

Multi-Source Data Fusion for Urban Traffic State Estimation: A Case Study of New York City

Jingqin Gao, Kaan Ozbay, Abdullah Kurkcu

This is the author's version of a work that has been accepted for presentation at the Transportation Research Board's 98th Annual Meeting, Washington D.C., 2019. The final version is available here: <http://amonline.trb.org/68387-trb-1.4353651/t0005-1.4505752/1214-1.4506969/19-04253-1.4494307/19-04253-1.4506990>

1 **MULTI-SOURCE DATA FUSION FOR URBAN TRAFFIC STATE ESTIMATION: A**
2 **CASE STUDY OF NEW YORK CITY**

3
4 **Jingqin Gao, M.Sc. (Corresponding author)**

5 Graduate Research Assistant, C2SMART Center,
6 Department of Civil and Urban Engineering,
7 Tandon School of Engineering, New York University
8 Six MetroTech Center, 4th Floor, Brooklyn, NY 11201, USA
9 Tel: (646) 717-3652
10 E-mail: jingqin.gao@nyu.edu

11
12 **Kaan Ozbay, Ph.D.**

13 Professor & Director, C2SMART Center (A Tier 1 USDOT UTC),
14 Department of Civil and Urban Engineering &
15 Center for Urban Science and Progress (CUSP),
16 Tandon School of Engineering, New York University (NYU)
17 Fifteen MetroTech Center, 6th floor, Brooklyn, NY 11201, USA
18 Tel: 1-(646) 997-3691; E-mail: kaan.ozbay@nyu.edu

19
20 **Abdullah Kurkcu, Ph.D.**

21 Research Associate, C2SMART Center (A Tier 1 USDOT UTC),
22 Department of Civil and Urban Engineering &
23 Center for Urban Science and Progress (CUSP),
24 Tandon School of Engineering, New York University (NYU)
25 370 Jay Street, 12th Floor, Brooklyn, NY 11201, USA
26 Tel: 1-(646)-997-0538
27 Email: ak4728@nyu.edu

28
29
30
31
32 Word count: 6,001 words text + 2 tables x 250 words (each) = 6,501 words
33 Resubmission Date: November 15th, 2018

34
35
36
37
38
39
40
41
42 Paper resubmitted for Presentation in the
43 *Transportation Research Board's 98th Annual Meeting*, Washington, D.C., 2019

1 ABSTRACT

2 Data fusion techniques are often used to enhance traffic state estimation. The objective of this
3 paper is to evaluate and validate the applicability of three different data fusion techniques on a
4 non-trivial urban transportation network. Simple weighted, machine learning and evidence theory
5 based approaches are applied to generate estimates of traffic states. Multiple data sources including
6 information collected from electronic toll collection tag readers, Global Positioning System-
7 equipped probe vehicles, and crowdsourcing map applications are utilized as primary data sources
8 in the paper. A case study is provided to illustrate an application of the proposed data fusion
9 techniques with data extracted in 2017 for a period of two weeks. Ground truth information
10 collected from real-time camera feeds are used for validation purposes. The evidence theory based
11 approach considering temporal evidence reliability is found to outperform other methods of cross
12 validation in model accuracy. In addition, this study proves that the information extracted from
13 web-based map services using “virtual sensors” can be an excellent supplementary data source for
14 current travel time monitoring systems at no additional cost for the installation and maintenance
15 of traditional infrastructure-based sensors.

16

17 *Keywords:* Data fusion, Traffic state estimation, Virtual sensors, Evidence theory

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

1 INTRODUCTION AND MOTIVATION

2 With the rapid growth of Intelligent Transportation Systems (ITS), sensing technologies and the
3 variety and volume of data collected, providing an enhanced and accurate interpretation of
4 monitored traffic state is becoming a major challenge for public agencies and private companies
5 (1). In the era of big data, researchers are faced with a diversity of information from different data
6 sources that have different representations, scales, and density. Recent technology developments
7 also make it possible to use different acquisition methods to obtain various heterogeneous traffic
8 data (2). Such data includes, but is not limited to information from microwave sensors, real-time
9 cameras, Wi-Fi/Bluetooth tracking devices (3, 4), and crowdsourcing mobile applications. For
10 example, Morgul, Ender, *et al.* (5, 6) developed a “virtual sensor” methodology using open traffic
11 data sources from web-based map providers to measure travel time from probe vehicles.

12 However, the traditional way of using a single data source without cross-validation cannot
13 achieve good performance in data mining (7) because the single data source may not be a good
14 representative of the whole population. Simply treating features of different datasets equally are
15 usually not useful under many circumstances (7). Advanced techniques that can fuse the
16 knowledge from different but potentially connected datasets in an intelligent way are needed to
17 provide a unified and global view of our transportation system. This can also enable a reliable
18 interpretation of the monitored traffic situation.

19 As a crucial component in traffic operations and planning, Traffic State Estimation (TSE)
20 usually plays an important role in day-to-day traffic monitoring and evaluation. TSE refers to “the
21 process of inference of traffic state variables, such as traffic flow or speed on road segments using
22 partially observed and noisy traffic data” (2). For example, link travel times collected from
23 Electronic Toll Collection (ETC) tag readers are used to estimate traffic states and reveal prevailing
24 congestion levels in New York City (NYC). The information is fed into a real-time adaptive signal
25 control system, named “Midtown in Motion” (MIM) (8), to manage traffic congestion. However,
26 traffic states are not observed at all locations all the time, and single-source measurements are
27 usually noisy (2).

28 Seizing the potential of ‘big data’ and discovering smart ways to integrate them would
29 provide new insights in estimating dynamic traffic states of an urban road network, but how to
30 fully make use of multi-source data to produce a better inference remains a major challenge. The
31 objectives of our study include: 1) Assess weaknesses and strengths of different data fusion
32 approaches using reliable ground truth data over a non-trivial traffic network, 2) conduct real-time
33 traffic state estimation by using this fused traffic information from multiple sources to demonstrate
34 the usefulness of a robust data fusion approach, 3) validate the virtual sensor approach with
35 frequently used ITS data collection technology and evaluate whether it can be used as a good
36 supplementary data source. The focus of this research is concentrated in testing the validity of
37 three data-driven data fusion techniques that rely on historical-data and statistical or machine
38 learning methods and their applicability on a complex transportation network in Midtown NYC.
39 Real-time Automatic vehicle location (AVL), travel time, and speed information collected through
40 in-vehicle GPS devices or ETC tag readers from three different data sources, namely, the MIM
41 system, virtual sensors, and 1,336 buses (9), are extracted and fused. The fused data are compared
42 with “ground truth” information collected from video cameras.

43
44
45

1 LITERATURE REVIEW

2 Data Fusion Techniques

3 Generally, data fusion methodologies can be split into three categories: Statistical approaches,
4 probabilistic approaches, and artificial intelligence approaches (10). Statistical approaches include
5 weighted combination, multivariate statistical analysis and so on (11). Among them, the arithmetic
6 mean approach is the simplest method for information combination (10). Probabilistic approaches
7 such as Bayesian approach with Bayesian network, maximum likelihood and Kalman filter based
8 methods (12, 13), evidence theory (13-16) are widely used for the multi-sensor data fusion. For
9 example, Kong *et. al* utilized federated Kalman filter and evidence theory to estimate urban traffic
10 states. They assigned dynamic evidence reliability and fused link mean speed from both
11 underground loop detectors and GPS-equipped probe vehicles. Their results showed that the
12 proposed approach could well be used in urban traffic applications on a large scale. However, their
13 approach did not consider the temporal and spatial interactions and dependencies of the data
14 between different links.

15 Artificial intelligence has also become a very popular approach in recent years. Machine
16 learning techniques such as neural networks, clustering, and artificial cognition are applied in
17 various transportation applications (7, 17-20). For example, Xu *et. al* applied K-means clustering
18 analysis to classify freeway traffic flow into five different states. They built a conditional logistic
19 regression model to study the relationship between traffic states and crash risks. Each traffic state
20 was compared to identify the underlying traffic flow characteristics that made certain states more
21 hazardous than others. Likewise, Bartin *et. al* conducted several studies related to travel time
22 estimation (21, 22). One contribution of their work was the use of k-means clustering of sample
23 space-time vehicle trajectory data. Its purpose is to find the optimal roadway segment
24 configuration to estimate travel times with minimum errors. Moreover, this approach was able to
25 determine the optimum number of segments by using the percent-gain of clustering technique as
26 well.

27 Furthermore, models inspired by other disciplines such as the Mer-Gesh fusion framework
28 proposed by Zhang *et. al* (23) are also new approaches in data fusion. Their model was based on
29 the idea of optimizing the working process of transmission gears to fuse data that were collected
30 from multiple sensors with different location and features. This data fusion system consisted of
31 multiple measurements for the input data including data source effectiveness, consistency, and
32 variation measurement. The key advantage of this method is its ability to accept any types of
33 sensors including taxi trajectories, bus trajectories, fixed detectors, and cameras.

34 In short, data fusion techniques appear promising in the context of the estimating traffic
35 states. Nevertheless, several key questions related to data fusion remain to be addressed including
36 how to assess the input data reliability and credibility of fusion system (24) and how to fuse them
37 smartly when the input data do not have same representations, scales, or density. There is a need
38 to investigate a different data fusion approach, especially when it comes to a highly complex and
39 congested traffic environment.

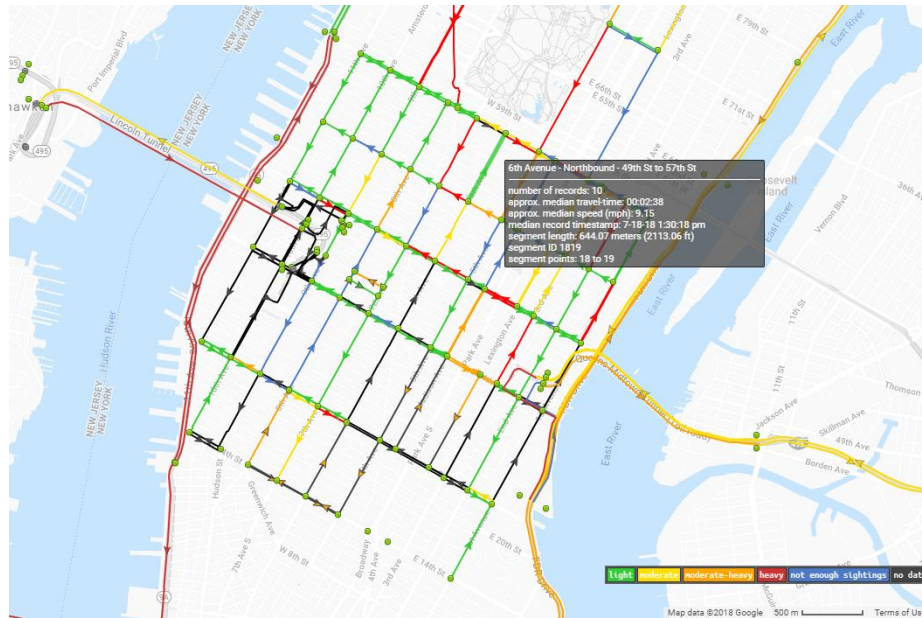
40

41 DATA PREPARATION

42 Midtown in Motion (MIM)

43 Midtown in Motion, as an integral part of the NYC Department of Transportation (NYCDOT)
44 enhanced mobility strategy, has improved travel times on the avenues in Midtown NYC area by
45 10% since its first implementation in 2011 (25). The MIM system collects traffic flow and
46 occupancy from microwave sensors, travel times from ETC tag readers and uses video cameras

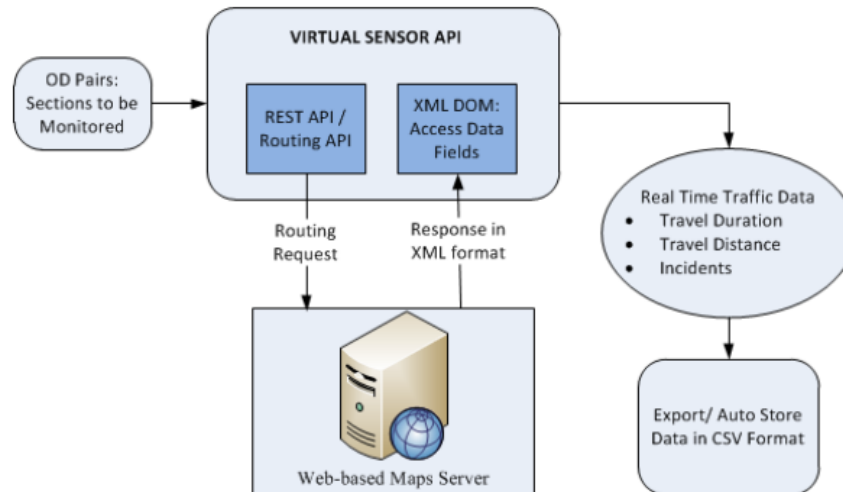
1 for verification and monitoring the field conditions. By using a unique hierarchical two-level
 2 control, traffic engineers at the Traffic Management Center (TMC) are able to quickly respond to
 3 congestion issues and smooth the flow of traffic remotely using adaptive signal control (25).
 4 Typically, travel times in segments are recorded and aggregated about every two to three minutes.
 5 The advantage of the MIM system is that travel times collected from ETC tag readers work well
 6 under congested conditions. However, they can only provide traffic states among long segments
 7 (typically 8-block segments on north-south avenues). FIGURE 1 illustrates current MIM
 8 implementation in Midtown Manhattan.
 9



10
 11 **FIGURE 1 Midtown in Motion real time traffic speed (<http://nyctmc.org/>).**
 12

13 Virtual Sensors

14 “Virtual sensor” methodology was first purposed by Morgul *et al.* (5, 6) in 2013. Open traffic data
 15 sources from web-based map providers namely, Bing Maps™ and MapQuest™ are used to
 16 measure travel time from probe vehicles that are already in the traffic stream to estimate traffic
 17 conditions in real-time across large networks. FIGURE 2 shows the framework of the virtual
 18 sensor concept. After selecting geographical coordinates of origin-destination (OD) pairs for the
 19 road segments to be monitored, Bing Maps’ REST API and MapQuest Open Data Map API
 20 services are utilized for extracting real-time traffic-based routing information. Routing requests
 21 for each of the defined OD pairs are sent automatically using the developed Python code and the
 22 responses are received from the server in Extensible Markup Language (XML) format every five
 23 minutes. In this study, travel times in segments are recorded and aggregated about every five
 24 minutes. This virtual sensor methodology comes with almost no additional cost while the quality
 25 of obtained data is proved to be quite satisfactory compared with traditional sensors such as loop
 26 detectors on highways (5). One of the goals of this study aims to conduct a further validation of
 27 Virtual Sensor methodology by comparing it with travel time information collected from
 28 NYCDOT’s MIM system.



1
2 **FIGURE 2 Virtual Sensor framework (5).**

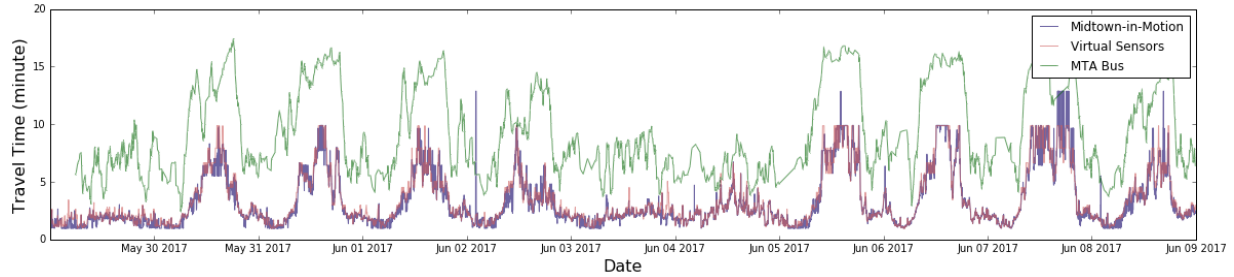
3
4 **MTA Bus Time**

5 Data from GPS-integrated cellphones or in-vehicle devices becomes a new cost-effective way of
6 monitoring transportation system. Numerous technologies was developed to identify vehicle
7 locations from vehicle-embedded smart devices using roadside sensors such as Bluetooth (5).
8 Many public agencies installed in-vehicle GPS devices to collect information such as spot speed
9 or AVL as well. In NYC, a service called MTA Bus Time by the Metropolitan Transportation
10 Authority (MTA), has been providing massive GPS based information from MTA buses since
11 2011 (9). Rich data such as vehicle location, expected arrival/departure time, next stop names of
12 all buses are collected at 30 seconds intervals.

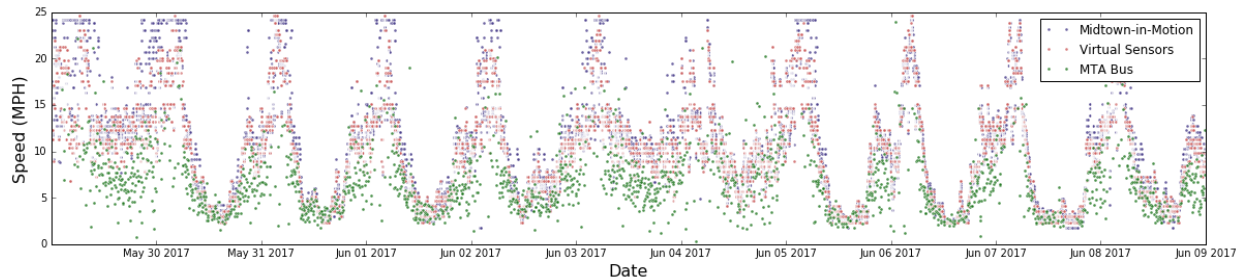
13
14 **Video Cameras**

15 Traffic video camera recordings also provide a vast volume of information including traffic
16 volume, travel time, driver behavior, incidents, occupancy and detailed traffic operations. It often
17 serves as a tool for traffic monitoring and congestion/incident verification. Open-source video
18 feeds from NYCDOT's closed circuit television (CCTV) cameras (26) installed on major arteries
19 allow road users to view real-time traffic movements from frequently updated still images
20 (typically every 3-5 seconds).

21
22 To evaluate different data fusion technologies, two-week long data was manually collected for 6th
23 Avenue between 49th Street and 57th Streets from May 29th to June 9th in 2017. Travel time and
24 speed information for this 8-block segment were extracted from the MIM system and virtual
25 sensors. MIM data contains 6,105 records and virtual sensor has 7,920 records. For MTA Bus
26 Time, mean travel time and mean speed for each street block among 6th Avenue between 49th Street
27 and 57th Streets were estimated through haversine estimation (great-circle distance) between two
28 geographic coordination. The raw data has 42,768 data points that record the locations of each
29 individual bus before aggregating them at a 5-minute time interval. Travel time and speed for each
30 individual vehicle on a reference street block on 6th Avenue from 56th Street to 57th Street are
31 recorded by manual processing from CCTV camera videos for AM peak hour (8:00AM – 9:00AM)
32 from May 30th to June 8th in 2017 and are used as “ground truth” information in this study. FIGURE
33 3 shows travel time and traffic speed information from the MIM system, virtual sensor, and MTA
34 Buses for the same road segment (6th Avenue between 49th Street and 57th Streets).



(a) Travel time



(a) Traffic Speed

FIGURE 3 Travel time and traffic speed information extracted from MIM, virtual sensor and MTA Bus Time during May 29th, to Jun 9th, 2017.

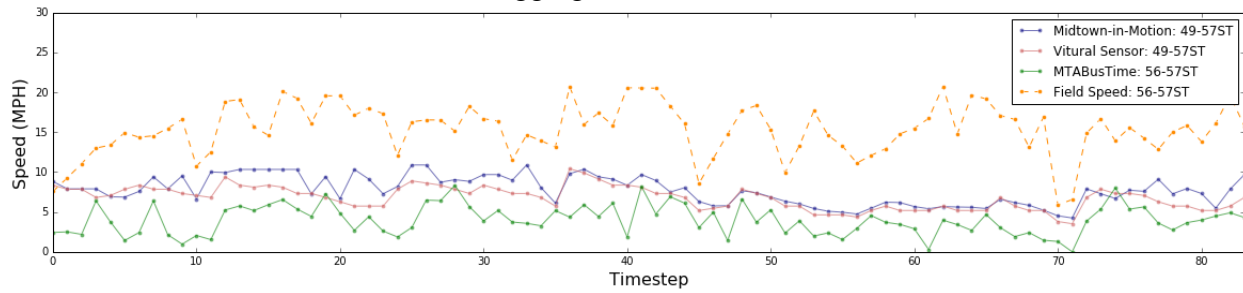
PROBLEM STATEMENT

As mentioned earlier, traffic state estimation is a crucial component of day-to-day traffic operations, monitoring, and evaluation. The recent developments in ITS and emerging technologies have made two types of travel time/ traffic speed data available to the public. The first type of data is relatively reliable yet has low-resolution. For example, the MIM system has been implemented in the field for 7 years but is only available for limited number of road segments. Each road segment usually contains multiple street blocks. The second type of data which comes from GPS-equipped vehicles like buses that has a wider coverage and higher-resolution and it can be aggregated at a street block level. However, this type of data might be noisy. For instance, bus travel time data obtained from GPS-equipped buses can be biased at street blocks that have bus stations. The key research question is that whether it is possible to reasonably fuse these different data sources to accurately estimate traffic states at the street block level. Another important research question is to test whether crowd sourcing data such as the one obtained by the “virtual sensors” approach proposed by the co-authors of this paper are accurate enough to be served as an additional data source for an existing system such as MIM. The next sections attempt to provide answers to these research questions.

METHODOLOGY

Three data fusion techniques: 1) simple weighted, 2) Random Forest, and 3) evidence theory with improved reliability, are investigated in this study. To compare with available ground truth data collected from video feeds, the three data fusion techniques are applied for the morning peak hours. The data for the implementation of these techniques come from seven weekdays from May 30th to June 8th, 2017 on one reference street block (56th Street to 57th Street on 6th Avenue). Traffic speed is used as the performance measure. FIGURE 4 illustrates mean traffic speed during AM peak hour from MIM and virtual sensor on 6th Avenue from 49th Street to 57th Street, and MTA Buses

1 and field (Ground truth) for the reference street block on 6th Avenue from 56th Street to 57th Street.
 2 The raw data from each data source is aggregated at a 5-minute interval.



3
 4 **FIGURE 4 Mean traffic speed between 8AM to 9AM from MIM, virtual sensor, MTA Buses**
 5 **and field (ground truth) between May 30th and June 8th, 2017 (Time step is 5 minutes).**

7 Traffic State Estimation

8 TSE is an important component of traffic control and operation and is often used in emergency
 9 response or congestion mitigation. For example, New York City uses estimated traffic states from
 10 the MIM system to determine its adaptive signal strategies that helps relieve congestion problems
 11 in its central core business area. However, this type of system, such as MIM, can only measure
 12 traffic state variables at certain locations or at low resolutions (i.e. MIM typically has a
 13 measurement at every 8-block). When it scales down to street block level, such system will need
 14 further assistants from additional data sources like occupancy information from microwave sensors
 15 or CCTV cameras (25). Unfortunately, these devices are not installed everywhere due to
 16 technological or financial limitations. Our study proposes to process the inference of traffic speed
 17 information from the three data sources with different resolutions to obtain TSE at street block
 18 level. The proposed approaches are tested for the reference street block. Five traffic states $\{S_1, S_2,$
 19 $S_3, S_4, S_5\}$ are assumed using partition method based on a previous study (13). S_1 to S_5 means
 20 “very congested” ($<1/6$ maximum speed), “congested” ($1/6-1/3$ maximum speed), “moderate”
 21 ($1/3-1/2$ maximum speed), “smooth” ($1/2-3/4$ maximum speed), and “very smooth” ($>3/4$
 22 maximum speed).

24 Simple Weighted method and Random Forest Classifier

25 Simple weighted method assigns static weights to each data source. The weights are generated
 26 based on the historical data. As a machine learning approach, Random Forest is utilized as a
 27 supervised nonlinear classification algorithm in this study. We first label our ground truth mean
 28 traffic speed data for the reference link into five traffic states S_1 to S_5 . A 70-30 split is used to
 29 generate the training and test data in time order. This train-test split applies for all three methods.
 30 A random forest estimator constructs various decision trees on many sub-samples of the dataset.
 31 Each decision tree starts with a root node and every non-terminal node has two child nodes. It
 32 applies a binary test to every node of the input data and propagates the node to either of the child
 33 nodes depending on the test outcome (27). Mean decrease Gini impurity (27) is used as the criteria
 34 to split the node. Next, the trained model is applied to the test set. Each decision tree of the test set
 35 obtains an estimated traffic state and the final traffic state is voted by multiple trees. This allows
 36 the predictive accuracy to be improved, and over-fitting problem to be controlled (28). Python
 37 Scikit-learn package (28) is used with 1000 decision tree estimators, a minimum two samples to
 38 split an internal node and bootstrap strategies.

39

1 Evidence Theory with Improved Reliability

2 Information-fusion technology has been employed by several researchers in the context of the
3 traffic state estimation problem in the last two decades (13, 29, 30). It is a powerful tool employed
4 to fuse on-line or off-line data from multiple sources to obtain more accurate and complete traffic
5 state estimation than just using a single source. Based on the literature review, Dempster–Shafer
6 (DS) theory or evidence theory is used because of its advantages in dealing with incompleteness
7 and inaccuracy of traffic data. The initial work introducing DS theory is found in Dempster
8 (14) and Shafer (15). The DS theory is a generalization of Bayesian inference and is a great tool
9 for probabilistic reasoning based on a formal calculus for combining evidence (31). Let's denote
10 the DS environment as follows:

$$11 \theta: \theta = \{\theta_1, \theta_2, \theta_3, \dots, \theta_n\} \quad (1)$$

12 θ is a set of possible conclusions in which all elements are assumed to be mutually exclusive. θ
13 is exhaustive. Each subset of θ can be interpreted as a possible answer to a question and only one
14 answer is correct. The “frame of discernment” or Power set of θ , $P(\theta)$, has all subsets in θ and is
15 denoted by:

$$16 P(\theta) = \{A | A \subseteq \theta\} \quad (2)$$

17 In DS Theory, the Degree of Belief in evidence is analogous to the mass of a physical object. The
18 Basic Probability Assignment (BPA) is the evidence that measures the amount of mass and mass
19 is a function that maps each element of the Powerset into a real number in the [0,1] interval:
20 $m: P(\theta) \rightarrow [0,1]$ (14). For example, mass function of \mathbf{A} , $m(\mathbf{A})$ is the proportion of all evidence
21 that supports this element of the power set. The belief function (*bel*) is the belief measure that
22 represents the sum of masses in all subsets of \mathbf{A} and the plausibility function (*pls*) represents the
23 sum of masses committed to those subsets that do not discredit \mathbf{A} (13). They are defined as:

$$24 \begin{cases} bel(\mathbf{A}) = \sum_{\emptyset \neq B \subseteq A} m(\mathbf{B}) \quad \forall A \subseteq \theta \\ pls(\mathbf{A}) = \sum_{B \cap A \neq \emptyset} m(\mathbf{B}) \quad \forall A \subseteq \theta \end{cases} \quad (3)$$

25 The range [*Bel*, *Pls*] is called range of belief or evidential interval. In general, $0 \leq Bel \leq Pls \leq 1$.
26 Multiple evidences can be fused using Dempster's combination rules that is expressed in terms of
27 orthogonal sums of the masses of a set and all its subsets (13):

$$28 m(\mathbf{C}) = \begin{cases} 0, & \mathbf{B} \cap \mathbf{A} = \emptyset \\ \frac{1}{1-K} \sum_{\mathbf{B} \cap \mathbf{A} = \mathbf{C}, \forall \mathbf{A}, \mathbf{B} \subseteq \theta} m_i(\mathbf{A}) \cdot m_j(\mathbf{B}), & \mathbf{B} \cap \mathbf{A} \neq \emptyset \end{cases} \quad (4)$$

$$29 K = \sum_{\mathbf{B} \cap \mathbf{A} = \emptyset, \forall \mathbf{A}, \mathbf{B} \subseteq \theta} m_i(\mathbf{A}) \cdot m_j(\mathbf{B}) \quad (5)$$

30 where the term K is called the conflict factor between two evidences, which reflects the conflict
31 degree of them.

32 Moreover, previous studies in transportation and internet of thing have emphasized the
33 importance of “Evidence Reliability” in terms of improving fusion results (13, 24, 32). When the
34 information extracted from the evidence is missing or not totally reliable, a coefficient α can be
35 used to discount the belief. When $\alpha = 0$, it indicates that the information is completely not reliable;
36 on the contrary, $\alpha = 1$ denotes that the evidence is absolutely reliable. In this study, two factors are
37 considered for constructing a reliability matrix. The first one is the importance of evidence. Since
38 MIM and virtual sensors both provide segment mean travel time and speed and do not necessarily
39 represent the mean travel time and speed on a single street block segment, they are assigned with
40 a lower importance compared with MTA Bus Time data. Moreover, the accuracy of the aggregated
41 mean travel time and speed from GPS-equipped vehicles depends on the sample size as well, the
42 number of data samples before aggregation at each time step is used as an approximation for
43 reliability.

1 One of the key challenges in applying D-S theory is how to obtain the mass function (33).
 2 Here we construct the mass function using the following negative exponential function proposed
 3 by Denoeux (34):

$$4 \quad m_i(\mathbf{A}) = \exp(-\alpha_i D_i^\beta) \quad (6)$$

5 Where D_i is the distance between the data detected by the evidence and the prototype of the
 6 evidence. α and β are coefficients and the prototype are obtained by historical data.

7 If we let $v_{i,j'}$ denote the time-dependent reliability matrix, and $m_i(\mathbf{A}_{i,t})$ ($i = 1, 2, \dots, M$) represent
 8 the BPA extracted from the evidence i at time t , then the modified time-dependent BPA can be
 9 rewritten as:

$$10 \quad \begin{cases} m'_i(\mathbf{A}_{i,t}) = v_{i,j'} \cdot m_i(\mathbf{A}_{i,t}) \\ m'_i(\emptyset) = 1 - \sum_{\mathbf{A}_{i,t} \subset \emptyset} v_{i,j'} \cdot m_i(\mathbf{A}_{i,t}), \quad \mathbf{A}_{i,t} \subset \emptyset \end{cases} \quad (7)$$

11 Once all of the BPAs from different evidence for each timestep (every 5 minutes) are obtained,
 12 they can be fused using the combining method from DS theory. The integrated result $m'(\mathbf{C}_t)$ of
 13 the fusion system at time t is (13):

$$14 \quad m'(\mathbf{C}_t) = m'_1(\mathbf{B}_{1,t}) \oplus m'_2(\mathbf{B}_{2,t}) \oplus \dots \oplus m'_M(\mathbf{B}_{M,t}) \\ 15 \quad = \frac{\sum_{\cap_{i=1}^M \mathbf{B}_{i,t} = \mathbf{C}_t} [\prod_{i=1}^M m'_i(\mathbf{B}_{i,t})]}{1 - \sum_{\cap_{i=1}^M \mathbf{B}_{i,t} = \emptyset} [\prod_{i=1}^M m'_i(\mathbf{B}_{i,t})]} \quad (8)$$

16 Three evidences E_1 , E_2 and E_3 denote the information provided by MIM, virtual sensors and MTA
 17 Bus Time. Maximum belief is used as the decision strategy of the output state.

18

19 **Comparison of MIM and Virtual Sensors**

20 To validate the virtual sensor methodology over a complex and highly congested urban arterial
 21 network, we compare it with current-used Midtown in Motion system using the data collected for
 22 a period of two weeks. Although both methodologies are using a sample of the true population
 23 which may not represent the true travel time or speed among the studied road segments, virtual
 24 sensors can be a good supplementary data source at no additional cost if it can be shown that the
 25 information they provide is similar in accuracy to that of the MIM system.

26 Comparing time series data is usually not straightforward. Therefore, three different
 27 methods, namely Kullback-Leibler Divergence (KL Divergence), Mean Absolute Percentage
 28 Error (MAPE) and graphic representation of absolute difference are used. Graphic representation
 29 of absolute difference between MIM and virtual sensors data is first investigated for temporal
 30 patterns. Time series data that is originated from diverse natural processes can be conveyed by
 31 probabilistic distribution functions (PDFs) (35). KL Divergence from information entropy is often
 32 used to measure how much information is lost when comparing two PDFs. For two normalized,
 33 discrete PDF p and q , the KL Divergence can be represented as follows:

$$34 \quad D_{KL}(p \parallel q) = \sum_{i=1}^n \ln\left(\frac{p_i}{q_i}\right) p_i, \quad \text{with } D_{KL}(p \parallel q) \text{ equals to zero if } p=q. \quad (92)$$

35 Mean Absolute Percentage Error is popular in trend estimation and is applied in this study
 36 because of its intuitive interpretation in terms of relative error. Its formulation is as follows:

$$37 \quad MAPE = \frac{1}{n} \sum_{t=1}^n \frac{|A_t - B_t|}{A_t} \quad (13)$$

38 Where A_t and B_t are two observations from data source A and B at time t .

39

1 RESULTS AND DISCUSSIONS

2 Data Fusion Results

3 D_i , the prototype of the evidence (equation (9)) is obtained from the training set. After obtaining
 4 prototypes, the initial BPA of each evidence for each time step can be constructed. TABLE 1
 5 provides a representation of the initial BPAs for one-hour time-periods (12 time steps). As
 6 mentioned before, BPAs are an essential part of the D-S theory since they represents the degree of
 7 belief that a single traffic state in the subset is true, given the source of evidence (33). For instance,
 8 information extracted from the MIM system at time step 1 indicates a 50.16% degree of belief that
 9 the traffic on the reference link is in a traffic state S4 “Smooth” (TABLE 1). Next, the initial BPAs
 10 are discounted with the reliability matrix and are combined to obtain the final fusion results.
 11 Maximum belief is used as the decision strategy of the final fusion results.

12 **TABLE 1 Example of the Initial BPAs of the the Input Evidence for One Hour**

13

Time step*	Evidence	S1 (very congested)	S2 (congested)	S3 (moderate)	S4 (Smooth)	S5 (very smooth)
1	MIM	0.0428	0.0428	0.0428	0.5016	0.3699
1	Virtual sensors	0.0536	0.0536	0.0536	0.6235	0.2157
1	MTA Bus Time	0.0326	0.0326	0.0326	0.5613	0.3409
2	MIM	0.0906	0.0906	0.0906	0.0906	0.6377
2	Virtual sensors	0.0468	0.0468	0.0468	0.0468	0.8126
2	MTA Bus Time	0.0135	0.0135	0.0135	0.0135	0.9460
3	MIM	0.0614	0.0614	0.0614	0.5763	0.2395
3	Virtual sensors	0.0338	0.0338	0.0338	0.5484	0.3503
3	MTA Bus Time	0.0013	0.0013	0.0013	0.5352	0.4610
4	MIM	0.0376	0.0376	0.0376	0.4797	0.4074
4	Virtual sensors	0.0359	0.0359	0.0359	0.5064	0.3858
4	MTA Bus Time	0.0030	0.0030	0.0030	0.5271	0.4638
5	MIM	0.0344	0.0344	0.0344	0.4934	0.4033
5	Virtual sensors	0.0343	0.0343	0.0343	0.4506	0.4465
5	MTA Bus Time	0.0034	0.0034	0.0034	0.7979	0.1918
6	MIM	0.0360	0.0360	0.0360	0.4500	0.4420
6	Virtual sensors	0.0507	0.0507	0.0507	0.3531	0.4947
6	MTA Bus Time	0.0439	0.0439	0.0439	0.5267	0.3415
7	MIM	0.0457	0.0457	0.0457	0.4564	0.4065
7	Virtual sensors	0.0723	0.0723	0.0723	0.5030	0.2800
7	MTA Bus Time	0.0612	0.0612	0.0612	0.5668	0.2496
8	MIM	0.0336	0.0336	0.0336	0.3740	0.5252
8	Virtual sensors	0.0648	0.0648	0.0648	0.4513	0.3541
8	MTA Bus Time	0.0271	0.0271	0.0271	0.6376	0.2811
9	MIM	0.0327	0.0327	0.0327	0.3880	0.5139
9	Virtual sensors	0.0842	0.0842	0.0842	0.4246	0.3229
9	MTA Bus Time	0.0276	0.0276	0.0276	0.5679	0.3494
10	MIM	0.1187	0.1187	0.1187	0.2892	0.3546

10	Virtual sensors	0.0782	0.0782	0.0782	0.3441	0.4215
10	MTA Bus Time	0.0144	0.0144	0.0144	0.7733	0.1835
11	MIM	0.0381	0.1449	0.0381	0.3345	0.4444
11	Virtual sensors	0.0432	0.1533	0.0432	0.2956	0.4647
11	MTA Bus Time	0.0071	0.0302	0.0071	0.2735	0.6822
12	MIM	0.0374	0.0374	0.1169	0.2875	0.5207
12	Virtual sensors	0.0463	0.0463	0.1347	0.3441	0.4286
12	MTA Bus Time	0.0185	0.0185	0.0185	0.3604	0.5842

*Each time step represents 5 minutes.

TABLE 2 presents the out-of-sample accuracy of the three fusion techniques. Evidence Theory with improved reliability performs best among three approaches tested in this paper. Both the simple weighted and Random Forest models indicate that MTA Bus Time data have the largest impact on estimating accurate traffic state. Simple weighted model assigns the largest weight (7.24) to MTA data and Random Forest model ranks MTA data the top 1 feature contributing to the Mean Decrease Impurity of the model. Evidence theory with improved reliability, even though more computationally complex than a simple weighted model, has much more control over the spatio-temporal features of the traffic characteristics if such information is known. For example, a street block with many truck loading and unloading activities taking place during the mid-day period over weekdays will have an associated reliability matrix with lower values assigned to these time periods for the estimation of traffic states.

TABLE 2 Results of Fusion Accuracy

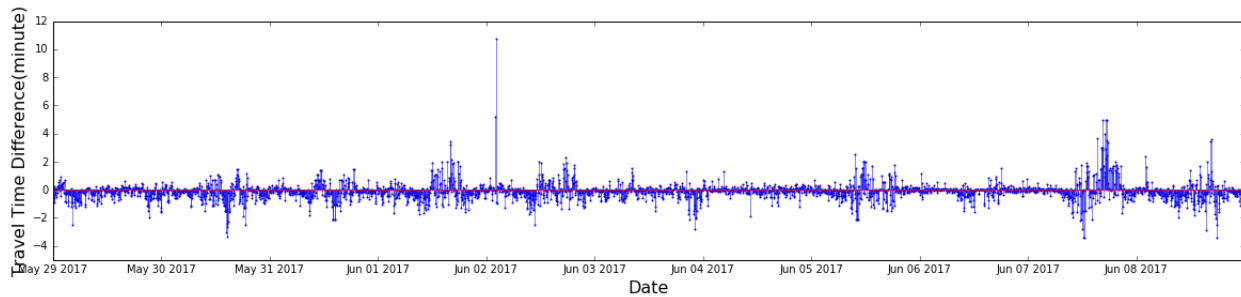
Method	Out-of-sample Cross-validate Accuracy
Simple Weighted	0.62 (Estimated weights: MIM: 0.17, Virtual Sensor: 0.57, MTA Bus Time: 7.24)
Random Forest	0.56 (Feature importance: MIM: 0.30, Virtual Sensor: 0.17, MTA Bus Time: 0.52)
Evidence Theory with improved reliability	0.75

Random Forest has the lowest accuracy among the three proposed approaches. One possible reason is that machine learning models depend on large historical data to learn the underlying patterns. It can be hypothesized that the seven days of data is too small to properly train a Random Forest classifier. This also leads to another challenge when conducting the study: obtaining the ground truth data representing the whole traffic population is extremely challenging and labor intensive. With the developments in machine learning/deep learning and image processing techniques, video cameras have gained importance and potential to be leveraged to automatically extract valuable information via image-processing techniques in the future.

Comparison of Results of MIM and Virtual Sensors

Two PDFs are constructed for MIM and virtual sensor data using Gaussian kernel estimation and KL Divergence. The result of KL Divergence indicates a small value close to zero (0.01393). Therefore, the information loss when using PDF of virtual sensor data as an approximation of PDF of MIM is small. MAPE test shows that a mean absolute percentage error of 14.23% of two datasets. However, it should be noted that using MAPE to evaluate the overall time period may not be the best approach because errors at night can be much less than the errors during congested

1 periods. To further investigate the temporal difference between the two datasets, the difference of
 2 travel time by date and time was plotted in the following figure. 95.8% of the records have a travel
 3 time difference less than one minute and 99.2% of the records have a travel time difference less
 4 than two minutes.



5
 6 **FIGURE 5 Mean travel time difference between MIM and Virtual Sensor.**

7 8 **CONCLUSIONS AND FUTURE WORK**

9 This paper leveraged travel time and traffic speed information from three different data sources,
 10 MIM, virtual sensors, and MTA Bus Time, to evaluate multiple fusion technologies on traffic state
 11 estimation. Two weeks of data extracted from the three data sources are utilized along with
 12 “ground truth” information collected from video cameras.

13 Virtual sensors have been proven to be a good additional data source over freeways like
 14 New Jersey Turnpike as shown in a previous study (5). This paper further validates its usefulness
 15 over a complex and highly congested urban arterial network. After applying MAPE, KL
 16 Divergence and Chi-Square test of independence, the information extracted from virtual sensors
 17 is shown to be a good data source that can improve NYCDOT’s current MIM system at no
 18 additional cost. Moreover, incidents or extreme weather conditions such as natural disasters have
 19 relatively fewer impacts on the virtual sensor methodology since it is based on GPS and map
 20 services instead of sensors installed on the road. This may overcome the missing data challenge
 21 caused by damaged or malfunctioning infrastructure sensors.

22 To investigate the applicability of various data fusion techniques, following three methods,
 23 namely: 1) simple weighted, 2) Random Forest, and 3) Evidence theory are applied using morning
 24 peak hour data for a reference street block in Midtown Manhattan, NYC. Although off-line data is
 25 used in the case study, all three techniques can be implemented in real-time. A reliability matrix
 26 is estimated for the evidence theory approach to discount the bias caused by small sample size and
 27 long road segments with multiple street blocks. The final result indicates that the evidence theory
 28 with improved reliability performs the best among all three tested models. The results also show a
 29 relatively heavier weight assigned to MTA Bus Time when fusing the information to estimate
 30 traffic states, and the underlying reason may be that MTA Bus Time has a better representation on
 31 the street block level.

32 However, the whole approach still has some deficiencies. Firstly, it should be noted that
 33 the results are obtained from a specific street block, and they might not be generalized to other
 34 urban areas with significantly different demand and infrastructure characteristics. Secondly, the
 35 reliability of the evidence is stochastic and has a spatio-temporal dynamic behavior. For example,
 36 street blocks with bus stations may have lower reliability in terms of MTA Bus Time data due to
 37 loading and unloading of the passengers. Therefore, extending the study area to cover more street
 38 blocks is one of the main future research objectives. In addition, the research team is developing
 39 an object detection and vehicle tracking application that can help to automate the image processing

1 of the video feeds to collect more relevant information to improve state estimations described in
2 this paper.

3 **ACKNOWLEDGMENTS**

4 The work in this paper is partially funded C2SMART, a Tier 1 University Transportation Center at
5 New York University. The authors acknowledge Dr. Mohamad Talas from New York City
6 Department of Transportation for sharing Midtown in Motion data. The contents of this paper only
7 reflect views of the authors who are responsible for the facts and do not represent any official
8 views of any sponsoring organizations or agencies.
9

10 **AUTHOR CONTRIBUTION STATEMENT**

11 The authors confirm contribution to the paper as follows: study conception and design: Jingqin
12 Gao, Kaan Ozbay; data collection: Jingqin Gao, Abdullah Kurkcu; analysis and interpretation of
13 results: Jingqin Gao, Kaan Ozbay; draft manuscript preparation: Jingqin Gao, Abdullah Kurkcu,
14 Kaan Ozbay. All authors reviewed the results and approved the final version of the manuscript.
15

16 **REFERENCES**

- 17 1. Federal Highway Administration (FHWA). *ITS system architecture and standards, FHWA ruling on*
18 *architecture standards and conformance of projects to ITS architecture standards*. 2001.
- 19 2. Seo, T., A.M. Bayen, T. Kusakabe, and Y. Asakura, Traffic state estimation on highway: A
20 comprehensive survey. *Annual Reviews in Control* 43, 2017, pp. 128-151.
- 21 3. Kurkcu, A. and K. Ozbay, Estimating Pedestrian Densities, Wait Times, and Flows with Wi-Fi and
22 Bluetooth Sensors. *Transportation Research Record: Journal of the Transportation Research*
23 *Board(2644)*, 2017, pp. 72-82.
- 24 4. Shlayan, N., A. Kurkcu, and K. Ozbay, Exploring pedestrian Bluetooth and WiFi detection at public
25 transportation terminals. In: *Proceedings of the ITSC, 2016*. pp. 229-234.
- 26 5. Morgul, E., H. Yang, A. Kurkcu, K. Ozbay, B. Bartin, C. Kamga, and R. Salloum, Virtual Sensors:
27 Web-Based Real-Time Data Collection Methodology for Transportation Operation Performance
28 Analysis. *Transportation Research Record: Journal of the Transportation Research Board(2442)*,
29 2014, pp. 106-116.
- 30 6. Kurkcu, A., M. Eng, E.F. Morgul, and K. Ozbay, Extended Implementation Methodology for Virtual
31 Sensors: Web-based Real Time Transportation Data Collection and Analysis for Incident
32 Management. In the compendium of the Transportation Research Board 94th Annual Meeting, 2015.
- 33 7. Bengio, Y., A. Courville, and P. Vincent, Representation learning: A review and new perspectives.
34 *IEEE transactions on pattern analysis and machine intelligence* 35(8), 2013, pp. 1798-1828.
- 35 8. Xin, W., J. Chang, S. Muthuswamy, and M. Talas, "Midtown in Motion": A New Active Traffic
36 Management Methodology and Its Implementation in New York City. In the compendium of the
37 Transportation Research Board 92nd Annual Meeting, Washington DC, 2013.
- 38 9. The Metropolitan Transportation Authority, MTA Bus Time, <http://bustime.mta.info/>, Accessed
39 March 3rd, 2018.
- 40 10. El Faouzi, N.-E., H. Leung, and A. Kurian, Data fusion in intelligent transportation systems:
41 Progress and challenges—A survey. *Information Fusion* 12(1), 2011, pp. 4-10.
- 42 11. Han, J., J. Pei, and M. Kamber, *Data mining: concepts and techniques*, Elsevier, 2011.
- 43 12. Huang, D. and H. Leung, An expectation-maximization-based interacting multiple model
44 approach for cooperative driving systems. *IEEE Transactions on Intelligent Transportation Systems*
45 6(2), 2005, pp. 206-228.
- 46 13. Kong, Q.-J., Z. Li, Y. Chen, and Y. Liu, An approach to urban traffic state estimation by fusing
47 multisource information. *IEEE Transactions on Intelligent Transportation Systems* 10(3), 2009, pp.
48 499-511.
49

- 1 14. Dempster, A.P., Upper and lower probabilities induced by a multivalued mapping. *The annals of*
2 *mathematical statistics*, 1967, pp. 325-339.
- 3 15. Shafer, G., *A mathematical theory of evidence*, Princeton university press Princeton, 1976.
- 4 16. El Faouzi, N.-E., L. Klein, and O. De Mouzon, Improving travel time estimates from inductive
5 loop and toll collection data with dempster-shafer data fusion. *Transportation Research Record:*
6 *Journal of the Transportation Research Board(2129)*, 2009, pp. 73-80.
- 7 17. Hashem, S., Optimal linear combinations of neural networks. *Neural networks 10(4)*, 1997, pp.
8 599-614.
- 9 18. Jiang, G.-y., J.-f. Wang, X.-d. Zhang, and L.-h. Gang, The study on the application of fuzzy
10 clustering analysis in the dynamic identification of road traffic state. In: *Proceedings of the Intelligent*
11 *Transportation Systems*, 2003. *Proceedings. 2003 IEEE*, 2003. pp. 1149-1152.
- 12 19. Xu, C., P. Liu, W. Wang, and Z. Li, Evaluation of the impacts of traffic states on crash risks on
13 freeways. *Accident Analysis & Prevention 47*, 2012, pp. 162-171.
- 14 20. Foerster, G., Traffic state estimation using hierarchical clustering and principal components
15 analysis: a practical approach. In: *Proceedings of the the European Conference of Transport Research*
16 *Institutes Young Reseachers Seminar 2007*, Brno, Czech Republic, 2007.
- 17 21. Bartin, B., K. Ozbay, and C. Iyigun, Clustering-based methodology for determining optimal
18 roadway configuration of detectors for travel time estimation. *Transportation Research Record:*
19 *Journal of the Transportation Research Board(2000)*, 2007, pp. 98-105.
- 20 22. Bartin, B. and K. Ozbay, Determining the optimal configuration of highway routes for real-time
21 traffic information: A case study. *IEEE Transactions on Intelligent Transportation Systems 11(1)*,
22 2010, pp. 225-231.
- 23 23. Zhang, S., B. Du, and N. Du, Mer-Gesh: A New Data Fusion Framework to Estimate Dynamic
24 Road Travel Time, *Geo-Informatics in Resource Management and Sustainable Ecosystem*. Springer,
25 2013, p.^pp. 1-15.
- 26 24. El Faouzi, N.-E. and L.A. Klein, Data fusion for ITS: techniques and research needs.
27 *Transportation Research Procedia 15*, 2016, pp. 495-512.
- 28 25. O'Connor, A.T., B. Schaller, M. Talas, J. Tipaldo, and S. Muthuswamy, New York City Active
29 Traffic Management, Midtown in Motion,
30 https://www.pcb.its.dot.gov/t3/s130418/s130418_mim_presentation.pdf, Accessed July 30, 2018.
- 31 26. NYCDOT, Real-Time Traffic Cameras, <http://www.nyc.gov/html/dot/html/motorist/atis.shtml>,
32 Accessed July 20, 2018.
- 33 27. Hänsch, R. and O. Hellwich, When to fuse what? random forest based fusion of low-, mid-, and
34 high-level information for land cover classification from optical and SAR images. In: *Proceedings of*
35 *the Geoscience and Remote Sensing Symposium (IGARSS)*, 2016 IEEE International, 2016. pp. 3587-
36 3590.
- 37 28. Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.
38 Prettenhofer, R. Weiss, and V. Dubourg, Scikit-learn: Machine learning in Python. *Journal of machine*
39 *learning research 12(Oct)*, 2011, pp. 2825-2830.
- 40 29. El Faouzi, N.-E., Data fusion in road traffic engineering: An overview. In: *Proceedings of the*
41 *Defense and Security*, 2004. pp. 360-371.
- 42 30. Choi, K. and Y. Chung, A data fusion algorithm for estimating link travel time. *ITS journal 7(3-*
43 *4)*, 2002, pp. 235-260.
- 44 31. Martin, R., J. Zhang, and C. Liu, Dempster-Shafer theory and statistical inference with weak
45 beliefs. *Statistical Science*, 2010, pp. 72-87.
- 46 32. Esposito, C., A. Castiglione, F. Palmieri, M. Ficco, C. Dobre, G. Iordache, and F. Pop, Event-
47 based sensor data exchange and fusion in the Internet of Things environments. *Journal of Parallel and*
48 *Distributed Computing*, 2018.
- 49 33. Ali, T. and P. Dutta, Methods to obtain basic probability assignment in evidence theory.
50 *International Journal of Computer Applications 38(4)*, 2012, pp. 46-51.

- 1 34. Denoeux, T., A k-nearest neighbor classification rule based on Dempster-Shafer theory. *IEEE*
- 2 *transactions on systems, man, and cybernetics* 25(5), 1995, pp. 804-813.
- 3 35. Kowalski, A.M., M.T. Martin, A. Plastino, and G. Judge, On extracting probability distribution
- 4 information from time series. *Entropy* 14(10), 2012, pp. 1829-1841.
- 5 36. McHugh, M.L., The chi-square test of independence. *Biochemia medica: Biochemia medica*
- 6 23(2), 2013, pp. 143-149.
- 7
- 8
- 9
- 10
- 11