method and the Exponentially Weighted Moving Average (EWMA) method. The real-time car speed of current, five minute and ten minute are used to be the training data. We propose SRHTCP model to explore the feature values of different car speed so that the car speed in next time period can be predicted accurately. The training data collected from May 1, 2015 to June 9, 2015. The prediction result records from June 10 to 16, 2015. The prediction results of SRHTCP and EWMA are captured in Table 8. Note that the lower MARE and MSE stand for the higher prediction accuracy. It can be observed that the proposed SRHTCP forecasts the car speed in next week accurately, which the MARE and MSE is less than 4.27% and 88.89 respectively.

On the other hand, the SRHTCP method is compared with the Exponentially Weighted Moving Average (EWMA) method. The EWMA method is also used to forecast the car speed of next time period, which is calculated by

$$EWMA_n = a * price + (1 - a) * EWMA_{n-1}, \quad (18)$$

where  $EWMA_n$  is the prediction result of next time period,  $a$  is a weighted value and price is current value, and  $EWMA_{n-1}$  is the prediction result of previous time period. Note that we

TABLE 7. Prediction result of SRHTCP and EWMA.

Date	SRHTCP MARE	EWMA MARE	SRHTCP MSE	EWMA MSE
June 10	2.96 %	4.39 %	29.99	39.01
June 11	4.14 %	4.24 %	38.10	37.80
June 12	4.15 %	4.22 %	41.03	37.25
June 13	2.32 %	3.32 %	19.07	27.30
June 14	4.27 %	4.18 %	88.89	93.25
June 15	3.87 %	6.96 %	55.59	106.44
June 16	1.46 %	3.83 %	12.06	31.56
Average	3.31 %	4.45 %	40.68	53.20

TABLE 8. Prediction result of case 1 and case 2.

Date	Case 1 MARE	Case 2 MARE
June 10	3.85%	3.02%
June 11	5.14%	5.18%
June 12	4.15%	4.66%
June 13	3.44%	2.95%
June 14	4.50%	5.52%
June 15	7.31%	6.00%
June 16	3.78%	0.39%
Average	4.60%	3.93%

set the weight  $a$  equals to 0.125, which is a common value in computer networks. The used training data in EWMA method is the same with the training data in SRHTCP method. It can be observed that the prediction accuracy of EWMA method is worse than that of the proposed SRHTCP, no matter in what kind of estimation criteria, i.e., MARE and MSE. Thereby we can say that the proposed SRHTCP method is superior to the EWMA method in terms of prediction accuracy.

**G. EXPERIMENT 5: CAR SPEED PREDICTION BY USING THE TRAINING DATA IN DIFFERENT TIME PERIODS**

In this experiment, we use the proposed SRHTCP model to predict the car speed of 30 minutes later. Two sets of training data are used and named as case 1 and case 2. In case 1, the SRHTCP model is trained by the data in current, previous 5 and 10 minutes. In case 2, training data is obtained from the data in current, previous 10, 20 and 30 minutes. The SRHTCP retrieves feature values from these two training data sets and forecasts the car speed of 30 minutes later. The training data of this experiment starts from May 1 to June 9, 2015. The prediction results of SRHTCP in case 1 and case 2 are captured in Table 9. It can be observed that the SRHTCP model yields the higher prediction accuracy in case 2 than case 1. The average MARE value of SRHTCP in case 1 and case 2 is 4.6% and 3.93% respectively. The proposed SRHTCP model improves 14.57% prediction accuracy in case 2 compared with case 1. This result is mainly attributed to the fact that the SRHTCP can find feature value more accurately in case 2 because there are more training data.

**V. CONCLUSIONS**

Unlike existing literatures used batch method to predict car traffic, we utilize Apache Storm platform to achieve real-time traffic prediction. The constructed platform integrates

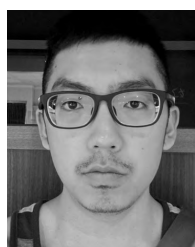
different kinds of open data, i.e., traffic data from Taiwan Area National Freeway Bureau, weather data from Central Weather Bureau, and social media data from Police Broadcasting Service. By analyzing the great quantity of traffic data, we found that there are two sorts of traffic pattern in Taiwan, i.e., weekdays from Monday to Thursday and weekend from Friday to Sunday. We analyzed social media data and found that drivers in inner line inform traffic jam report when car speed is lower than 60 km per hour. In addition, drivers in south Taiwan inform traffic jam when car speed is lower than 90 km per hour. It implies that drivers in south Taiwan have less tolerance of car speed. In experiments, we not only utilized fuzzy theory to analyze real-time traffic and congestion level but also proposed SRHTCP model to forecast the car speed of next time period. It has been shown that the SRHTCP model is superior to the EWMA method in terms of prediction accuracy no matter in MARE or MSE analysis. In the future, we will try to verify the used open data sets with t-test method.

## REFERENCES

- [1] C.-W. Tsai, C.-F. Lai, H.-C. Chao, and A. V. Vasilakos, "Big data analytics: A survey," *J. Big Data*, vol. 2, no. 21, pp. 1–32, Oct. 2015.
- [2] D. Laney, "3D data management: Controlling data volume, velocity and variety," META Group, Stamford, CT, USA, Tech. Rep. 949, Feb. 2001.
- [3] P. Chandarana and M. Vijayalakshmi, "Big data analytics frameworks," in *Proc. CSCITA*, Mumbai, India, 2014, pp. 430–434.
- [4] OICA. (May 1, 2018). *Organisation Internationale des Constructeurs d'Automobiles*. [Online]. Available: <http://www.oica.net/>
- [5] (May 1, 2018). *Ministry of Transportation and Communications R.O.C.* [Online]. Available: <http://www.motc.gov.tw/ch/index.jsp>
- [6] *Apache Storm*. (Dec. 24, 2015). [Online]. Available: <http://storm.apache.org/>
- [7] (Jul. 19, 2017). *Apache Storm Rationale Version 1.1.0*. [Online]. Available: <http://storm.apache.org/releases/current/Rationale.html>
- [8] (Jul. 19, 2017). *DATA.GOV.TW*. [Online]. Available: <http://data.gov.tw/>
- [9] (Jul. 19, 2017). *Tisvcloud.Freeway.Gov.Tw*. [Online]. Available: <http://tisvcloud.freeway.gov.tw/>
- [10] (Jul. 19, 2017). *OPEN DATA*. [Online]. Available: <http://opendata.cwb.gov.tw/index>
- [11] (Jul. 19, 2017). *Police Broadcasting Service of Taiwan*. [Online]. Available: <http://rtr.pbs.gov.tw/pbsmg/RoadAll.html>
- [12] B. D. Greenshields, "A study of traffic capacity," in *Proc. Highway Res. Board*, vol. 14, 1935, pp. 448–477.
- [13] A. Nakayama et al., "Detailed data of traffic jam experiment," in *Traffic & Granular Flow*. Berlin, Germany: Springer, 2009, pp. 389–394.
- [14] University of Bristol. (Jul. 19, 2017). *Phantom Traffic Jams*. [Online]. Available: <http://www.bristol.ac.uk/news/2010/6948.html>
- [15] L. A. Zadeh, "Fuzzy sets," *Inf. Control*, vol. 8, no. 3, pp. 338–353, Jun. 1965.
- [16] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Proc. ACM 5th Annu. Workshop Comput. Learn. Theory (COLT)*, Pittsburgh, PA, USA, Jul. 1992, pp. 144–152.
- [17] F.-H. Tseng, X. Chen, L.-D. Chou, H.-C. Chao, and S. Chen, "Support vector machine approach for virtual machine migration in cloud data center," *Multimedia Tools Appl.*, vol. 74, no. 10, pp. 3419–3440, May 2015.
- [18] C.-Y. Chen, K.-D. Chang, and H.-C. Chao, "Transaction-pattern-based anomaly detection algorithm for IP multimedia subsystem," *IEEE Trans. Inf. Forensics Security*, vol. 6, no. 1, pp. 152–161, Mar. 2011.
- [19] W. Wu and H. Zhou, "Data-driven diagnosis of cervical cancer with support vector machine-based approaches," *IEEE Access*, vol. 5, pp. 25189–25195, Oct. 2017.
- [20] S. U. Jan, Y. D. Lee, J. Shin, and I. Koo, "Sensor fault classification based on support vector machine and statistical time-domain features," *IEEE Access*, vol. 5, pp. 8682–8690, May 2017.
- [21] M. Milojevic and V. Rakocevic, "Distributed vehicular traffic congestion detection algorithm for urban environments," in *Proc. IEEE Veh. Netw. Conf. (VNC)*, Boston, MA, USA, Dec. 2013, pp. 1–6.
- [22] Y. Cheng, M. Zhang, and D. Yang, "Automatic incident detection for urban expressways based on segment traffic flow density," *J. Intell. Transp. Syst.*, vol. 19, no. 2, pp. 205–213, Jan. 2015.
- [23] H. Ji, A. Xu, X. Sui, and L. Li, "The applied research of Kalman in the dynamic travel time prediction," in *Proc. 18th Int. Conf. Geoinformat.*, Beijing, China, Jun. 2010, pp. 1–5.
- [24] H. Feng, C. Liu, Y. Shu, and O. W. W. Yang, "Location prediction of vehicles in VANETs using a Kalman filter," *Wireless Pers. Commun.*, vol. 80, no. 2, pp. 543–559, Jan. 2015.
- [25] H.-F. Wang, "A freeway automatic incident detection algorithm using grey prediction," M.S. thesis, Dept. Comput. Sci. Inform. Eng., Nat. Central Univ., Taoyuan, Taiwan, 2003.
- [26] K.-W. Kuo, "A hybrid travel-time prediction approach based on macroscopic and microscopic methodologies," M.S. thesis, Dept. Comput. Sci. Inform. Eng., Nat. Central Univ., Taoyuan, Taiwan, 2013.
- [27] D. C. Li, L.-D. Chou, L.-M. Tseng, Y.-M. Chen, and K.-W. Kuo, "A bipolar traffic density awareness routing protocol for vehicular ad hoc networks," *Mobile Inf. Syst.*, vol. 2015, Sep. 2015, Art. no. 401518, doi: [10.1155/2015/401518](https://doi.org/10.1155/2015/401518).
- [28] W. Yang, X. Liu, L. Zhang, and L. T. Yang, "Big data real-time processing based on storm," in *Proc. 12th IEEE Int. Conf. Trust, Secur. Privacy Comput. Commun. (TrustCom)*, Melbourne, VIC, Australia, Jul. 2013, pp. 1784–1787.
- [29] T. Chardonens, P. Cudre-Mauroux, M. Grund, and B. Perroud, "Big data analytics on high velocity streams: A case study," in *Proc. IEEE Int. Conf. Big Data*, Silicon Valley, CA, USA, Oct. 2013, pp. 784–787.
- [30] *Twitter*. (Dec. 24, 2015). [Online]. Available: <https://twitter.com/>
- [31] *Bitly*. (Dec. 24, 2015). [Online]. Available: <https://bitly.com/>
- [32] A. Bifet and G. De F. Morales, "Big data stream learning with SAMOA," in *Proc. IEEE Int. Conf. Data Mining Workshop (ICDMW)*, Shenzhen, China, Dec. 2014, pp. 1199–1202.
- [33] T. Chai and R. R. Draxler, "Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against avoiding RMSE in the literature," *Geosci. Model Develop.*, vol. 7, no. 3, pp. 1247–1250, Jun. 2014.

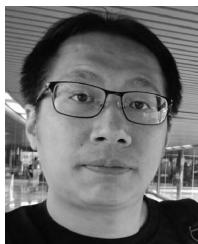


**FAN-HSUN TSENG** (S'12–M'18) received the Ph.D. degree in computer science and information engineering from National Central University, Taiwan, in 2016. He is currently an Assistant Professor with the Department of Technology Application and Human Resource Development, National Taiwan Normal University, Taipei, Taiwan. His research interests include mobile edge computing, 5G mobile networks, and artificial intelligence. He received the 2018 Young Scholar Fellowship Program from the Ministry of Science and Technology, Taiwan, for his dedication to research of engineering and technologies. He has served as an Associate Editor-in-Chief of the *Journal of Computers* and an Associate Editor of *Human-Centric Computing and Information Sciences*.



**JEN-HAO HSUEH** received the M.S. degree in computer science and information engineering from National Central University, Taoyuan, Taiwan, in 2015. His research interests include big data, vehicle ad-hoc network, fuzzy theory, and support vector machine.

886  
887  
888  
889  
890  
891  
892  
893  
894



**CHIA-WEI TSENG** (S'17) received the M.S. degree in computer science and information engineering from National Dong Hwa University, Hualien, Taiwan, in 2006. He is currently pursuing the Ph.D. degree with the Department of Computer Science and Information Engineering, National Central University, Taiwan. His research interests include software-defined networking, vehicle ad-hoc network, and IPv6 protocol.



**HAN-CHIEH CHAO** (SM'04) received the M.S. and Ph.D. degrees in electrical engineering from Purdue University, West Lafayette, IN, USA, in 1989 and 1993, respectively. He is currently a Professor with the Department of Electrical Engineering, National Dong Hwa University, where he also serves as the President. He is also with the Department of Computer Science and Information Engineering and the Department of Electronic Engineering, National Ilan University, Taiwan, the College of Mathematics and Computer Science, Wuhan Polytechnic University, Wuhan, China, and the Fujian University of Technology, Fuzhou, China. He is a fellow of IET (IEE) and a Chartered Fellow of the British Computer Society. He serves as the Editor-in-Chief for the Institution of Engineering and Technology Networks, the *Journal of Internet Technology*, the *International Journal of Internet Protocol Technology*, and the *International Journal of Ad Hoc and Ubiquitous Computing*.

904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920

895  
896  
897  
898  
899  
900  
901  
902  
903



**YAO-TSUNG YANG** (S'10) received the M.S. degree in computer science and information engineering from National Central University, Taoyuan, Taiwan, in 2008, where he is currently pursuing the Ph.D. degree with the Department of Computer Science and Information Engineering. His research interests include software-defined networking, vehicle ad-hoc network, and network management.



**LI-DER CHOU** (M'95) received the M.S. and Ph.D. degrees in electronic engineering from the National Taiwan University of Science and Technology, Taiwan, in 1991 and 1995, respectively. He is currently a Distinguished Professor with the Department of Computer Science and Information Engineering and the Director of Computer Center, National Central University, Taiwan. He was a Director of the Board of Taiwan Network Information Center. He was also the Deputy Director General of the National Center for High-performance Computing, Taiwan, from 2013 to 2016. He is the holder of five U.S. and 16 Taiwan invention patents. His research interests include SDN/NFV/SFC, vehicular networks, network management, broadband wireless networks, and Internet services. He has published over 200 papers in these areas. He was a recipient of seven best paper awards and four excellent paper awards from international and domestic conferences. He was also a recipient of two Gold Medal Awards and four Silver Medal Awards in international invention shows held in Geneva, Moscow, London, and Taipei.

921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940

• • •