

1 **A TRAIN DATA COLLECTION AND ARRIVAL TIME PREDICTION**
2 **SYSTEM FOR HIGHWAY-RAIL GRADE CROSSINGS**

3
4
5 **Yifeng Chen, Ph.D.**

6 (corresponding author)

7 Traffic Engineer, AECOM

8 27777 Franklin Road, Suite 2000

9 Southfield, Michigan, 48034

10 Tel: 402-570-8771

11 Email: chaserchen98@hotmail.com

12
13
14 **Laurence R. Rilett, Ph.D., P.E.**

15 Keith W. Klaasmeyer Chair in Engineering and Technology

16 Director, Nebraska Transportation Center

17 262D Prem S. Paul Research Center at Whittier School

18 2200 Vine Street

19 University of Nebraska–Lincoln

20 Lincoln, NE 68583-0815, P.O. Box 830851

21 Telephone: 402-472-1992

22 Email: lrilett2@unl.edu

23
24
25
26
27 Word Count: 5,728 + 4 Figures + 3 Tables =7,478

28
29
30
31
32 Revision date: November 15, 2016

33

1 ABSTRACT

2 A key component to the safe and efficient operation of traffic signals near highway-rail grade
3 crossings (HRGCs) is an accurate estimate of a given train's arrival time at each crossing. An
4 improvement in the accuracy and timing of the predicted train arrival would allow for
5 improvements in traffic signal preemption algorithms which would, in turn, lead to increased
6 driver and pedestrian safety and/or reduced network delay. This paper introduces a system for
7 train data collection and arrival time prediction at an HRGC using real-time data from a
8 combined video-Doppler radar vehicle detection system. A dual-track test bed system in Lincoln,
9 Nebraska was chosen for train speed data collection. Both kinematic equation-based and multiple
10 linear regression models were developed and used to predict train arrival time at the HRGC. The
11 best models, based on average absolute errors (AAE), were identified. It was found that on
12 average, regression models were more accurate than kinematic models by approximately 13%.
13 As detection time increases, AAEs of both the kinematic and regression models decrease. In
14 addition, the confidence intervals (CI) about mean prediction errors were obtained using a
15 bootstrap method. The average prediction errors of the regression models were approximately 50%
16 smaller than the kinematic models, all else being equal.

17

18

19

1 INTRODUCTION

2 Highway-rail grade crossings (HRGCs) are locations where a roadway and a railroad intersect
3 (1). When a highway-rail grade crossing (HRGC) is located close, e.g., within 61 meters (200
4 feet), to a signalized intersection, the standard practice is to interconnect the railroad signal
5 control equipment and the highway signal control equipment and use a preemption signal control
6 strategy when a train is approaching (1). According to the *Manual on Uniform Traffic Control*
7 *Devices* (MUTCD) (1) and the *AREMA/AAR Signal Manual* (3), a minimum of 20 seconds of
8 warning time should be provided to the traffic control system before a train arrives at a crossing.
9 This allows the railroad crossing warning devices to be activated prior to the arrival of a train
10 and to clear any vehicles that may be stopped on the track(s) as part of the signal preemption
11 strategy. To accomplish this task, a train detection system with train arrival forecasting
12 functionality is required to determine when the railroad warning devices and the highway signal
13 preemption will be activated.

14 The current preemption system uses railroad track circuits to detect arriving trains and to
15 activate railroad warning devices and signal preemption at the intersection near HRGCs
16 (IHRGC). This is known as first generation train detection technology (4, 5). There are two types
17 of first-generation prediction technologies in practice: fixed distance and constant warning time
18 (5, 6). The fixed distance system activates the warning devices whenever a train's presence is
19 detected at a fixed point in the track. The location is calibrated so that the warning devices and
20 signal controller will be notified at least 20 seconds before the fastest possible train arrives at the
21 HRGC. However, the warning time would double for a train operating at half of the maximum
22 speed. The constant warning time (CWT) system measures the train speed upon detection, and
23 uses a simple distance-speed relationship to forecast train arrival time, assuming that the train
24 will travel at this constant speed (4, 5, 6, 7). The resulting warning time algorithm ensures a
25 minimum warning time of 20 seconds, and is relatively consistent compared to the fixed distance
26 system. However, if the trains are accelerating (decelerating) then the predicted arrival time may
27 be in error and traffic signal preemption may start earlier (later) than the ideal preemption
28 warning time (5). For example, a study by L-P Tardif & Associates, Parviainen & Associates,
29 and CANAC International, Inc. showed that the CWT system warning time varied from 20
30 seconds to 38 seconds for a sample of 68 trains traveling at speeds greater than 48 km/h (30
31 mph), with 93% of times falling between 20 and 33 seconds (5).

32 Second generation detection technologies, including acoustic sensors, infrared sensors,
33 Autoscope video detection, and Doppler radar detectors, are now available for train-detection
34 applications. These detectors are non-intrusive devices mounted off the railroad right-of-way and
35 can detect train movements continuously upstream from the railroad circuits. With this setup,
36 trains will be detected considerably earlier than 20 seconds (e.g., 60 seconds). Larger train
37 warning times mean that new advanced preemption strategies, such as transition preemption
38 strategies (TPS1, TPS2 and TPS3) (2), can be used, which allows signal timings to be optimized
39 at an earlier stage of train arrival, and consequently improves safety and efficiency at IHRGCs (2,
40 8). It should be noted that the second train detection technologies are intended to supplement, not
41 replace, the current train preemption system (5).

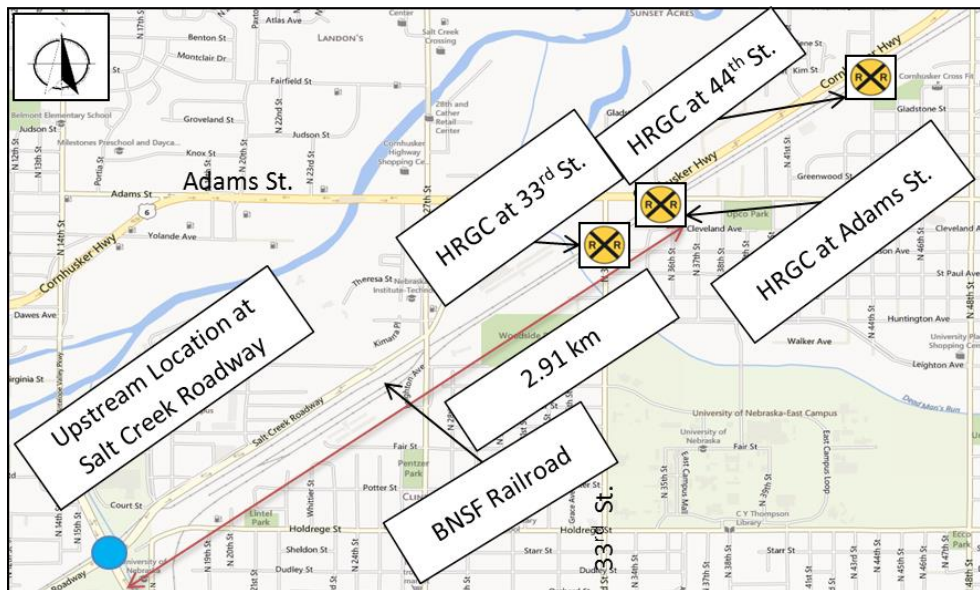
42 Third generation technology is based on on-board computers, Geographic Positioning
43 Systems (GPS), and two-way communication between the train and the signal controller/warning
44 devices. These systems have the ability to provide continuous and more accurate information on
45 the location and speed of the approaching train than either the first or second generation
46 technologies. To date, third-generation technology is mostly in the development or testing phase

1 and has not been widely implemented. Closer cooperation between the railroad and roadway
 2 agencies needs to be established before the technology will have widespread deployment.

3 This paper focuses on using the second-generation train detection technologies (e.g.,
 4 video and radar detectors) to measure train speeds, and subsequently uses this information to
 5 predict train arrival times and their associated prediction intervals at HRGCs. The organization
 6 of this paper is as follows: the next section introduces the test bed and the data collection site for
 7 the train prediction system, followed by the architecture of the train data collection and
 8 prediction system. Then, the prediction models used in the system are described mathematically.
 9 Evaluation and comparison of the models are performed. Next, error bounds of the prediction
 10 models are estimated using a bootstrap method. Lastly, concluding remarks complete the paper.

12 TEST BED

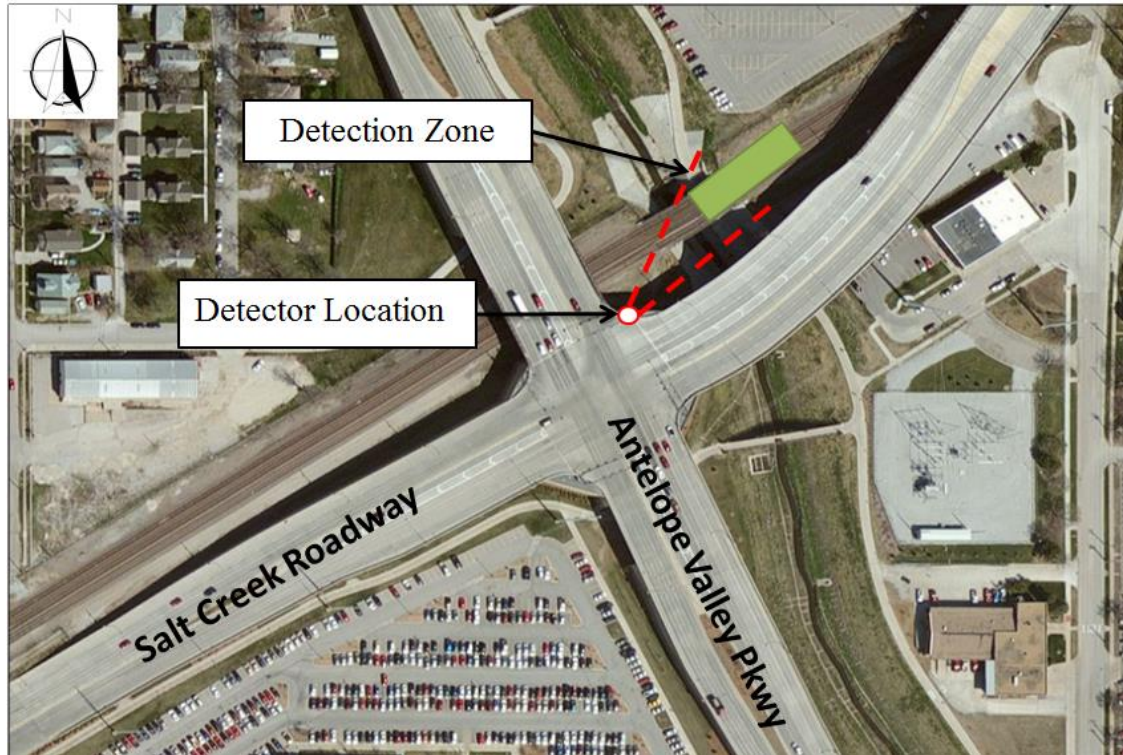
13 The BNSF railroad in Lincoln, Nebraska was chosen as the test bed for the study and is shown in
 14 **FIGURE 1**. It is a dual-track mainline railroad with the tracks running in a northeast-southwest
 15 direction. The railroad intersects the urban roadway network, at grade, at 33rd Street, Adams
 16 Street, and 44th Street. The HRGC at Adams Street was chosen as the target HRGC because: 1)
 17 it has considerable train traffic, on the order of 50-70 trains per day; 2) it has been identified as
 18 one of the most hazardous HRGCs in Lincoln, NE by the FRA, based on the accident history
 19 dating back to 1975(9); and 3) it is near a signalized intersection (35th Street and Cornhusker
 20 Highway), which has a unique geometry and high potential for traffic to back up on to the
 21 railroad tracks (8, 9, 10). As **FIGURE 1** shows, the test bed system also includes an upstream
 22 location located at Salt Creek Roadway and Antelope Valley Parkway, which is approximately
 23 2.91 km (1.81 miles) southwest from the Adams Street HRGC. This upstream location is used
 24 for eastbound (EB) train detection and speed measurement, using two second-generation
 25 detectors: a Doppler radar detector and an Autoscope video detector.



27
 28 **FIGURE 1 Test bed in Lincoln, Nebraska.**

29
 30 **FIGURE 2** shows the layout of the upstream location at Salt Creek Roadway. The
 31 railway and roadway are grade-separated, and the intersection of Salt Creek Roadway and

1 Antelope Valley Parkway is comprised of a bridge that stands approximately 12.19 meters (40
 2 feet) above the railroad tracks. A Doppler radar and a day/night video camera with an Autoscope
 3 video image processor (10) are installed on a traffic signal pole at the northeast corner of the
 4 intersection, as shown in **FIGURE 2**. The Doppler radar and the video camera are aimed at the
 5 same detection zone on the railroad, which is located approximately 213 meters (698 feet) east of
 6 the bridge. This detection zone is represented in **FIGURE 2** as a green rectangle.
 7



8
 9

10 **FIGURE 2** Upstream location at Salt Creek roadway.

11

12 **ARCHITECTURE OF THE PROPOSED SYSTEM**

13 The data collection and arrival time prediction system consists of three parts: (1) online data
 14 collection, (2) offline data processing, and (3) offline detector calibration. **FIGURE 3** illustrates
 15 the architecture of the speed measurement and arrival time prediction system.

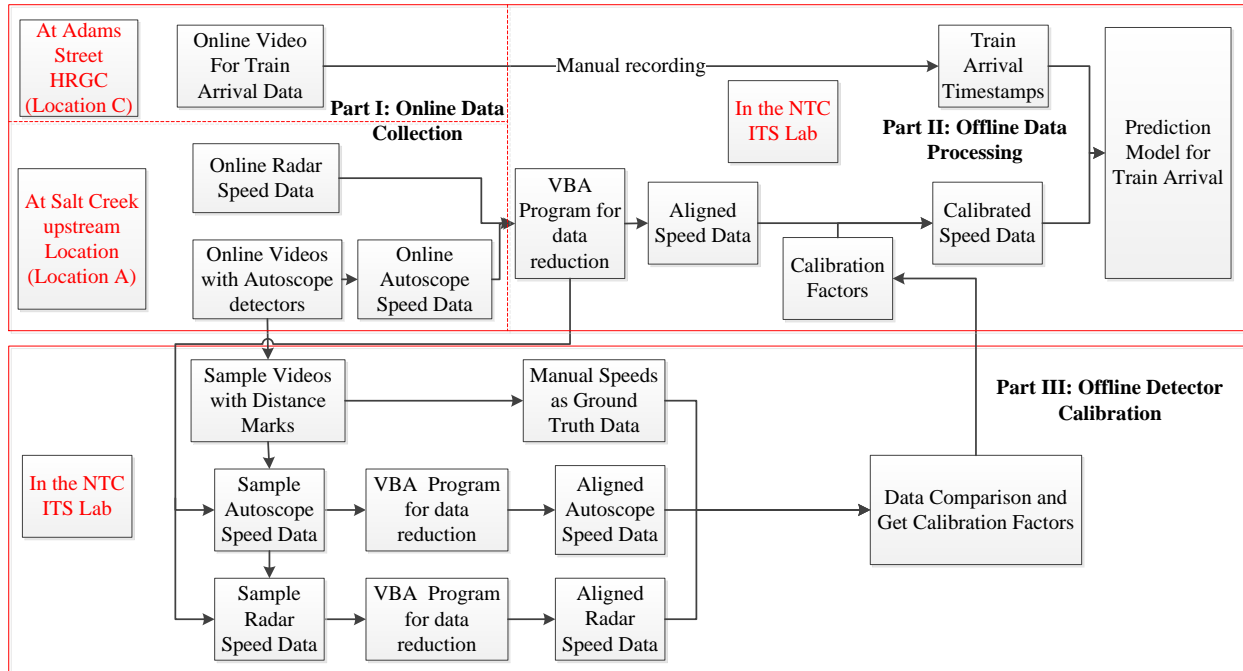


FIGURE 3 Architecture of the data collection and arrival time prediction system.

At the upstream Salt Creek location and the Adams Street HRGC, real-time train videos are recorded on a digital video recorder (DVR) and sent back to the Nebraska Transportation Center (NTC) ITS lab through the City of Lincoln's fiber optic network. Simultaneously, train speed data from radar and Autoscope video detectors are collected at the Salt Creek location and transmitted back to the ITS lab through the same network. This constitutes the online data collection listed as Part I in **FIGURE 3**.

In the offline detector calibration (Part III), train videos are randomly sampled from the videos transmitted back to the NTC ITS lab. Subsequently, the corresponding online speed data from radar and Autoscope speed detectors are identified. The sampled train videos are used to obtain ground truth train speed data. The aligned radar and Autoscope video speed data are subsequently used, along with the ground truth data, to calibrate the detectors. A detailed description of the detector calibration methodology can be found elsewhere (8, 10, 11). The calibration factors obtained offline are then applied to the online data collection.

With regard to offline data processing (Part II), a program for data reduction was developed in Visual Basic for Applications (VBA), which eliminates outliers and false calls in the speed data. The other function of the program is to align speed measurements of the two speed sensors to one-second increments, as these two detectors measure train speed at different resolutions (10).

The aligned speed data were then adjusted using the calibration factors obtained from the offline detector calibration process. The train arrival time stamps at the Adams Street HRGC were recorded manually by observing the videos transmitted from the field. In Part II, the calibrated train speed data and train arrival timestamps were used to develop the train arrival prediction models. It was found that the data from the radar detector were more comprehensive and more robust than those from the Autoscope detector, and consequently, radar data were used

1 in the prediction model development. The videos and Autoscope speed data were used to verify
2 and validate the radar speed data.

3

4 **TRAIN ARRIVAL PREDICTION MODELS**

5 During the data collection process, a total of 138 EB trains were randomly sampled from May
6 2013 through September 2013. Of these, 100 were randomly selected and used for developing
7 the train arrival prediction models. The remaining 38 trains were used to validate the prediction
8 models.

9 Intuitively, in order to provide the most accurate train arrival estimates, the predicted
10 train arrival time should be continuously updated as the train approaches the HRGC (6).
11 Consequently, the model for predicting train arrival time should also be a continuous model that
12 is updated at specific time intervals as the train approaches the HRGC (6). Previous studies (4, 6)
13 have found that it is advantageous to predict train arrival time using models that are updated
14 every 10 seconds. Therefore, it was decided to use the speed reading from every 10 seconds to
15 develop the prediction models.

16 The average absolute error (AAE) at the end of every 10-second interval is used as the
17 performance index of prediction models. It is calculated as shown in Equation 1.

18

$$AAE_t = \frac{\sum_{x=1}^{N_t} |T_t^x - P_t^x|}{N_t} \quad \forall t = 10, 20, \dots, 300 \quad (1)$$

where:

AAE_t average absolute error for prediction time t

T_t^x actual train arrival time for train x at prediction time t

P_t^x predicted train arrival time for train x at prediction time t

N_t number of trains in detection at prediction time t

19 There are two categories of models used in the analysis. The first category contains
20 kinematic models that are based on first or second order motion equations and can be used
21 directly. The second category consists of regression models that attempt to identify a statistical
22 relationship between train speed measurements and the train's arrival time at the HRGC. This
23 category of models requires extensive datasets so that the models can be calibrated and validated.
24 In the next section, both categories of models are presented.

25

26 **Kinematic Models (Model 1- Model 6)**

27

28 *Model 1 Basic train prediction method*

29

30 This model uses the instantaneous speed of a train to calculate the arrival time at the target
31 downstream crossing. Equation 2 and Equation 3 show the calculation of the remaining

- 1 estimated distance d_t^x and the predicted arrival time p_t^x for train x for every 10-second interval.
 2 This approach assumes that the train will travel at the current speed until it arrives at the crossing.
 3 This model is the most straightforward and is based on the CWT forecasting procedure (2).
 4 However, if the train changes speed over time, this model will result in fairly large errors (4, 5).
 5

$$d_t^x = D - \sum_{i=1}^t v_i^x \quad \forall t = 10, 20, \dots, 300 \quad (2)$$

$$p_t^x = \frac{d_t^x}{v_t^x} \quad \forall t = 10, 20, \dots, 300 \quad (3)$$

where:

- d_t^x estimated distance from the head of train x to the downstream grade crossing t seconds after detection (m)
- D distance between the upstream detection location and the HRGC. For this testbed D is equal to 2910 m
- p_t^x predicted arrival time for train x, t seconds after detection (s)
- v_t^x train speed t seconds after detection (m/s)
- v_i^x train speed i seconds after detection (m/s), $i = 1, 2, 3 \dots t$

6
 7 *Model 2 Basic prediction method with average train speeds*

- 8
 9 Model 2 is a similar model to Model 1 and uses a distance-speed equation to predict train arrival
 10 times. The difference is that in Model 2, average train speed of train x at t seconds in detection
 11 (the end of the current 10-second interval) is used instead of the instantaneous speed at t seconds
 12 in detection, as shown in Equation 5.
 13

$$p_t^x = \frac{d_t^x}{\overline{v_t^x}} \quad \forall t = 10, 20, \dots, 300 \quad (4)$$

$$\overline{v_t^x} = \frac{\sum_{i=1}^t v_i^x}{t} \quad \forall t = 10, 20, \dots, 300 \quad (5)$$

where:

\overline{v}_t^x average speed of train x, t seconds after detection (m/s)

1 *Model 3 Basic prediction method with moving average train speeds*

2

3 In Model 3, the moving average speed over the current 10-second interval is used, as shown in
4 Equation 6 and Equation 7.

5

$$p_t^x = \frac{d_t^x}{\overline{v}_t^x} \quad \forall t = 10 \text{ s}, 20 \text{ s}, \dots, 300 \text{ s} \quad (6)$$

$$\widehat{v}_t^x = \frac{\sum_{j=1}^9 v_{t-j}^x}{10} \quad \forall t = 10 \text{ s}, 20 \text{ s}, \dots, 300 \text{ s} \quad (7)$$

where:

\widehat{v}_t^x average train speed for train x, t seconds after detection (m/s)

v_{t-j}^x train speed j seconds ahead of v_t^x , $j=1,2,\dots,9$

6 *Model 4 Train prediction method using second order kinematic equation and instant speed*

7

8 This model calculates the arrival time by using the kinematic equation with the moving average
9 acceleration rate, instantaneous speed of a train, and remaining distance to the crossing, as shown
10 in Equation 8 and Equation 10. The moving average acceleration is calculated as the speed
11 difference of the two speeds by the end of the current and last 10-second intervals divided by 10
12 seconds, as shown in Equation 9. The principle behind this model is that by including
13 information on whether the train is increasing or decreasing its speed, the resulting model would
14 be more accurate.

15

$$d_t^x = \frac{\widehat{a}_t^x * (p_t^x)^2}{2} + v_t^x * p_t^x \quad \forall t = 10, 20, \dots, 300 \quad (8)$$

$$\widehat{a}_t^x = \frac{v_t^x - v_{t-10}^x}{10} \quad \forall t = 10, 20, \dots, 300 \quad (9)$$

$$p_t^x = \frac{-v_t^x + \sqrt{v_t^{x2} + 2 * \widehat{a}_t^x * d_t^x}}{\widehat{a}_t^x} \quad \text{when } \widehat{a}_t^x > 0 \quad \forall t = 10, 20, \dots, 300 \quad (10)$$

$$p_t^x = \frac{d_t^x}{v_t^x} \quad \text{when } \widehat{a}_t^x \leq 0$$

where:

\widehat{a}_t^x moving average acceleration/deceleration rate of train x, t seconds after detection (m/s^2)

1 The limitation of this model is that the one part in the square root (i.e. $v_t^{x2} + 2 * \widehat{a}_t^x * d_t^x$)
 2 in Equation 10 may yield a negative value when \widehat{a}_t^x is negative. Therefore, Equation 10 will
 3 work only when the acceleration rate is positive. As an alternative, \widehat{a}_t^x is set to 0 when it is
 4 negative in Equation 9. In this case, Equation 10 reduces to Equation 3, and consequently Model
 5 4 has the same results as Model 1 when the estimated acceleration is zero or negative.

6
 7 *Model 5 Train prediction method using second order kinematic equation and average speed*

8
 9 Model 5 uses average speed in place of instantaneous speed, as shown in Equation 11. Similar to
 10 Model 4, the value of \widehat{a}_t^x is set to 0 if $\widehat{a}_t^x \leq 0$. In this situation, Model 5 is equivalent to Model 2.
 11

$$p_t^x = \frac{-\overline{v}_t^x + \sqrt{\overline{v}_t^{x2} + 2 * \widehat{a}_t^x * d_t^x}}{\widehat{a}_t^x} \quad \text{when } \widehat{a}_t^x > 0 \quad \forall t = 10, 20, \dots, 300 \quad (11)$$

$$p_t^x = \frac{d_t^x}{\overline{v}_t^x} \quad \text{when } \widehat{a}_t^x \leq 0$$

12
 13 *Model 6 Train prediction method using second order kinematic equation and moving average*
 14 *speed*

15
 16 Model 6 uses moving average speeds in place of instantaneous speed, as Equation 12 shows.
 17 Similar to Model 4 and Model 5, the value of \widehat{a}_t^x is set to 0 when $\widehat{a}_t^x \leq 0$. In this situation, Model
 18 6 is equivalent to Model 3.
 19

$$p_t^x = \frac{-\overline{v}_t^x + \sqrt{\overline{v}_t^{x2} + 2 * \widehat{a}_t^x * d_t^x}}{\widehat{a}_t^x} \quad \text{when } \widehat{a}_t^x > 0 \quad \forall t = 10, 20, \dots, 300 \quad (12)$$

$$p_t^x = \frac{d_t^x}{\overline{v}_t^x} \quad \text{when } \widehat{a}_t^x \leq 0$$

1 The speed data from 100 trains were used to test the above six model types. **TABLE 1**
 2 lists the AAEs of these models at each time in detection (time during which a train dwells in the
 3 detection zone) from 10 seconds to 300 seconds. The number of trains present in each interval is
 4 also included. It may be seen that the number of trains decreases as the time in detection
 5 increases because some trains will have arrived at the crossing. Obviously, once this occurs it is
 6 not necessary to include them in the calibration.

7 In **TABLE 1**, the best prediction models (e.g., the ones with the lowest AAE) for each
 8 prediction interval are marked with an asterisk. Generally, the AAEs for each model type
 9 decreased as time in detection increased. It can be seen that the second order models (Model 4,
 10 Model 5, and Model 6) have lower AAEs than the first order models (Model 1, Model 2, and
 11 Model 3) for all the detection intervals except for the first one. Thus, it was concluded that the
 12 second order models are better than the first order models for predicting train arrival times. It can
 13 also be seen from **TABLE 1** that Model 4 using instantaneous speed, and Model 6 using moving
 14 average speed, yield the lowest AAEs in 27 out of the 30 prediction intervals. Model 6 yields the
 15 lowest AAEs during most of the detection intervals between 20 seconds and 150 seconds.
 16 Although Model 4 yields most of the lowest AAEs at detection times between 160 seconds to
 17 300 seconds, the differences in AAEs between Model 4 and Model 6 are less than 1 second. If
 18 only one model were to be used, for consistency, then Model 6 is recommended. Note that
 19 Model 6 is 30% better, on average, than Model 1, which is based on current practice.
 20

21 **TABLE 1 Kinematic Model AAE Results**
 22

Time in Detection k (s)	Model1	Model2	Model3	Model4	Model5	Model6	# of Trains
10	138.7	138.7	138.7*	172.4	172.4	172.4	100
20	111.7	114.8	114.8	76.6	75.8	75.8*	100
30	107.4	111.7	108.9	73.5	71.6*	71.6*	100
40	103.3	109.2	105.1	69.0	68.0*	68.2	100
50	99.8	106.9	101.3	70.7	69.6*	69.6*	100
60	94.7	104.4	96.7	65.7	66.0	64.7*	100
70	90.6	101.8	92.5	55.4	56.7	55.3*	99
80	86.8	99.2	88.2	55.4	56.5	55.0*	99
90	83.2	96.6	84.6	54.9	57.8	54.3*	99
100	80.7	95.0	81.5	54.6	57.4	53.9*	97
110	77.5	92.6	78.1	49.2	54.1	48.9*	97
120	73.7	90.5	75.4	43.3	46.8	43.3*	95
130	71.1	88.2	72.1	49.8	55.1	49.6*	95
140	68.1	85.7	69.2	45.3	51.6	45.1*	94
150	65.7	83.7	67.1	46.0	52.0	45.9*	93
160	64.4	81.5	65.5	43.1*	50.8	43.2	90
170	62.3	79.6	63.4	40.3*	47.4	40.5	89
180	60.7	78.2	61.7	46.0*	54.5	46.1	87
190	59.1	77.4	60.5	43.3*	51.7	43.6	84
200	57.0	75.6	58.5	41.4*	50.9	41.7	83
210	59.3	78.5	60.0	42.6*	55.2	42.9	74

220	59.3	79.7	60.4	41.6*	52.5	41.9	68
230	63.7	86.7	64.2	47.3	61.7	47.2*	56
240	64.4	88.9	65.4	44.5*	59.4	44.5*	51
250	66.2	93.0	67.7	43.6*	57.6	44.3	45
260	69.9	97.3	70.9	52.8*	70.7	53.2	39
270	71.9	102.4	73.1	45.5*	62.5	46.0	34
280	73.3	108.8	73.9	51.6	69.2	51.2*	30
290	75.9	115.6	78.1	42.2*	62.6	43.1	25
300	75.0	118.7	77.5	47.2*	73.4	48.3	23
Average	77.8	96.9	79.1	55.0	63.5	54.9	

* Best prediction model for the given interval

Regression Models (Model 7-Model 21)

Models 7 through 21 are regression models. As in the kinematic model section, the same 100 trains were used as training data for model calibration, and the remaining 38 trains were used for model validation.

Model 7, Model 8, and Model 9 use speed as the only independent variable, as shown in Equations 13, 14, and 15. Similar to Model 1 through Model 3, instant speed, average speed, and moving average speed were used to represent speed in Model 7, Model 8, and Model 9, respectively.

Model 7

$$p_t^x = c_1 + v_t^x * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (13)$$

Model 8

$$p_t^x = c_1 + \overline{v_t^x} * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (14)$$

Model 9

$$p_t^x = c_1 + \widehat{v_t^x} * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (15)$$

Model 10 through Model 12 include speed and squared speed as independent variables. The squared speed variable was chosen in order to examine if there were any second order relationships between speed and arrival time.

Model 10

$$p_t^x = c_1 + v_t^x * c_2 + (v_t^x)^2 * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (16)$$

Model 11

$$p_t^x = c_1 + \overline{v_t^x} * c_2 + \overline{v_t^x}^2 * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (17)$$

1 **Model 12**

$$p_t^x = c_1 + \widehat{v}_t^x * c_2 + \widehat{v}_t^{x^2} * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (18)$$

2 Models 13 through 15 use speed, squared speed, and estimated remaining distance to the
 3 crossing as independent variables, as shown in Equations 19 to 21. As the train approaches the
 4 crossing, arrival time should decrease given that train speed does not change dramatically.
 5 Therefore, it is hypothesized that the estimated remaining distance is negatively correlated to
 6 arrival time. As before, the three different speed definitions were used in these models.

7
8 **Model 13**

$$p_t^x = c_1 + v_t^x * c_2 + v_t^{x^2} * c_3 + d_t^x * c_4 \quad \forall t = 10, 20, \dots, 300 \quad (19)$$

9 **Model 14**

$$p_t^x = c_1 + \overline{v}_t^x * c_2 + \overline{v}_t^{x^2} * c_3 + d_t^x * c_4 \quad \forall t = 10, 20, \dots, 300 \quad (20)$$

10 **Model 15**

$$p_t^x = c_1 + \widehat{v}_t^x * c_2 + \widehat{v}_t^{x^2} * c_3 + d_t^x * c_4 \quad \forall t = 10, 20, \dots, 300 \quad (21)$$

11
 12 Model 16 through Model 21 are natural logarithm-based models. Specifically, Model 16,
 13 Model 17, and Model 18 use the natural logarithm of speed as independent variables, while
 14 Models 19 through 21 also include the logarithm of the estimated remaining distance to the
 15 HRGC. The natural logarithm was chosen to see if there was a log-linear relationship between
 16 the independent variables and the arrival time.

17
18 **Model 16**

$$\ln(p_t^x) = c_1 + \ln(v_t^x) * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (22)$$

19 **Model 17**

$$\ln(p_t^x) = c_1 + \ln(\overline{v}_t^x) * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (23)$$

20 **Model 18**

$$\ln(p_t^x) = c_1 + \ln(\widehat{v}_t^x) * c_2 \quad \forall t = 10, 20, \dots, 300 \quad (24)$$

21 **Model 19**

$$\ln(p_t^x) = c_1 + \ln(v_t^x) * c_2 + \ln(d_t^x) * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (25)$$

22

23

1 **Model 20**

$$\text{Ln}(p_t^x) = c_1 + \text{Ln}(\overline{v_t^x}) * c_2 + \text{Ln}(d_t^x) * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (26)$$

2 **Model 21**

$$\text{Ln}(p_t^x) = c_1 + \text{Ln}(\widehat{v_t^x}) * c_2 + \text{Ln}(d_t^x) * c_3 \quad \forall t = 10, 20, \dots, 300 \quad (27)$$

where:

c_1 intercept coefficient

c_2, c_3 coefficients for the independent variables

3 The overall results for each 10-second time period are shown in **TABLE 2**. This table
4 shows the “best” models, which are defined as the models with the lowest AAE for the given ten
5 second detection interval, for each of the 30 detection intervals. The equations, the AAEs, and
6 the coefficient of determination (i.e. R^2) value of these models are shown in **TABLE 2** as well.

7 It can be seen that the best models in each interval are either natural logarithm-based
8 models or second order models. This indicates that the relationship between measured speed and
9 arrival time is non-linear. This would be expected because, in general, as the trains leave Lincoln,
10 they are allowed to travel at a faster speed and this section of track will experience trains with
11 fairly high acceleration rates. It was found that as the time in detection increases, the coefficient
12 of determination of the best model increases, and the AAE decreases slightly. In other words, as
13 more information is gathered about the train, the model accuracy improves. This makes intuitive
14 sense because as the train travels closer to the HRGC, it would be expected that model accuracy
15 would increase.

16
17 **TABLE 2 Best Regression Model AAE Results**

Time in detection (s)	Best Regression Model	Equations of the Regression Models *	AAE (s)	R^2
10	Model 19	$\text{Ln}(p)=104.78-1.02*\text{Ln}(v)-12.21\text{Ln}(d)$	33.6	0.64
20	Model 19	$\text{Ln}(p)=82.22-1.18*\text{Ln}(v)-9.36*\text{Ln}(d)$	32.6	0.65
30	Model 19	$\text{Ln}(p)=68.32-1.31*\text{Ln}(v)-7.61*\text{Ln}(d)$	31.7	0.66
40	Model 21	$\text{Ln}(p)=55.51-1.37*\text{Ln}(\tilde{v})-5.995*\text{Ln}(d)$	32.0	0.67
50	Model 19	$\text{Ln}(p)=41.63-1.336*\text{Ln}(v)-4.26*\text{Ln}(d)$	30.3	0.70
60	Model 19	$\text{Ln}(p)=34.33-1.35*\text{Ln}(v)-3.34*\text{Ln}(d)$	30.0	0.71
70	Model 21	$\text{Ln}(p)=26.948-1.31*\text{Ln}(\tilde{v})-2.42*\text{Ln}(d)$	30.0	0.71
80	Model 15	$p=1864.6-121.52*\tilde{v}+3.77*\tilde{v}^2-0.39*d$	31.2	0.72
90	Model 16	$\text{Ln}(p)=7.24-0.88*\text{Ln}(v)$	30.9	0.72
100	Model 16	$\text{Ln}(p)=7.28-0.92*\text{Ln}(v)$	30.4	0.73
110	Model 16	$\text{Ln}(p)=7.35-0.97*\text{Ln}(v)$	29.8	0.73
120	Model 16	$\text{Ln}(p)=7.39-1.01*\text{Ln}(v)$	29.7	0.74

130	Model 18	$\text{Ln}(p)=7.46-1.06*\text{Ln}(\tilde{v})$	29.2	0.75
140	Model 18	$\text{Ln}(p)=7.51-1.10*\text{Ln}(\tilde{v})$	28.7	0.75
150	Model 14	$p=5259.86-324.38*\bar{u}+3.18*\bar{u}^2-1.58*d$	27.6	0.70
160	Model 14	$p=3845.27-265.26*\bar{u}+3.44*\bar{u}^2-1.10*d$	28.1	0.70
170	Model 14	$p=4039.1-286.78*\bar{u}+3.34*\bar{u}^2-1.17*d$	28.1	0.71
180	Model 14	$p=6333.22-434.97*\bar{u}+2.73*\bar{u}^2-1.96*d$	26.3	0.72
190	Model 16	$\text{Ln}(p)=8.01-1.46*\text{Ln}(v)$	28.7	0.78
200	Model 12	$p=652.76-88.87*\tilde{v}+3.32*\tilde{v}^2$	28.0	0.75
210	Model 10	$p=657.203-90.94*v+3.42*v^2$	27.6	0.76
220	Model 10	$p=657.19-90.34*v+3.46*v^2$	29.0	0.74
230	Model 10	$p=670.546-98.64*v+3.96*v^2$	28.8	0.74
240	Model 10	$p=656.23-95.90*v+3.79*v^2$	29.8	0.76
250	Model 10	$p=643.91-94.03*v+3.67*v^2$	31.5	0.75
260	Model 10	$p=639.29-93.45*v+3.60*v^2$	32.1	0.76
270	Model 15	$p=1150.6-147.13*\tilde{v}+5.20*\tilde{v}^2-0.17*d$	30.1	0.84
280	Model 15	$p=1062-138.97*\tilde{v}+5.0*\tilde{v}^2-0.15*d$	32.7	0.82
290	Model 7	$p=491.46-47.71*v$	33.2	0.75
300	Model 16	$\text{Ln}(p)=8.0-1.60*\text{Ln}(v)$	29.6	0.81

* v: instantaneous speed, \bar{u} : average speed, \tilde{v} : moving average speed

Model Comparison and Validation

TABLE 3 shows the validation results of the regression models and the kinematic models using the data set of 38 trains. In 22 of the 30 intervals, the regression models have lower AAEs than the kinematic models, indicating that overall the best regression models are more accurate than the best kinematic models in train arrival time prediction. On average, the AAE values for the regression models are 13% less than the AAE values for the kinematic models. The regression models in the 27th through 30th intervals have higher AAEs than the kinematic models when using the validation data. It is hypothesized that this is because the decrease in the number of observations in these latter intervals decreases the ability of the regression models to capture the wide variety of train behavior.

TABLE 3 Model Validation AAEs

Time in detection (s)	AAE (s)		Difference (%)	# of Trains
	Best Kinematic Model	Best Regression Model	Kinematic – Regression Regression	
10	101.6	47.3	53.5%	38
20	58.6	48.3	17.5%	38
30	62.2	50.4	18.9%	38
40	54.1	47.6	11.9%	38
50	58.1	49.7	14.5%	38
60	52.9	50.3	4.9%	38
70	49.7	50.1	-0.9%	38

80	53.1	46.3	12.7%	38
90	52.4	46.1	11.9%	38
100	48.7	46.8	3.8%	38
110	46.7	45.7	2.0%	38
120	39.2	42.7	-8.9%	37
130	47.8	41.8	12.5%	36
140	45.8	42.0	8.3%	36
150	39.3	34.1	13.3%	36
160	43.6	37.3	14.4%	35
170	45.4	35.4	21.9%	34
180	45.8	30.6	33.2%	33
190	42.7	43.1	-1.0%	31
200	44.6	47.6	-6.6%	31
210	46.4	45.1	3.0%	26
220	48.5	34.6	28.5%	23
230	48.6	39.9	17.8%	23
240	50.2	36.5	27.3%	20
250	51.6	29.5	42.9%	19
260	50.6	30.3	40.0%	18
270	45.7	44.9	1.6%	13
280	41.2	52.9	-28.5%	10
290	40.1	51.1	-27.5%	8
300	43.5	54.0	-24.3%	7
Average	49.9	43.4	13.1%	

1

2 **ESTIMATION OF PREDICTION ERROR INTERVALS**

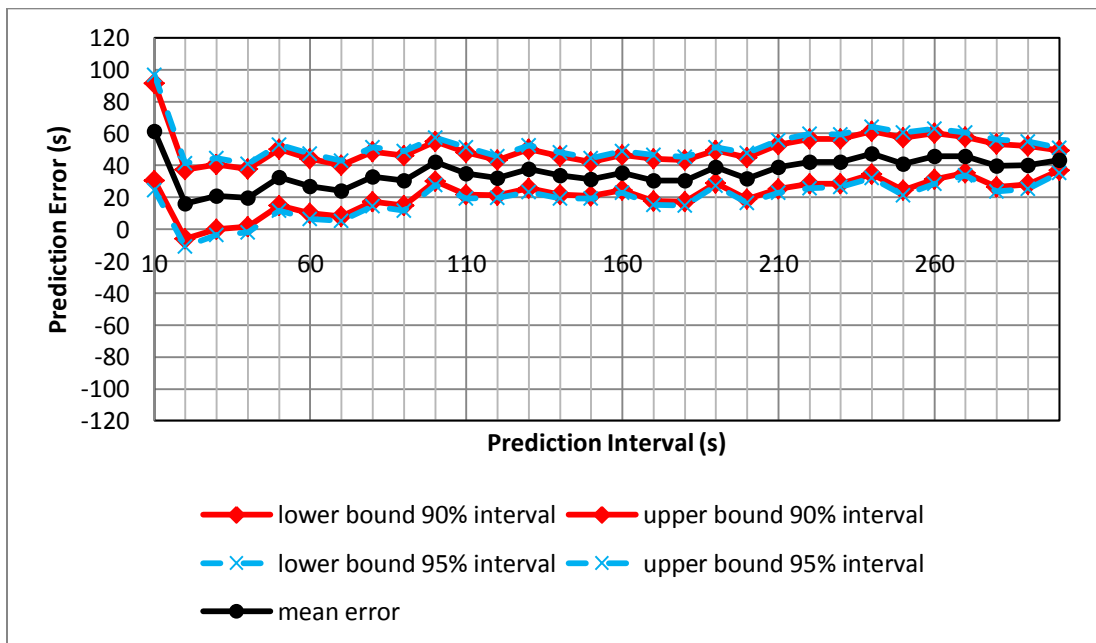
3 In addition to the earlier predicted train arrival time, the confidence interval (CI) about the
4 predicted train arrival is also very useful in the advanced traffic signal preemption algorithms to
5 adjust the predicted arrival times and mitigate the prediction errors (2). Both standard statistical
6 methods (12) and bootstrap methods (13) can be used to estimate confidence intervals. The
7 former are used for multiple linear regression (MLR) modeling because the theory and equations
8 are well known. For non-linear models, such as the kinematic models and the local polynomial
9 models (14), standard statistical methods do not apply because they rely on prior information
10 related to the models and parameter distributions (12). The bootstrap method is a computer-based
11 numerical method for statistical inference, especially for statistics such as sample median and
12 sample correlation, which are difficult to infer using an analytical approach (12, 13). Therefore,
13 the bootstrap method was used to calculate the confidence intervals of mean prediction errors for
14 both kinematic models and regression models.

15 **FIGURE 4 a)** illustrates the mean prediction error (actual arrival time minus predicted
16 arrival time) and the bootstrap CIs at 90 and 95 percentiles from the prediction using the best
17 kinematic models, while **FIGURE 4 b)** illustrates those from the prediction using the best
18 regression models. It may be seen that both the kinematic models and regression models can
19 predict the train arrival time to within approximately ± 60 s of the true arrival time for most

1 prediction intervals. This accuracy is adequate for applying advanced preemption algorithm for
 2 an average train arrival time of 200 seconds and above (6). On average, the regression models
 3 had a prediction error of 17.3 seconds, which is 51.4% lower than that of the kinematic models
 4 (35.6 seconds).

5 However, the performance of the regression models degrades after 260 seconds in
 6 detection, while the performance of the kinematic models remains stable. It is hypothesized that
 7 this is due to the small size of the training data sample in the regression models. It is suggested
 8 that the best regression models be used prior to 260 seconds in detection. After 260 seconds of
 9 detection, either the best kinematic models can be used or the last predicted arrival time can be
 10 used as the final estimate. As the train data set increases, this range may vary. For the application,
 11 the analysis can be redone at set intervals (e.g., every few months), so that the fittest models can
 12 be applied.

13

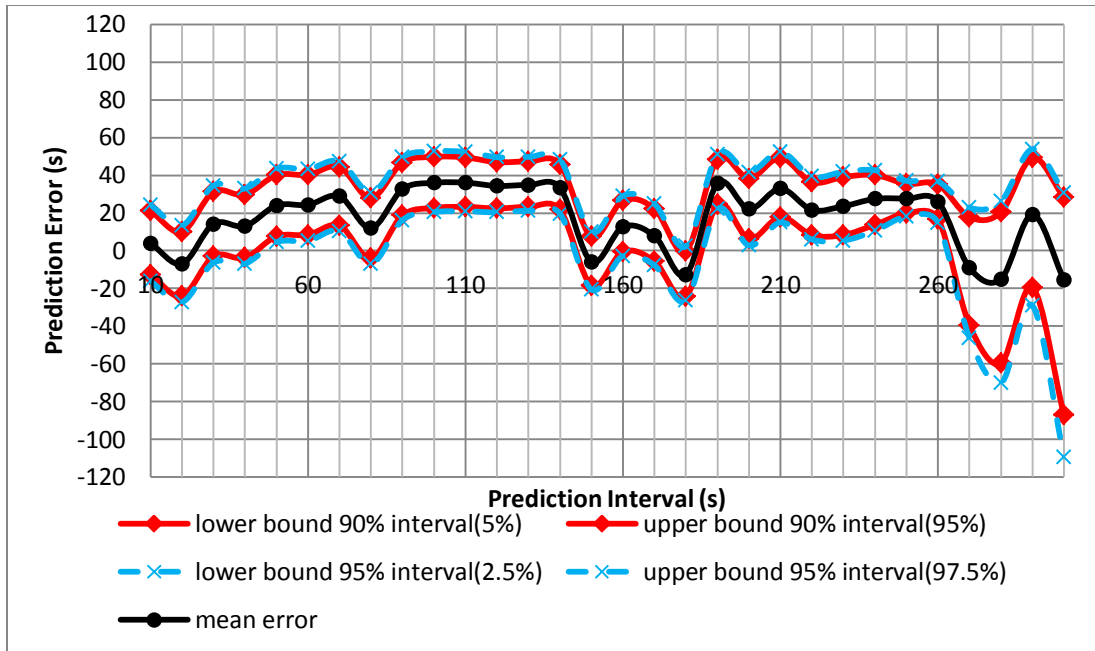


14

15

16

a) Kinematic models



b) Regression models

FIGURE 4 Mean errors and the bootstrap CIs.

CONCLUDING REMARKS

This paper developed a system for train speed measurements and train arrival time prediction at HRGCs using second-generation train detection technologies and competitive kinematic and regression models. On average, the regression models were proven to be 13% more accurate than the kinematic models based on the validation data. The bootstrap method was used to estimate the error bounds of the mean estimates of the models, which, along with the predicted arrival times, are useful to new advanced traffic signal preemption algorithms (2, 8). It was found that the average prediction error for the regression models was approximately 50% smaller than the kinematic models, all else being equal.

With respect to the practicality of the models, the regression approach needs sufficient data to calibrate and validate the models at regular intervals. The advantage of the kinematic models is that they do not require calibration and are ready to implement in the field. Intuitively, they would be especially useful at locations with low train volumes and/or locations where train speeds do not vary much. The trade-off between the accuracy and ease-of-use of the different model categories would be application dependent, and the transportation agencies would have to decide which best fits their needs. In situations where there is not enough sample data, a combined use of kinematic and regression models would be a solution. It should also be pointed out that although the prediction models in the paper are developed using only a freight train sample, the models can be extended to passenger and/or light rail systems.

Future work will focus on streamlining the components of this system and integrating the system in new advanced traffic signal preemption algorithms (2, 8). In this paper, the methodology was used to predict train arrival times for EB trains at a single HRGC. This methodology can be extended to predict train arrival times for both EB and WB trains at multiple HRGCs (e.g., the 33rd Street, Adams Street, and 44th Street HRGC along the test bed).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

1 **ACKNOWLEDGEMENTS**

2 The authors would like to thank the U.S. Department of Transportation Office of the Assistant
3 Secretary for Research and Technology and the University Transportation Center for Railway
4 Safety for their sponsorship and support of this research. The contents of this paper reflect the
5 views of the authors, who are responsible for the facts and the accuracy of the information
6 presented herein, and are not necessarily representative of the sponsoring agency.

7

1 REFERENCES

- 2
- 3 1. FHWA. *Manual on Uniform Traffic Control Devices: 2009 Edition*, U.S. Department of
- 4 Transportation, 2009.
- 5 2. Cho, H. Preemption Strategy for Intersections near Highway-Railroad Crossings. Ph.D.
- 6 Dissertation, Texas A&M University, College Station, TX, 2003.
- 7 3. Association of American Railroads. *Communication and Signal Section, Signal Manual of*
- 8 *Recommended Practice*, Vol. 1, Section 3, Washington, D.C., 1995.
- 9 4. Cho, H. and Rilett, L.R. Forecasting train travel times at At-Grade Crossings, In:
- 10 *Transportation Research Record, No.1844*, TRB, Washington, D.C., 2003, pp. 94-102.
- 11 5. Korve, H.W. *Traffic Signal Operations Near Highway-Railroad Grade Crossings*,
- 12 NCHRP Synthesis 271, TRB, National Research Council, Washington, D.C., 1999.
- 13 6. Esters, R.M. and Rilett, L.R. Advanced Prediction of Train Arrival and Crossing Times at
- 14 Highway-railroad Grade Crossings, In *Transportation Research Record, No.1844*,
- 15 Transportation Research Board, Washington, D.C. 2003, pp. 94-102.
- 16 7. Richards, S.H., K.W. Heathington, and D.B. Fambro. Evaluation of Constant Warning
- 17 Times Using Train Predictors at a Grade Crossing with Flashing Light Signals. In
- 18 *Transportation Research Record 1254*, TRB, National Research Council, Washington,
- 19 D.C., 1990, pp. 60-71.
- 20 8. Chen, Y. An Adaptive Corridor-wide Signal Timing Optimization Methodology for
- 21 Traffic Networks with Multiple Highway-Rail Grade Crossings. PhD Dissertation.
- 22 Department of Civil Engineering, University of Nebraska-Lincoln, 2015.
- 23 9. Jones, Elizabeth G., Khattak, Aemal H., and Rilett, L.R. The University Nebraska-
- 24 Lincoln's Highway Rail Grade Crossing Test Bed System. In *Transportation Research*
- 25 *Record: Journal of the Transportation Research Board*. CD-ROM. TRB of the National
- 26 Academics, Washington, D.C., 2009.
- 27 10. Franca, D. Estimation Train Arrival Times at Highway-Railroad Grade Crossings Using
- 28 Multiple Sensors, Master Thesis, University of Nebraska, Lincoln, NE, 2009.
- 29 11. Chen, Y.F. and Rilett, L. R. A train speed measurement and arrival time prediction system
- 30 for highway-rail grade crossings (HRGCs). *Proceedings of 21ST World Congress on*
- 31 *Intelligent Transport System*, Detroit, Michigan, September 8, 2014.
- 32 12. Spiegelman, Clifford H., Park, Eun Sug, and Rilett, L.R.: *Transportation Statistics and*
- 33 *Microsimulation*. CRC Press, 2011.
- 34 13. Efron, B. and Tibshirani, R.J. *An Introduction to Bootstrap*. Chapman and Hall, New York,
- 35 1993.
- 36 14. Hang Yue, E. Jones, and Peter Revesz. Local Polynomial Regression Models for Average
- 37 Traffic Speed Estimation and Forecasting in Linear Constraint Databases. *17th*
- 38 *International Symposium on Temporal Representation and Reasoning (TIME'10)*, IEEE,
- 39 Paris, France, September, 2010.
- 40