

# **National Transportation Center**

# Project ID: NTC2014-SU-R-02 COMBINING DIFFERENT DATA SOURCES TO PREDICT ORIGIN-DESTINATIONS AND FLOW PATTERNS FOR TRUCKS IN LARGE NETWORKS

### **Final Report**

by

Dr. Mecit Cetin MCetin@odu.edu (757) 683-6700 Old Dominion University

Olcay Sahin Ilyas Ustun Old Dominion University

for

National Transportation Center at Maryland (NTC@Maryland) 1124 Glenn Martin Hall University of Maryland College Park, MD 20742

#### April, 2015

#### ACKNOWLEDGEMENTS

This project was funded by the National Transportation Center @ Maryland (NTC@Maryland), one of the five National Centers that were selected in this nationwide competition, by the Office of the Assistant Secretary for Research and Technology (OST-R), U.S. Department of Transportation (US DOT). We also would like to thank VDOT for providing the datasets.

#### DISCLAIMER

The contents of this report reflect the views of the authors, who are solely responsible for the facts and the accuracy of the material and information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation University Transportation Centers Program in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof. The contents do not necessarily reflect the official views of the U.S. Government. This report does not constitute a standard, specification, or regulation.

# **TABLE OF CONTENTS**

| EXEC | CUTIVE SUMMARY                                   |              |
|------|--|--------------|
| 1.0  | INTRODUCTION                                     |              |
| 2.0  | DATA SOURCES AND CREATION OF TRAINING DATA FOR M | <b>IODEL</b> |
| DEVE | LOPMENT AND TESTING                              | 5            |
| 2.1  | CONTINUOUS COUNT STATIONS AND CAMERA RECORDINGS  | 5            |
| 2.2  | INRIX DATA                                       | 7            |
| 2.   | .2.1 Instantaneous and Experienced Travel Time   | 7            |
| 3.0  | LITERATURE SURVEY AND THE METHODOLOGY            | 11           |
| 3.1  | BAYESIAN MODEL FOR RE-IDENTIFICATION             |              |
| 3.2  | DETERMINING THE SEARCH SPACE                     |              |
| 3.3  | BAYESIAN METHOD                                  | 12           |
| 4.0  | EXPERIMENTS AND RESULTS                          |              |
| 4.1  | SCENARIO 1                                       | 15           |
| 4.2  | SCENARIO 2                                       |              |
| 4.3  | SCENARIO 3                                       | 19           |
| 5.0  | SUMMARY AND CONCLUSIONS                          |              |
| 6.0  | REFERENCES                                       |              |
| 7.0  | APPENDICES                                       |              |

#### LIST OF TABLES

| Table 1: Recorded video dates and times  | 5         |
|--|-----------|
| Table 2: Number of total upstream and downstream vehicles (FHWA Classes 4, 6, 7, | 8, 9, 10, |
| 11)  | 7         |
| Table 3: Testing results for scenario 1  |           |
| Table 4: Testing results for scenario 2 and scenario 3                           |           |

### LIST OF FIGURES

| Figure 1: CCS, Camera, and INRIX TMC locations   | 6    |
|--|------|
| Figure 2: Graphical representations of instantaneous and experienced travel time             | 8    |
| Figure 3: Length of the vehicle (a&c) and spacing between axle 1 and 2 (b&d) at two stations | for  |
| matched and mismatched trucks  | . 14 |
| Figure 6: Search window with 14 minutes minimum travel time & 10 minutes time window         | / in |
| Scenario 1   | . 16 |
| Figure 7: Percent of accurate matches with static travel times (September 17, 2014 results)  | . 16 |
| Figure 8: Percent of accurate matches with static travel times (September 22, 2014 results)  | . 17 |
| Figure 9: INRIX travel time vs. observed travel time   | . 18 |
| Figure 10: Scenario 2 - Dynamic travel time with a 10-minute window                          | . 18 |
| Figure 11: Percent of accurate matches with dynamic search space (September 17, 2014 resu    | lts) |
|  | . 19 |

| Figure 12: Percent of accurate matches with dynamic search space (September 22, 2014 | 4 results) |
|--|------------|
|  |            |
| Figure 13: Scenario 3 - Dynamic travel time & 2 conditions time window               |            |
| Figure 14: Percent of accurate matches with dynamic search space and varying         | windows    |
| (September 17, 2014 results).  |            |
| Figure 15: Percent of accurate matches with dynamic search space and varying         | windows    |
| (September 22, 2014 results)   |            |
| Figure 16: Sensitivity analysis of fixed TT when the best TT is changed by 2 minutes |            |
| Figure 17: Maximum % accurate match for each scenario of all data set                |            |
|  |            |

### **EXECUTIVE SUMMARY**

Ground transportation is essential for national and international freight movement. In terms of tons and value of the transported goods, trucking is the dominant method used in the USA. Thus, efficient management of freight transportation is essential. Monitoring freight movement and the performance of the system as a whole is important for making informed decisions. One of the key aspects of monitoring freight over the highways has to do with determining the flow patterns of trucks, which can be achieved by uniquely identifying trucks at specific points along the roads or by tracking individual trucks using technology such as GPS. Both methods require investment in technology. Maintenance cost of equipment is another significant factor to be considered in the case of infrastructure-based sensing. For the case of GPS tracking, since most of the trucks are owned by private parties, they may not be willing to share these data due to privacy concerns.

In this research a method is proposed that is capable of tracking trucks by using anonymously collected data from sensors already in place. The data collected from existing vehicle count and classification stations are utilized. The attributes collected such as length of truck, number of axles, and axle spacings provide valuable information for matching the same truck passing through two stations. The variance in these attributes between different trucks provides a means for re-identifying the same truck. Although there will be measurement errors between two stations due to speed, weather, interference, calibration of devices at stations; measurements from matched trucks still exhibit a distinct pattern in which the difference of measurements between two stations will be less compared to data from non-matched trucks. The feasibility of matching trucks anonymously based on axle data have been demonstrated in previous studies.

In this project, the previously developed models are enhanced to investigate the value of incorporating travel times provided by private companies (e.g., INRIX) into the vehicle reidentification algorithms. Therefore, the source data for this project consists of both INRIX data and attribute data from vehicle classification sites. For this project, the needed data are collected from two vehicle classification sites along the I-64 corridor in Hampton Roads, VA. Per vehicle data from the classification sites include a timestamp, vehicle class, speed, number of axles, axle to axle spacing, and overall length for each vehicle. Trucks crossing upstream and downstream sites are manually identified from the recorded video files so that the results from the vehicle reidentification algorithms can be validated. The re-identification algorithms are applied with different options for incorporating INRIX travel times. Since the selected I-64 corridor experience recurrent congestion, the collected datasets include varying levels of traffic conditions. Change in travel time between congested and free-flow conditions significantly impacts the performance of the re-identification model. It is found that using a dynamic travel time window informed by the INRIX data significantly improves the accuracy of the vehicle reidentification results. For some tested cases, the improvement in accuracy is up to 19% when compared to the results from static search windows. Results also show that dynamic search windows provide more robust results against small perturbations in travel times.

### **1.0 INTRODUCTION**

While other modes are clearly important for freight transportation, trucking is the dominant mode in terms of tons and value. Monitoring freight movement and freight transportation performance is essential in making effective policies and informed decisions to enhance and to efficiently manage the freight transportation system. One of the key aspects of monitoring freight over the highways has to do with determining the flow patterns of trucks, which can be achieved by uniquely identifying trucks at specific points along the roads or by tracking individual trucks using technology such as GPS. However, not all trucks are equipped with tracking devices. While point sensors along the highways allow determining the truck volumes, they do not provide much information about the paths and origin-destinations for trucks. However, by exploiting vehicle-specific attributes (e.g., axle spacings, length) collected by such sensors vehicles can be re-identified (matched) to enable prediction of paths taken by trucks. Data from other infrastructure-based sensors (e.g., Bluetooth readers, AVI sensors) can also be utilized for the same purpose. Furthermore, such data elements can be combined with freight generators in a network (e.g., ports, distribution centers) to better determine origins and destinations. Developing such a system where data from all these sources are assimilated and synthesized to predict freight patterns will be useful for planning and performance monitoring of the national freight network.

In this project, re-identification models for matching vehicles between two Continuous Count Stations (CCSs) are developed. At a typical CCS, total vehicle length, and axle spacings are measured per vehicle basis. Such data are then archived for future use. In addition to data from CCS sensors, it is assumed that travel time information (or variation) between the two sites is available. Such information can be obtained from various sources, including private companies (e.g., INRIX) or estimated from point sensors (e.g., loops, radar) installed along the corridor. The travel time information along with CCS data is incorporated into algorithms to re-identify trucks. In previous models, travel time between the sites is usually assumed to be constant and the variation in travel time is ignored (Cetin et al., 2011a, Cetin and Nichols, 2009, Cetin et al., 2011b), which is not realistic especially along corridors through congested urban areas. By incorporating the travel time variation, the re-identification algorithms have proven to produce more accurate matching, which is explored in this project. Therefore, the main objectives of this project include:

- Developing vehicle re-identification algorithms that can integrate travel time information with CCS data for matching trucks between two sites
- Assessing the accuracy of the re-identification algorithms as a function of the reliability of the travel time information

In order to conduct the proposed research, both travel time and truck attribute data are needed. In a previous project, the Principal Investigator worked with the weigh-in-motion (WIM) data for twenty stations across Oregon (Cetin et al., 2011b). The WIM sites in Oregon are equipped with sensors that can measure axle weights, axle spacing, and gross vehicle weight estimates that are uniquely matched to each truck (Elkins and Higgins, 2008). Since some of the trucks (20- to 35%) are carrying radio-frequency identification (RFID) transponders, these measured attributes

are also uniquely matched to transponder-equipped trucks. These particular trucks provided the needed truck attribute data for model development and testing. This project is built on the previous re-identification methodologies (Cetin et al., 2011a, Cetin et al., 2011b) but employs new datasets collected in Hampton Roads, VA. These datasets include vehicle length and axle spacing data from CCSs and travel time information from INRIX. In addition, the previous methodologies are enhanced by incorporating travel time information into the re-identification algorithms.

The following section describes the study site and the data sources. Section 3 provides a brief literature review and discusses the re-identification algorithms. Empirical results are given in Section 4 which is followed by conclusions in Section 5.

## 2.0 DATA SOURCES AND CREATION OF TRAINING DATA FOR MODEL DEVELOPMENT AND TESTING

This project makes use of 3 different sources of datasets: Vehicle attribute data from Continuous Count Stations (CCSs), video data, and speed data. The datasets are provided by the Virginia Department of Transportation (VDOT). First dataset contains truck attribute data for approximately 5 months from June 2014 to October 2014. The second dataset contains video recordings of the highway for a week of 7 days starting on Wednesday 17<sup>th</sup> of September, ending Tuesday September 23<sup>rd</sup> 2014. This dataset is used for ground truth validation and model training. The third source is the INRIX speed data which will serve as dynamic travel time input to the re-identification algorithm. These are described below.

#### 2.1 CONTINUOUS COUNT STATIONS AND CAMERA RECORDINGS

The continuous count stations have sensors to capture vehicle attributes for each lane. The per vehicle records (PVR) from CCSs include a timestamp, vehicle class, speed, number of axles, axle to axle spacing, and overall length for each vehicle. The two continuous count stations that are used in this project will be referred to as upstream and downstream stations throughout the report. The locations of these stations are shown on the map in Figure 1 with a red star icon. Both sites are on I-64 EB along the Hampton Roads Bridge Tunnel (HRBT) corridor, one of the most congested major freeways in the Hampton Roads. The downstream site is at the exit of the HRBT tunnel.

The data recorded at these stations do not have unique identifier for the trucks. Thus, in order to develop, validate, and test the proposed re-identification technique, a labelled dataset is needed. In order to find the unique matches of upstream trucks to downstream ones, videos recorded from VDOT's traffic surveillance cameras are used. The locations of traffic cameras can be seen on the map in Figure 1 as a red camera icon. Table 1 provides information on the time frames of the recorded videos for the days that were present in this dataset.

| Camera     | 9/17/2014 | 9/18/2014 | 9/19/2014 | 9/20/2014 | 9/21/2014 | 9/22/2014 | 9/23/2014 |
|------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| CAM<br>802 | 7AM-9AM   | -         | -         | ALL DAY   | ALL DAY   | ALL DAY   | -         |
|            | 4PM-6PM   |           |           |           |           |           |           |
| CAM<br>822 | 7AM-9AM   |           | 4PM-MN    | 7AM-9AM   | 7AM-9AM   | 7AM-9AM   | · 7AM-9AM |
|            | 4PM-6PM   | -         |           | 4PM-6PM   | 4PM-6PM   | 4PM-6PM   |           |

#### Table 1: Recorded video dates and times



Figure 1: CCS, Camera, and INRIX TMC locations

Through a manual process, the video footage captured by cameras 802 and 822 was used to identify and extract trucks that cross both stations. Camera 802 is directly facing the upstream station and had a good angle to visually recognize trucks. The video captured by a different camera at the downstream station was not clear therefore camera 822 was used instead to identify vehicles. This camera is located at the entrance of the HRBT and there is no exit or entrance once the trucks enter the bridge. Thus, the video captured from this camera provides the same information as the camera on the bridge. Since these cameras capture the common vehicles traveling from the upstream to downstream, they provide a means for matching the trucks. If the two trucks match based on visual inspection, they are uniquely labeled and then associated with the corresponding trucks in the PVR data from CCSs by utilizing the timestamps from video and CCSs. The labeled PVR data now include unique identifiers and provide the ground truth data for model development and testing. The labeling is done for only those vehicles that belong to FHWA vehicle classes of 4, 6, 7, 8, 9, 10, and 11.

Among the one week video data, the recordings from September 17 and 22 are selected to create the ground truth labeled data. The area of study exhibits congested and free flow travel conditions based on the time of day. Since free flow travel time can be known based on the distance between stations and since the variation between travel times of vehicles is very low during free flow, congested traffic condition is of primary interest which exhibits a large variation. Hence, morning and afternoon rush hours were taken to create the labeled subset data. Table 2 shows the number of upstream and downstream vehicles (FHWA Classes 4, 6, 7, 8, 9, 10, 11) observed at the two stations for September 17, 7:00-9:00 AM and 4:00-6:00 PM and

September 22, 7:00-9:00 AM and 4:00-6:00 PM. Based on the visually observed trucks from the videos, there are 265 matched trucks. However, some outliers were found to be present, perhaps, due to sensor measurement errors. Once these outliers were removed, the remaining subset of matched trucks has 219 vehicles.

|                             | 1              | 1              |                | 1              |  |
|-----------------------------|----------------|----------------|----------------|----------------|--|
| Station                     | 9/17/2014      | 9/17/2014      | 9/22/2014      | 9/22/2014      |  |
|                             | 7:00 – 9:00 AM | 4:00 – 6:00 PM | 7:00 – 9:00 AM | 4:00 – 6:00 PM |  |
| Upstream                    | 498            | 302            | 458            | 233            |  |
| DownStream                  | 248            | 131            | 214            | 105            |  |
| <b>#of Matched vehicles</b> | 06             | 52             | 60             | 17             |  |
| (with outliers)             | 90             | 55             | 09             | 47             |  |
| #of Matched vehicles        | 70             | 22             | 65             | 12             |  |
| (w/o outliers)              |                | 32             | 03             | 43             |  |

Table 2: Number of total upstream and downstream vehicles (FHWA Classes 4, 6, 7, 8, 9, 10, 11)

The variability of travel time between two stations plays a critical role in vehicle re-identification since the potential matches are searched within expected travel time windows. The objective of this research is the assimilation of a third party data source to get travel time dynamically for re-identification of vehicles between two stations. As mentioned earlier this information is obtained from INRIX which is explained in detail in the next section.

#### 2.2 INRIX DATA

VDOT provided the research team access to INRIX speed data for the analyses. INRIX data are being used by VDOT to provide travel times on dynamic message signs and has been used in previous VDOT studies. The Virginia Center for Transportation Innovation and Research (VCTIR) evaluated the quality of the INRIX data by comparing it to travel times generated by Bluetooth readers<sup>1</sup>. Based on this VCTIR study, INRIX travel time data were found to meet the accuracy and availability VDOT's benchmarks.

INRIX data were aggregated into five-minute intervals and associated with road segments (i.e., TMCs) along the study corridor. There were 22 links (i.e., TMCs) with lengths ranging from 0.1 to 2.4 mile in length between the upstream and downstream. Figure 1 illustrates the INRIX links (or TMC segments) with the blue lines and a link's starting point with the yellow pin icon. As explained in the next section the experienced travel time has been calculated to estimate the time traveled between the two stations based on the INRIX data.

#### 2.2.1 Instantaneous and Experienced Travel Time

Given a spatio-temporal speed profile of traffic, two types of travel times can be computed: 1) instantaneous travel time and 2) experienced or dynamic travel time (Chen and Rakha, 2014, Mazaré et al., 2012, Tu, 2008). Instantaneous travel time is the time required for a vehicle to travel through a particular route, provided that the traffic conditions and speed of the vehicles

<sup>&</sup>lt;sup>1</sup> Fontaine, M.D., Evaluation of INRIX Travel Time Data in Virginia, VCTIR Report, October 2013.

remain unchanged over a specified amount of time. Dynamic travel time is the time a vehicle would actually require traveling along a given route considering the changes in speed of traffic as it travels through the road segments. If there is no change in speed, both instantaneous and experienced travel times provide similar results; however, when there is a sudden change in speed, the travel time provided by these two methods is not the same. Figure 2 is a graphical representation of instantaneous and experienced travel times overlaid on a speed heat-map created from INRIX traffic data.



Figure 2: Graphical representations of instantaneous and experienced travel time

INRIX provides the travel time to cover a segment over a course of time. In this study, experienced travel time is computed as follows:

Let *i* represent link number

 $t_i$  is the time when a vehicle exits at link i

 $TT_i^t$  is the travel time to traverse link *i* 

Assuming a vehicle exits link *i* at time  $t_i$ , the time the vehicle exits the previous link *i*-1 is given by the current time  $t_i$  minus the travel time of the link  $TT_i^{ti}$ , as shown in equation (1).

$$t_{i-1} = t_i - TT_i^{ti} \tag{1}$$

INRIX travel time data aggregated every 5 minute is used in this study. Therefore, linear interpolation between the two successive time intervals is used to determine the travel time  $TT_i^{ti}$  at any time *t*. Assuming  $t_n$  and  $t_{n-1}$  are the INRIX time intervals which are the immediate before and after time  $t_i$  (i.e.,  $t_n$ -1 <  $t_i$  <  $t_n$ ),  $TT_i^{ti}$  is given by equation (2):

$$TT_i^{t_i} = \frac{(TT_i^{t_n} - TT_i^{t_{n-1}}) * (t_i - t_{n-1})}{(t_n - t_{n-1})} + TT_i^{t_i}$$
(2)

This way, the time  $t_i$  when a vehicle arrives at each segment *i* can be easily determined. Therefore, the experienced travel time as a vehicle travels from segment *i* to segment *i*-*n* is given by the difference in the time when the vehicle arrives at segment *i* and the time the vehicle departs segment *i*-*n*, as shown in equation (3).

$$Experienced \ Travel \ Time = \ t_i - t_{i-n} \tag{3}$$

### **3.0 LITERATURE SURVEY AND THE METHODOLOGY**

Vehicle re-identification methods rely on the variability within the vehicle population and the ability to accurately identify the pairs of measurements collected at upstream and downstream stations that are generated by the same vehicle. These measurements can either be the actual physical attributes of vehicles such as length (Coifman and Cassidy, 2002) and axle spacing (Cetin and Nichols, 2009) or some characteristics of the sensor waveform or inductive vehicle signature (Sun et al., 1999). Researchers have developed various methods, such as lexicographic optimization (Sun et al., 1999, Oh et al., 2007), decision trees (Tawfik et al., 2004), etc, to reidentify vehicles. In a typical implementation of these methods, a downstream vehicle is matched to the most "similar" upstream vehicle (or vice versa) based on some defined metric (e.g., Euclidian distance). The resulting accuracy of these methods depends on several factors including the variation of the attribute data from vehicle to vehicle, number of attributes, the distance between data collection stations, variability of travel time, and type of the re-identification algorithm used.

Vehicle re-identification methods can been used to anonymously match vehicle crossing two different locations based on vehicle attribute data measured by sensors at each location. Let U and D be two nonempty sets that denote the vehicle crossing the upstream and downstream, respectively. Depending on various factors including the station locations, record validity, and types of activity between the sensors, four general cases arise:

i.  $U \subset D$  and  $U \neq D$ ii.  $D \subset U$  and  $U \neq D$ iii. U = Div.  $D \not\subset U$ ,  $U \not\subset D$ , and  $U \cap D \neq 0$ 

In the first three cases there is always a match for a vehicle in the smaller set (or either case for iii). However in the fourth case one needs to consider the possibility that a vehicle taken from one set might not have a match in the other set. The case that this project falls in would be iv, as there are entrances and exits between the upstream and downstream locations. However, since the testing is performed with the manually matched vehicles, case ii above is more pertinent to this project. For case iv, certain thresholds need to be used to eliminate vehicles for which a match does not exist (Cetin et al., 2011b).

In the following section, the Bayesian re-identification algorithm with a fixed travel time is explained. This is the technique used in previous research (Cetin et al., 2011b). This method does not take into account the INRIX speed data or other means of travel time information; hence the travel time variation is not captured. The search-space identification step of this method is modified to incorporate travel time information.

#### **3.1 BAYESIAN MODEL FOR RE-IDENTIFICATION**

Let  $X^U$  and  $X^D$  be two matrices with the same number of columns that denote the data collected at an upstream station and a downstream station, respectively; and  $X^U_i$  and  $X^D_j$  denote rows of these two matrices that correspond to the measurements (e.g., axle spacings) taken for vehicle *i* at the upstream station and for vehicle *j* at the downstream station. Further, assume that the time stamps indicating arrival times of vehicles at each station are given and denoted by  $t^U_i$  for the upstream vehicles and  $t^D_j$  for the downstream vehicles. Given  $X^U$ ,  $X^D$ ,  $t^U_i$  and  $t^D_j$  the vehicle matching problem involves determining  $X^U_i$  and  $X^D_j$  that are generated by the same vehicle. Let  $\delta_{ij}$  be a binary variable that equals 1 if  $X^U_i$  and  $X^D_j$  belong to the same vehicle and equals zero otherwise. The main objective of the matching algorithms is to estimate all  $\delta_{ij}$ 's with minimum error.

#### **3.2 DETERMINING THE SEARCH SPACE**

For re-identification, each vehicle in D needs to be matched to the most similar vehicle in U. A reasonable "search space" from the upstream vehicle records (U) can be identified based on expected travel time and some window. Before the search starts to match a downstream vehicle j to an upstream vehicle i, a search space for vehicle j, denoted by  $S_j$ , is determined. The variability in travel time can be captured by specifying minimum and maximum values for travel times. The minimum value (*minTime*) can be easily predicted based on an assumed maximum travel speed and the distance between the two stations. The maximum value can exhibit a large variation depending on the individual vehicle speeds, and traffic flow interruptions such as congestion or incidents between the two stations. The maximum value (*maxTime*) can be taken as some window added to the *minTime*. The search space for a downstream vehicle j is then determined as follows:

$$S_{i} = \{i \in \boldsymbol{U} \mid \boldsymbol{t}^{\boldsymbol{p}}_{j} - maxTime \leq \boldsymbol{t}^{\boldsymbol{U}}_{i1} \leq \boldsymbol{t}^{\boldsymbol{p}}_{j} - minTime \}$$

$$\tag{4}$$

Depending on the difference between *maxTime* and *minTime* or time window, the number of vehicles among which a match to be found varies. Larger time windows will result in a larger number of vehicles in the search space, which can make the matching problem more difficult. Too small of a window might cause the actual match to be missed. Finding an optimum window size is a problem in itself. The width of the time window can be adjusted depending on external travel time information (e.g., INRIX) as suggested in this research.

#### **3.3 BAYESIAN METHOD**

The Bayesian re-identification method relies on calculating the posterior probability of a match between two vehicles given two sets of data points collected for a vehicle pair (i,j) at the upstream and downstream stations. A vehicle j at the downstream station is matched to the upstream vehicle i that yields the largest probability of a match. The steps of the Bayesian method are more formally explained below.

For each vehicle *j* in DIdentify a search space  $S_j \subset U$  For each  $i \in S_j$ Calculate  $P(\delta_{ij} = 1 | \text{data})$  $m = \operatorname{argmax} P(\delta_{ij} = 1 | \text{data})$ 

Match vehicle *j* to *m*, i.e.,  $\delta_{ij} = 1$  if i = m

Once a search space is identified,  $P(\delta_{ij} = 1 | \mathbf{x}_{ij})$ , the conditional probability that  $X_i^U$  and  $X_j^D$  belong to the same vehicle given data (i.e.,  $\mathbf{x}_{ij} = \mathbf{x}_i^U \cup \mathbf{x}_j^D$ ), can be computed by the Bayes' theorem as follows:

$$P(\delta_{ij} = 1|x_{ij}) = \frac{f(x_{ij}|\delta_{ij} = 1)P(\delta_{ij} = 1)}{f(x_{ij}|\delta_{ij} = 1)P(\delta_{ij} = 1) + f(x_{ij}|\delta_{ij} = 0)P(\delta_{ij} = 0)}$$
(5)

In order to calculate this posterior probability, both the two conditional probability density functions (i.e.,  $f(x_{ii}|\delta_{ii}=1)$  and  $f(x_{ii}|\delta_{ii}=0)$ ) and the prior probabilities (i.e.,  $P(\delta_{ii}=0)$  and  $P(\delta_{ii}=1)$ ) are needed. The functions  $f(x_{ii}|\delta_{ii}=1)$  and  $f(x_{ii}|\delta_{ii}=0)$  are the density functions that characterize the collected data at two stations when it belongs to the same vehicle and different vehicles, respectively. Figure 3a-b and Figure 3c-d illustrate how the data distribute for observations at two stations when  $\delta_{ii}=1$  and  $\delta_{ii}=0$ , respectively, for a simple case when only a single attribute is considered. As it can be observed from these figures, when vehicles match (i.e., upstream and downstream measurements belong to the same vehicle) there is high correlation between the measurements, which is critical for re-identification. On the other hand, when random data for upstream and downstream measurements are plotted the correlation disappears as expected and a roughly uniform distribution of points is observed (Figure 3c-d). Since this amounts to an approximately uniform value for the density function,  $f(x_{ii}|\delta_{ii}=0)$  in equation (5) can be replaced by a constant ( $\alpha$ ). Furthermore, the travel time information can be used to approximate the prior distribution  $P(\delta_{i}=1)$ , as opposed to assigning a fixed value to the prior. In other words, the travel times of the common vehicles that cross both upstream and downstream vehicles are used to create a density function to replace  $P(\delta_{ii}=1)$ . If the probability density function for the travel time is denoted by,  $h(t_{ij})$  then, the posterior probability in equation (5) can be simplified to:

$$P(\delta_{ij} = 1|x_{ij}) \sim \frac{f(x_{ij}|\delta_{ij} = 1)h(t_{ij})}{f(x_{ij}|\delta_{ij} = 1)h(t_{ij}) + \alpha}$$
(6)

where  $\alpha$  is a positive arbitrary constant accounting for  $f(x_{ij}|\delta_{ij}=0)$  and  $f(\delta_{ij}=0)$ . Since in matching vehicles only relative magnitude of this posterior probability is important, the selected value of  $\alpha$  is not critical. In Principal Investigator's previous work equation (5) was used to calculate the posterior probability (Cetin and Nichols, 2009). In this project the simplified version (equation 6) is used which gives essentially the same results as before but does not require the estimation of  $f(x_{ij}|\delta_{ij}=0)$ , which is an advantage in terms of model calibration and development.

The conditional density function,  $f(x_{ij}|\delta_{ij}=1)$ , is obtained by fitting a probability distributions to a training dataset in which all vehicles are correctly matched based on the unique ID's created by manual video analysis. The data in this case are comprised of the attributes of matched vehicles  $(x_{ij} = x^{U_i} \cup x^{D_j})$ .



Figure 3: Length of the vehicle (a&c) and spacing between axle 1 and 2 (b&d) at two stations for matched and mismatched trucks

### 4.0 EXPERIMENTS AND RESULTS

After model training is completed and the parameters of the Bayesian model are obtained the model testing is performed. Since there are only two days of true match data available, namely September 17 and 22, one day is used for model training and the other for testing and vice versa. Training is done on the whole day, while for testing the data for morning and afternoon rush hours are used separately. The numbers of vehicles pertaining to these datasets are shown previously in Table 2.

The travel time information from INRIX is used to adjust the search space as defined previously (see equation 4). The search space is defined by two parameters: minimum travel time and time window. Adding the time window to the minimum travel time gives the maximum travel time. Three scenarios are constructed and tested to analyze the benefits of using INRIX data and dynamic windows for the search space:

- 1. *Static travel time with fixed window*: In this base scenario, the minimum travel time is taken to be a constant throughout the analysis period. And the maximum travel time is always assumed to be 10 minutes longer than the minimum. To analyze the impact of various travel times, the minimum travel time is varied from 10 to 25 minutes. These travel times are based on the distance between the two stations, the speed limit, and observed travel times under congestion.
- 2. *Dynamic travel time with fixed window*: Using INRIX speed data, the minimum travel time is varied based on the prevailing conditions in the field. As in scenario one, a fixed window of 10 minutes is used.
- 3. *Dynamic travel time with a varying window:* Similar to scenario two, a dynamic minimum travel time based on INRIX speed data is used. However, the time window is adjusted as explained below.

#### 4.1 SCENARIO 1

In scenario 1, a static value is used for the minimum travel-time that is needed to specify the search window. However, since the matching results depend on the selected value, a reasonable range for the minimum travel-times is considered. The minimum travel-time is varied from 10 to 25 minutes in increments of 1 minute to find the best possible static value for this scenario. All test plots for scenario 1 can be found in appendix A. Figure 4 shows one case where a fixed travel time is used. Here, the green dots represent the observed travel time for a vehicle crossing the downstream station. The blue and red dots represent the minimum and maximum travel times (i.e., search space) for each downstream vehicle. The difference between the red and blue dots makes up the search window. As can be seen, in some cases this search window fails to capture the observed travel times. If the search window is enlarged in order to accommodate all the travel times, the probability of an inaccurate match increases since the number of vehicles in the search window also increases.

The scenario has been tested on data from two different days. First, the parameters of the Bayesian model are obtained by training the model on the 22<sup>nd</sup> of September data, and the model is tested on data from 17<sup>th</sup> of September morning and afternoon rush hours. The minimum travel times vary from 10 to 25 minutes while time window is fixed at 10 minutes. The accuracy, which is calculated as the number of true matches found divided by the actual total true matches in the dataset, is shown in Figure 5. Morning rush hour reaches its best accuracy of approximately 58% with a 13 minute minimum travel time, while the afternoon period reaches its best accuracy of approximately 65% with an 18 minutes travel time. This illustrates that the same time window is not optimal for different traffic conditions.



Figure 4: Search window with 14 minutes minimum travel time & 10 minutes time window in Scenario 1



Figure 5: Percent of accurate matches with static travel times (September 17, 2014 results).

Second, the Bayesian model is trained on the September17<sup>th</sup> data, and is tested on 22<sup>nd</sup> of September morning and afternoon rush hours. This day exhibits a totally different performance, see Figure 6. The overall accuracy is significantly lower than before. Morning rush hour reaches its best accuracy of approximately 60% with a 22 minute minimum travel time, while the afternoon period reaches its best accuracy of approximately 78% with a 12 minutes travel time. It also is interesting to see that the best travel time of 12 minutes for the afternoon is the worst for the morning. Figure 6 shows these results in a bar plot for each travel time tested with a fixed window search size of 10 minutes. These results also illustrates that the same time window or search space is not optimal for different traffic conditions.





#### 4.2 SCENARIO 2

In scenario 2, the experienced travel time as defined before is calculated from the INRIX data and is used to determine the minimum travel time for the trucks. As in the previous scenario, 10 minutes is added to this travel time in to obtain the search space in the upstream. When the observed travel times and the INRIX experienced travel times are compared it is discovered that the observed travel time is sometimes less than the INRIX travel time, see Figure 7. In order to accommodate this and to find and optimum search space, the minimum travel times are shifted down from the reference INRIX times ranging from 1 minute to 5 minutes. Figure 8 shows one case where the minimum dynamic travel times are found by lowering the INRIX travel times by 3 minutes. Here, the green dots represent the observed travel time for a vehicle taken in downstream. The blue and red dots represent the minimum and maximum travel times for each downstream vehicle. The difference between the red and blue dots makes up the search window for that vehicle. As can be seen, as opposed to fixed time travel shown in Figure 4, the dynamic travel time is able to track the actual travel times much better.

The same cases tested in scenario 1 apply for scenario 2 as well. All the test plots for scenario 2 can be found in appendix B. Figure 9 and Figure 10 show the percent of true matches for morning and afternoon datasets for the two days. There is an overall increase in accuracy

compared to scenario 1. Also there is no big discrepancy between the two cases, as they both have approximately the same accuracy for the morning and afternoon periods.



Figure 7: INRIX travel time vs. observed travel time



Figure 8: Scenario 2 - Dynamic travel time with a 10-minute window



Figure 9: Percent of accurate matches with dynamic search space (September 17, 2014 results)



Figure 10: Percent of accurate matches with dynamic search space (September 22, 2014 results)

#### 4.3 SCENARIO 3

In scenario 3, as in scenario 2 dynamic travel times are found from the INRIX data. In the analysis of scenario 2 it is noticed that when the traffic conditions approach free flow the variation in observed travel times from the INRIX travel times was much lower. Due to this fact taking a smaller window size at these time periods would potentially increase the accuracy. This is simply because of the fact that fewer vehicles will be present in the search space, thus decreasing the chance of false matches. Thus, instead of a fixed time window of 10 minutes, two different time windows are used. When the INRIX travel time is less than 16 minutes, a window

of 5 minutes is used, and when the INRIX travel time is more than 16 minutes a window of 10 minutes is used. As can be seen in Figure 11, the decreased window size at travel times less than 16 minutes towards the end is still able to capture the observed travel times.

The same cases tested in scenario 1, and 2 apply for scenario 3 as well. All the testing plots for this scenario can be found in Appendix C. Figure 12 and Figure 13 are bar plots showing the percentage of true matches for the morning and afternoon cases. Comparing the results with those of scenario 2, in some cases a significant improvement has been gained. In Figure 13 afternoon case, the percent of match increased from 79% to 86%.



Figure 11: Scenario 3 - Dynamic travel time & 2 conditions time window.



Figure 12: Percent of accurate matches with dynamic search space and varying windows (September 17, 2014 results).



Figure 13: Percent of accurate matches with dynamic search space and varying windows (September 22, 2014 results)

### 5.0 SUMMARY AND CONCLUSIONS

This section summarizes the results obtained after applying the Bayesian re-identification model in 3 different scenarios listed above. In Table 3 and Table 4 the percentage of accurately matched trucks are presented. In Table 3, the time column represents the minimum travel time. A fixed window size of 10 minutes is used. The best performing minimum travel time is highlighted in green. As can be seen from the table, the best performing travel times differ by a considerable amount for each dataset. For example, the best minimum travel time is 13 minutes for the morning rush hour on September 17 whereas it is 22 minutes for September 22. It is clear that a common best time cannot be found. Especially for 22<sup>nd</sup> of September the best travel time for one time period is the worst for the other. This clearly illustrates that a fixed or static search space window will not provide accurate results under the varying traffic conditions.

In Table 4, the accurately matched trucks are presented for scenarios 2 and 3. In scenarios 2 and 3 dynamic travel times obtained from INRIX speed data are used. As explained before, since trucks travel generally slower than the general traffic, the minimum travel times are adjusted by shifting the INRIX travel times downwards by several minutes. The time column in Table 4 represents this amount of time which is referred to as time shift. Here the best performing time shift is highlighted green. In all the cases of both scenarios and both morning and afternoon periods it can be seen that a common time shift of 2 or 3 minutes is performing the best. Using either time shift results in very good match percentages with as low as 2% loss of accuracy from the optimum amount. This shows that dynamic time windows are robust to small perturbations in search space.

In conclusion, in scenario 1 the best performing minimum travel time is varying for each dataset. There is no common minimum travel time (or search space) that can be used for all datasets. However, in scenario 2 and 3, a dynamic travel time produces reasonably high matching accuracy. Since, a common minimum travel time cannot be found for scenario 1 and the optimum minimum travel time is somewhat arbitrary; sensitivity analysis was performed by changing the best performing travel time by 2 minutes. Here the accuracy for each time of the day for the fixed travel time with optimum, 2 minutes less than optimum, and 2 minutes more than optimum alongside the dynamic travel time with 3 minute window shift is analyzed. As it is depicted in Figure 14 there are large differences between the optimum case and the travel times that are close to optimum for the fixed scenarios. This shows that the model which assumes fixed travel time is not robust.

Most importantly it can be seen from the results that dynamic travel time does have its benefits and helps the re-identification algorithm in performing better. The largest accuracies obtained from each scenario tested on each day and time period are presented in Figure 15. There is a gain of up to 19% in accuracy is some periods using dynamic travel time compared to fixed times. These results prove the value of using dynamic travel time for re-identification algorithms. The model also proves to be robust as the time shift from dynamic travel time can be held constant at 2 or 3 minutes with a very small amount of loss in accuracy. Using a dynamic search window together with dynamic travel time has also made some improvements in the model. Compared to

scenario 2, scenario 3 with a dynamic window size performs either the same or better in some cases. For future research, these scenarios can be tested on larger datasets to further demonstrate the value of using dynamic travel times for the search space.

|                 | 9/17/           | 2014       | 9/22/2014  |            |  |  |
|-----------------|-----------------|------------|------------|------------|--|--|
|                 | 7AM-9AM 4PM-6PM |            | 7AM-9AM    | 4PM-6PM    |  |  |
| Travel Time (m) | Scenario 1      | Scenario 1 | Scenario 1 | Scenario 1 |  |  |
| 10              | 35%             | 25%        | 20%        | 74%        |  |  |
| 11              | 43%             | 31%        | 20%        | 77%        |  |  |
| 12              | 49%             | 41%        | 22%        | 77%        |  |  |
| 13              | 58%             | 47%        | 25%        | 67%        |  |  |
| 14              | 52%             | 56%        | 25%        | 30%        |  |  |
| 15              | 49%             | 63%        | 28%        | 28%        |  |  |
| 16 52%          |                 | 63%        | 32%        | 26%        |  |  |
| 17              | 51%             | 63%        | 38%        | 26%        |  |  |
| 18              | 48%             | 66%        | 45%        | 19%        |  |  |
| 19              | 48%             | 66%        | 54%        | 7%         |  |  |
| 20              | 40%             | 59%        | 52%        | 0%         |  |  |
| 21              | 30%             | 53%        | 58%        | 0%         |  |  |
| 22              | 22%             | 47%        | 60%        | 0%         |  |  |
| 23              | 15%             | 44%        | 57%        | 0%         |  |  |
| 24              | 14%             | 28%        | 52%        | 0%         |  |  |
| 25              | 10%             | 28%        | 45%        | 0%         |  |  |

Table 3: Testing results for scenario 1

#### Table 4: Testing results for scenario 2 and scenario 3

|                | 9/17/2014 |            |           |            | 9/22/2014 |            |           |            |
|----------------|-----------|------------|-----------|------------|-----------|------------|-----------|------------|
|                | 7AM-9AM   |            | 4PM-6PM   |            | 7AM-9AM   |            | 4PM-6PM   |            |
| Time Shift (m) | <b>S2</b> | <b>S</b> 3 |
| 1              | 54%       | 56%        | 66%       | 66%        | 57%       | 57%        | 72%       | 79%        |
| 2              | 68%       | 69%        | 78%       | 78%        | 60%       | 58%        | 77%       | 86%        |
| 3              | 73%       | 73%        | 81%       | 81%        | 57%       | 55%        | 79%       | 84%        |
| 4              | 72%       | 70%        | 81%       | 81%        | 49%       | 49%        | 77%       | 72%        |
| 5              | 65%       | 62%        | 69%       | 69%        | 49%       | 49%        | 77%       | 33%        |

Note. S2 = Scenario 2; S3 = Scenario 3.



Figure 14: Sensitivity analysis of fixed TT when the best TT is changed by 2 minutes



Figure 15: Maximum % accurate match for each scenario of all data set

#### 6.0 **REFERENCES**

- CETIN, M., MONSERE, C., NICHOLS, A. & USTUN, I. 2011a. Key Factors Affecting the Accuracy of Reidentification of Trucks over Long Distances Based on Axle Measurement Data. *Transportation Research Record: Journal of the Transportation Research Board*, 2243, 1-8.
- CETIN, M., MONSERE, C. M. & NICHOLS, A. P. 2011b. Bayesian Models for Reidentification of Trucks Over Long Distances on the Basis of Axle Measurement Data. *Journal of Intelligent Transportation Systems*, 15, 1-12.
- CETIN, M. & NICHOLS, A. P. 2009. Improving the accuracy of vehicle re-Identification by solving the assignment problem. *Transportation Research Record: Journal of the Transportation Research Board*, 2129, 1-8.
- CHEN, H. & RAKHA, H. Agent-Based Modeling Approach to Predict Experienced Travel Times. Transportation Research Board 93rd Annual Meeting, 2014.
- COIFMAN, B. & CASSIDY, M. 2002. Vehicle reidentification and travel time measurement on congested freeways. *Transportation Research, Part A (Policy and Practice)*, 36A, 899-917.
- DEMPSTER, A. P., LAIRD, N. M. & RUBIN, D. B. 1977. Maximum likelihood from incomplete data via EM algorithm. *Journal of the Royal Statistical Society Series B-Methodological*, 39, 1-38.
- ELKINS, L. & HIGGINS, C. 2008. Development of Truck Axle Spectra from Oregon Weigh-in-Motion Data for Use in Pavement Design and Analysis. Final Report FHWA-OR-RD-08-06, Oregon Department of Transportation.
- MAZARÉ, P.-E., TOSSAVAINEN, O.-P., BAYEN, A. & WORK, D. Trade-offs between inductive loops and GPS probe vehicles for travel time estimation: A Mobile Century case study. Transportation Research Board 91st Annual Meeting (TRB'12), 2012.
- MCLACHLAN, G. & PEEL, D. 2000. Finite Mixture Models, John Wiley & Sons.
- OH, C., RITCHIE, S. G. & JENG, S. T. 2007. Anonymous vehicle reidentification using heterogeneous detection systems. *Ieee Transactions on Intelligent Transportation Systems*, 8, 460-469.
- SUN, C., RITCHIE, S. G., TSAI, K. & JAYAKRISHNAN, R. 1999. Use of vehicle signature analysis and lexicographic optimization for vehicle reidentification on freeways. *Transportation Research Part C: Emerging Technologies*, 7, 167-185.
- TAWFIK, A. Y., ABDULHAI, B., PENG, A. & TABIB, S. M. 2004. Using decision trees to improve the accuracy of vehicle signature reidentification. *Transportation Research Record*, 24-33.
- TU, H. 2008. Monitoring travel time reliability on freeways, TU Delft, Delft University of Technology.

### 7.0 APPENDICES



### **APPENDIX A**

















### **APPENDIX B**











