

Automated Vehicle Identification Based on Car-following Data with Machine Learning

Qianwen Li, Xiaopeng Li, Handong Yao, Zhaohui Liang, and Weijun Xie

Abstract— Vehicles with adaptive cruise control, i.e., SAE Levels 1 and 2 automated vehicles (AVs), have been operating on roads with a significant and rapidly growing penetration rate. Identifying these AVs is critical to understanding near-future mixed traffic characteristics and managing highway mobility and safety. This study identifies adaptive cruise control-equipped vehicles from human-driven vehicles (HVs) by constructing a set of learning-based models using car-following trajectories in a short time window. It is extendible to Level 3 and + AV identification when data is available. To compare model performance and draw physical insights, two physics-based models are proposed based on the premise that, in general, the car-following behavior of an AV is less volatile than an HV. Four car-following datasets, including AV makes from different manufacturers, are mixed to build a comprehensive identification model. Results show that physics-based approaches identify more than 80% AVs and 70% HVs. The identification accuracy of learning-based models is even higher. For example, the cluster-aware long short-term memory network identifies 98.79% of AVs and 95.45% of HVs. Learning-based identification models developed by this study can be integrated with the existing infrastructure (e.g., surveillance cameras), which have been used to extract car-following trajectories, to detect AVs in mixed traffic streams. This opens unparalleled data-driven opportunities to analyze and control mixed traffic to enhance safety (e.g., notifying surrounding traffic of the presence of AVs) and mobility (e.g., opening AV dedicated lanes when the percentage is great enough).

Index Terms— Automated vehicle identification, adaptive cruise control, machine learning model, physics-based model, car following.

I. INTRODUCTION

Automated vehicle (AV) technology is expected to enhance traffic safety, elevate roadway capacity, reduce fuel consumption, and mitigate congestion [1]–[5]. Around 10% of total vehicles sold in the second quarter of 2019 were commercial AVs, e.g., those with adaptive cruise control (ACC) functions [6]. 40% of vehicles on the road are visioned to be automated by the 2040s [7], [8]. Despite the wide presence and burgeoning growth of AV technology, astonishingly, quite rare

efforts have been made to identify AVs in mixed traffic.

Note that existing commercial AVs (e.g., ACC-equipped vehicles) cannot be easily identified by their appearances. Further, commercial AVs may not have a mechanism to notify the AV type to the surrounding vehicles or infrastructure units, given that connected vehicle technology may take time to be widely deployed. Without proper technologies to detect AVs on public roads, it will be hard to evaluate the performance of AVs and their impacts on surrounding traffic at a large and realistic scale. Most existing studies on AV behaviors and their impacts are simulations [9], [10], or small-scale tests involving a few AVs in experimental but not naturalistic settings [11], [12]. There are doubts about whether the findings from these studies perfectly match real-world traffic. Further, the corresponding safety risks and capacity concerns will remain unaddressed without AV identification technologies. Although AV technology is promising in decreasing the number of traffic accidents by reducing human errors, their share of rear-end crashes increases [13], [14]. The reason is that when human-driven vehicles (HVs) are unaware that the preceding vehicles are AVs, they behave the same way as they are following HVs [15]. In this case, the mismatch between AVs' actual driving behavior and HVs' expectations is likely to contribute to traffic accidents [16]. Theoretical studies claim that AVs could reduce headways via platooning and thus improve roadway capacity [17]. However, it has been observed that current commercial AVs drive conservatively and even decrease roadway capacity [18]. From the above, it can be concluded that identifying AVs in mixed traffic is demanded.

Accurately identifying AVs allows for a better understanding and prediction of their behavior on the road. This information can be used to improve safety by enabling other drivers and infrastructure systems to anticipate and react to the actions of AVs more effectively to reduce the potential for accidents. AV identification can complement traditional vehicle classification that classifies HVs into various types (e.g., the FHWA vehicle classes) based on their shapes and appearances [19]. AV identification technology considers a vehicle's dynamic

Manuscript received April 28, 2022. This research is sponsored by National Science Foundation through Grants CMMI #1558887 and CMMI #1932452 and Susan A. Bracken Fellowship, University of South Florida. (Corresponding author: Xiaopeng Li). Dr. Weijun Xie was supported by NSF grant 2246417.

Qianwen Li is with the School of Environmental, Civil, Agricultural and Mechanical Engineering, University of Georgia, USA (e-mail: cam.li@uga.edu).

Xiaopeng Li is with the Department of Civil and Environmental Engineering, University of Wisconsin-Madison, USA (e-mail: xli2485@wisc.edu).

Handong Yao is with the School of Environmental, Civil, Agricultural and Mechanical Engineering, University of Georgia, USA (e-mail: handong.yao@uga.edu).

Zhaohui Liang is with the Department of Civil and Environmental Engineering, University of Wisconsin-Madison, USA (e-mail: zhaohui.liang@wisc.edu).

Weijun Xie is with the H. Milton Stewart School of Industrial and Systems Engineering, Georgia Tech, USA (e-mail: wxie@gatech.edu).

performance. Roadway users could act appropriately once AVs are identified. When following AVs, HVs can drive cautiously to reduce crash risk, while AVs can shorten the car-following distance to save space and increase roadway capacity. By knowing the presence and behavior patterns of AVs within the traffic stream, transportation authorities can make informed decisions. This facilitates the development of infrastructure that supports AV operations, such as dedicated lanes or specialized intersections, enabling more seamless integration of AVs into the existing transportation network. It also allows policymakers to address the unique challenges and opportunities presented by AVs, such as defining specific operational requirements, liability frameworks, and safety standards. Furthermore, naturalistic AV trajectory data is available for assessing the effects of AV technologies on safety, mobility, and energy performance. This data can be used to refine existing AV algorithms, enhance sensor technologies, and further advance the capabilities and safety of AV systems.

Despite these potentials, there has been no published research on AV identification to our knowledge, with the exception of a patent granted to Ford Global Technologies LLC [20]. This patent enabled a vehicle to determine whether its surrounding vehicles were automated. Data of vehicles to be identified were collected, such as speed, acceleration, and steering. However, vehicle car-following dynamics that obviously reflect AV characteristics were not investigated. Besides, this patent was from a single vehicle's perspective rather than a roadside unit. Thus, the identification information only reflects a small view of particular individual vehicles rather than the overall traffic states. And AV data, in this case, was collected by a floating vehicle, which was not as comprehensive as a traffic surveillance system concerning data quality and quantity. More importantly, no technical information about the identification was provided.

This study examines the feasibility of identifying ACC-equipped vehicles (Level 1 and Level 2 AVs) using vehicle trajectories in the mixed traffic stream. It is extendible to Level 3 and + AV identification when data is available. Several learning-based models are constructed to identify AVs using car-following data. Further, a novel cluster-aware learning-based model is developed to identify the disparity among different clusters and train the learning-based model simultaneously, aiming to enhance the identification performance. Two physics-based models are also proposed to compare model performance and draw physical insights. This study is an extension of a conference paper published by the authors, which only tested the identification performance of a few well-established machine learning models with neither comparison with physics-based models nor methodological innovations [21].

The major contribution of this paper is proposing the idea of identifying AVs to enhance mixed traffic management rather than developing a brand-new model or building a comprehensive mixed traffic library. Results show that the cluster-aware-based model already achieved excellent

identification accuracy. A more comprehensive AV identification model can be developed as more datasets involving more AV makes are collected. Once the model is trained, it can be embedded in the existing surveillance system (e.g., video cameras) to identify AVs in real-time.

II. DATA

This section first introduces four datasets used in this study, then presents the data preprocessing.

A. Dataset Introduction

The reasons for choosing the following four datasets are two-fold. First, they are representative enough in a way that all possible traffic situations were tested. Second, they are publicly available. They include data collected in both controlled environments and on public roads with surrounding traffic. Thus, the identification model developed upon application to real-world mixed traffic.

1) HISTORIC data

HISTORIC data includes only HV car-following data [22]. The dataset is publicly available online¹. Data was collected on a controlled highway segment of National Highway G202 in Harbin, Heilongjiang, China, on October 24th, 2015, which was a sunny day. Testing vehicles are 12 identical Kia K5 with 12 GPS-RTK devices installed to collect trajectories without surrounding traffic. The location accuracy of the GPS device was about 1 m, and the speed accuracy was about 0.28 m/s. Further, the data frequency was 20 Hz.

Two traffic scenarios were tested: stationary and oscillated. The leading vehicle was controlled for the stationary scenario to keep a relatively constant speed. For the oscillated scenarios, the leading vehicle was controlled to generate traffic oscillation by decelerating and accelerating periodically. Ideally, the leading vehicle would accelerate to the maximum speed with a predetermined acceleration rate within each oscillation period, cruise at maximum speed, and then decrease to the minimum speed. The deceleration rate should be equal to the negation of the acceleration rate. The leading driver was required to match the intended speed profile as closely as possible. The 11 following vehicles just followed the preceding vehicles as usual without overtaking. The stationary scenario was conducted 7 times with different constant speed settings. The oscillated scenario was conducted 12 times with different oscillation parameters. For more data information, please refer to [22].

2) Vanderbilt ACC data

Vanderbilt ACC data includes only AV car-following data [12]. The data is publicly available online². 8 AVs from different manufacturers were tested. Each participating vehicle was equipped with an uBlox EVK-M8T GPS device to collect trajectories on the controlled public road (about 16 km) without surrounding traffic. The tests were conducted during the daytime, yet the weather information was unavailable. The location accuracy of the GPS device was about 0.24 m, and the speed accuracy was about 0.002 m/s. Further, the data frequency was 10 Hz.

¹ <https://github.com/CATS-Lab-USF/HISTORIC-data>

² <https://vanderbilt.box.com/v/accData>

Two sets of experiments were conducted. In two-vehicle tests, 7 ACC-engaged AVs followed the same cruise control (CC) engaged AV. Different leading vehicle speed profiles reflecting different traffic situations were tested, including oscillatory, low-speed steps, high-speed steps, and speed dips. In the eight-vehicle platoon test, the cognitive and autonomous test vehicle (the CAT vehicle) was used as the leading vehicle. It was controlled to follow a predefined precise speed profile. The other 7 AVs followed preceding vehicles in a single lane to form a platoon. 4 of the 7 following AVs were from the same manufacturer, and the rest 3 were from another manufacturer. For more data information, please refer to [12].

3) CATS Lab ACC data

The connected and autonomous transportation systems laboratory (CATS Lab) collected mixed traffic data, including AVs and HVs [23]. The data is publicly available online³. 5 vehicles were tested, including 2 AVs and 3 HVs. The two AV models were Lincoln MKZ 2016 and 2017. To incorporate all the possible preceding-following vehicle pairs in a mixed traffic context, vehicle arrangement from downstream to upstream was HV, AV, AV, HV, and HV.

Each testing vehicle was installed with a uBox C066-F9P GPS device to collect trajectories. The location accuracy of the GPS devices was about 0.26 m, and the speed accuracy was about 0.089 m/s. The data frequency was 10 Hz.

Data were collected at two locations on open public roads with surrounding traffic. Disturbances caused by surrounding traffic existed during data collection. The team screened the collected raw data to make sure the published datasets are suitable for academic research. Both tests were conducted on clear nights in the presence of street lighting and vehicle headlights. The first set of data with low speed was collected at Lizard Trail Road, Tampa, Florida, USA, on November 18th, 2020. The experiment segment was about 2.4 km. The maximum speed was about 15 m/s. The second set of data with high speed was collected at State Road 56, Tampa, Florida, USA, on November 24th, 2020. The experiment segment was about 8 km. The maximum speed reached roughly 29 m/s. The first vehicle was instructed to generate varied oscillation patterns across several runs by accelerating and decelerating regularly during the experiment. AVs followed preceding vehicles with ACC turned on. HVs followed preceding vehicles as usual. Overtaking was prohibited.

4) Open ACC data

Open ACC data was collected by the European Commission in a previous study [24]. It is publicly available online⁴. It contains pure AV, pure HV, and mixed traffic data. 4 experiment campaigns were conducted. The data from campaign 3 was used. Data was collected in the second quarter of 2019 for two days on the rural road of the AstaZero test track in Sweden. Disturbances caused by surrounding traffic existed during data collection. The team screened the collected raw data to make sure the published datasets are suitable for academic research.

The weather information was not available. The study

segment was about 5.7 km. Five vehicles participated. The leading vehicle was Audi A8. The following vehicles were Tesla Model 3, BMW X5, Mercedes A-Class, and Audi A6. Data was collected using an inertial navigation system with differential GNSS accuracy. The speed accuracy was about 0.02 m/s, and the location accuracy was about 0.02 m. The data collection frequency was 10Hz.

TABLE 1
DATA DESCRIPTIVE STATISTICS

Data	Maximum	Minimum	Mean	Standard deviation
HISTORIC data				
d (m)	121.23	3.23	24.73	15.14
v^p (m/s)	24.89	1.06	9.51	4.38
v^f (m/s)	24.37	1.03	9.51	4.39
a^f (m/s ²)	4.83	-4.89	0.00	0.16
ϕ (1: AV; 0: HV)	0	0	0	0
Vanderbilt ACC data				
d (m)	100.76	17.06	45.45	15.41
v^p (m/s)	33.54	12.92	23.46	4.23
v^f (m/s)	34.96	12.94	23.45	4.32
a^f (m/s ²)	4.08	-4.64	0.00	0.36
ϕ (1: AV; 0: HV)	1	1	1	0
CATS Lab ACC data				
d (m)	92.42	3.59	37.53	13.10
v^p (m/s)	28.58	5.52	20.93	4.65
v^f (m/s)	29.58	5.44	20.95	4.73
a^f (m/s ²)	3.2	-5.29	0.00	0.56
ϕ (1: AV; 0: HV)	1	0	0.52	0.50
Open ACC data				
d (m)	133.81	5.10	30.43	15.36
v^p (m/s)	33.79	5.12	18.72	4.07
v^f (m/s)	34.54	5.01	18.73	4.21
a^f (m/s ²)	3.23	-3.99	0.00	0.48
ϕ (1: AV; 0: HV)	1	0	0.78	0.41

Two sets of experiments were conducted. The leading vehicle always operated with ACC activated. The 4 following vehicles were driven by humans in the first set of experiments but operated with ACC activated in another set of experiments. Two car-following patterns were applied during the experiments: platoon with constant speed and platoon with perturbation of the target speed. For more information, please refer to the European Commission webpage.

B. Data Preprocessing

The following procedures are used for data preprocessing.

1. The longitude and latitude of the vehicle are utilized to determine its location
2. To fill in the missing data caused by the GPS device, linear interpolation is deployed on the vehicle location.
3. Car-following distance is calculated by subtracting vehicle length from vehicle location difference.
4. Given the less than satisfactory GPS accuracy, i.e., 0.28 m/s, moving average smoothing is applied to the collected speed in the HISTORIC dataset. The data is resampled to 10 Hz.

³ <https://github.com/CATS-Lab-USF/CATS-Lab-ACC-data>

⁴ <https://data.jrc.ec.europa.eu/dataset/9702c950-c80f-4d2f-982f-44d06ea0009f>

5. Vehicle acceleration is calculated from vehicle speed by taking the first-order derivative.
6. Only stable car-following periods are considered. Data at the beginning or end of the test runs are excluded.

After preprocessing, there are 886,853 data points in HISTORIC data, 348,425 in Vanderbilt ACC data, 76,069 in CATS Lab ACC data, and 525,655 in Open ACC data. Vehicle type is coded as a binary variable with 1 denoting AV and 0 denoting HV. Data descriptive statistics are provided in TABLE 1. d denotes the car-following distance, v^p denotes the preceding vehicle speed, v^f denotes the following vehicle speed, a^f denotes the following vehicle acceleration, and φ denotes vehicle type.

Vanderbilt ACC data is the one with the greatest average speed and average car-following distance, followed by CATS Lab ACC data, Open ACC data, and HISTORIC data. The speed and car-following distance oscillation magnitudes across different datasets are similar, indicated by the relatively similar car-following distance and speed standard deviations. As for the following vehicle acceleration, CATS Lab ACC data has the greatest oscillation, followed by Open ACC data, Vanderbilt data, and HISTORIC data.

III. METHODOLOGY

This section introduces different models used for AV identification, including 2 physics-based models, 7 learning-based models, and one cluster-aware learning-based model.

Each dataset is segmented by the identification time window Δt . Segments from four datasets are mixed to construct a comprehensive identification model. Each trajectory segment denotes an observation. All observations are shuffled and then divided into two subsets. The dividing ratio is 9:1. 10-fold cross-validation is adopted. 90% of observations are used for training and validation. The rest observations are used for testing. The average value across ten folds is used to indicate the final model performance.

A. Physics-based Models

This subsection presents physics-based models to identify AVs. Two popular car-following models are adopted, including the intelligent driver model (IDM), and the optimal velocity model (OVM) [22], [25].

The following vehicle's acceleration calculated based on the IDM is formulated as [25]:

$$a_i(t) = a^0 \left[1 - \left(\frac{v_i(t)}{v^0} \right)^\delta - \left(\frac{s_i^*(t)}{s_i(t)} \right)^2 \right], \quad (1)$$

where $s_i^*(t) = s^0 + \max \left(0, v_i(t)T^0 + \frac{v_i(t) \times [v_i(t) - v_{i-1}(t)]}{2\sqrt{a^0 b^0}} \right)$ is the desired gap, $s_i(t) = x_{i-1}(t) - x_i(t) - l_{i-1}$ is the car-following distance between vehicles, a^0 is the maximum acceleration, v_i is the following vehicle speed, v_{i-1} is the preceding vehicle speed, v^0 is the desired speed, s^0 is the minimum distance, T^0 is the time gap, b^0 is the comfortable deceleration, x_i is the following vehicle location, x_{i-1} is the preceding vehicle location, and l_{i-1} is the preceding vehicle length.

The acceleration calculated based on the OVM is given as follows.

$$a_i(t) = \frac{\max \left(0, \min \left(v^0, \frac{s_i(t) - s^0}{T^0} \right) \right) - v_i(t)}{\tau}, \quad (2)$$

where τ is the adaptation time and other notation has been introduced above.

The first half of the data in the time window Δt is used for model calibration. The calibration objective is to minimize the following formula.

$$\sum_{i=1}^N \left(\frac{\dot{v}_i - \hat{v}_i}{\dot{v}_i} \right)^2, \quad (3)$$

where $N = \left\lceil \frac{\Delta t/2}{0.1} \right\rceil$ is the number of data points, \dot{v}_i is the observed acceleration, and \hat{v}_i is the calibrated acceleration.

The interior-point method is used to find the optimal model parameters. The other half of the data is used for model validation. The root mean square error (RMSE) between the observed and validated car-following distance is calculated as:

$$\text{RMSE}_s = \sqrt{\frac{\sum_{i=1}^M (s_i - \hat{s}_i)^2}{M}}, \quad (4)$$

where $M = \frac{\Delta t}{0.1} - N$ is the number of data points, s_i is the observed car-following distance, and \hat{s}_i is the validated car-following distance.

The above calibration and validation are conducted for all observations. An RMSE threshold RMSE_s^* producing the highest identification accuracy without underfitting and overfitting is found for each fold during the 10-fold cross-validation, see Fig. 1. Given the stochastic nature of human driving, if the RMSE_s of observation is less than RMSE_s^* , it is identified as an AV; Otherwise, it is identified as an HV. Then, the optimal threshold RMSE_s^* is used to identify AVs from the testing set.

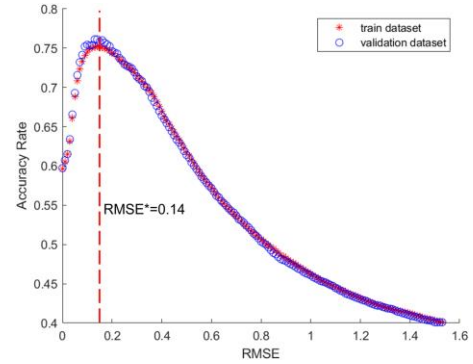


Fig. 1 Threshold selection for physics-based models.

B. Learning-based models

This subsection presents learning-based models. Various machine learning classification models have been proposed in the past, and each of them merits in different aspects [26]–[28]. To accomplish the most accurate AV identification, seven popular models are tested. Model inputs include the preceding vehicle speed, the following vehicle speed, the following vehicle acceleration, and the car-following distance. These inputs are standardized between -1 to 1. The vehicle type, either an AV or an HV, is the model output.

1) Long short-term memory network (LSTM)

As a type of recurrent neural network, the LSTM can learn

order dependence in time series data [28]. The LSTM model structure is illustrated in Fig. 2 (a). The dropout rates of each dropout layer are denoted by $\alpha_k \in \{\alpha_1, \alpha_2, \dots, \alpha_K\}$. The numbers of hidden layers N , dropout layers D , neurons N_{neuron} , batches N_{batch} , and epochs N_{epoch} and the dropout rate α_k are adjusted during model training.

The neuron structure is illustrated in Fig. 2 (b). X_t is the current input vector. h_t is the current neuron output. C_{t-1} is the memory from the last neuron. h_{t-1} is the output of the last neuron. C_t is the memory from the current neuron. The subscript t indexes the time step (0.1s is used here). \times is element-wise multiplication. $+$ is the element-wise summation/concatenation. σ is the sigmoid layer. \tanh is the hyperbolic tangent layer. $\{b_1, b_2, b_3, b_4\}$ are biases.

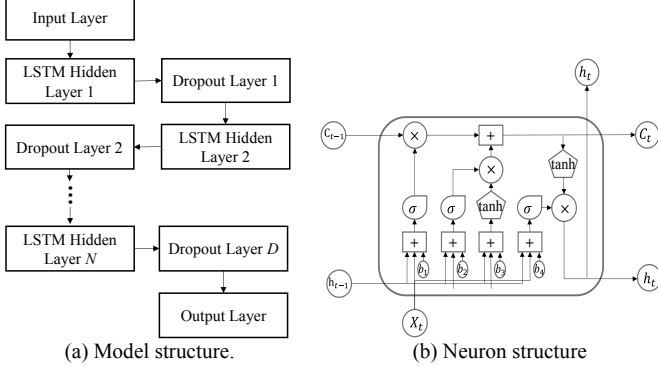


Fig. 2. LSTM structure.

2) Support vector machine (SVM)

The SVM model finds the best decision boundary, i.e., decision hyperplane, to separate different classes. The distance from the best hyperplane to the nearest data point of each class is the greatest [26]. Different kernels are tested to produce the best results.

3) k -nearest neighbors (KNN)

For the KNN model, we first calculate the Euclidean distance from the query observation to the classified observations. The classified observations are ordered by increasing distance. The class of the query observation is the majority voting of the top k observation classes [26]. The number of nearest neighbors k is adjusted for the best model results.

4) Fixed-radius near neighbors (FRNN)

The FRNN is a variant of the KNN. Instead of using only the k -neighbors, the FRNN locates all observations in the training set that are within a given radius r of the query observation. The selected radius neighbors are used to predict the query observation [29]. The radius value is adjusted to yield the best model results.

5) Random forest (RF)

The RF model has T decision trees. The maximum depth of each decision tree is N_T . Each decision tree produces each result. The final result is derived as the majority voting of all trees' results [26]. The number of decision trees T and the maximum depth T_N are tuned for the best model performance.

C. Cluster-aware learning-based model (CALM)

The above learning-based models assume that observations are homogeneous. However, it is natural that different observations belong to different clusters [30]. Ignoring such

heterogeneity within datasets may miss the opportunity to further improve the model performance. Therefore, we adopt the cluster-aware learning-based (CALM) technique to identify the disparity among different clusters and train the learning-based model simultaneously. The general cluster-aware learning framework was constructed in [31]. We adapt and materialize the framework with a classification application to identify AVs. However, the CALM is challenged with a considerably long training time. Motivated by this superiority, we further combine the CALM with transfer learning techniques, which have been proven to hold great potential in enhancing model training efficiency [32], [33]. Specifically, the well-tuned learning model weights are fed into the CALM as initial weights.

Suppose that all samples belong to K clusters. Within each cluster $k \in [K] := \{1, \dots, K\}$, the responses can be described by the data points with heterogeneity. The CALM is formulated as follows;

$$\min_{\theta, \delta, m} R_p(\delta, \theta, m) := \sum_{i \in [N]} \sum_{k \in [K]} \delta_{ik} \ell_k(y_i, \mathbf{x}_i, \theta_k) + \rho \sum_{i \in [N]} \sum_{k \in [K]} \delta_{ik} S_k(\mathbf{x}_i, \mathbf{m}_k) \quad (5)$$

s.t.,

$$\begin{aligned} \sum_{k \in [K]} \delta_{ik} &= 1, \forall i \in [N], \\ \sum_{i \in [N]} \delta_{ik} &\geq 1, \forall k \in [K], \\ \delta_{ik} &\in \{0, 1\}, \forall k \in [K], \forall i \in [N], \end{aligned}$$

where \mathbf{x}_i and y_i are the inputs and outputs, respectively. ρ is a non-negative tuning parameter controlling overlapping areas among different clusters. δ_{ik} denotes the cluster assignment decision, i.e., $\delta_{ik} = 1$ implies that the data point \mathbf{x}_i belongs to cluster k . $\ell_k(y_i, \mathbf{x}_i, \theta_k) = -\frac{1}{I} \sum_i (y_i \ln LSTM(\mathbf{x}_i, \theta_k) + (1 - y_i) \ln(1 - LSTM(\mathbf{x}_i, \theta_k)))$ denotes the loss function with an estimator θ_k . $LSTM(\mathbf{x}_i, \theta_k)$ denotes the function of the LSTM network. $S_k(\mathbf{x}_i, \mathbf{m}_k)$ denotes the dissimilarity function (i.e., the squared Euclidean distance-based function $S_k(\mathbf{x}_i, \mathbf{m}_k) = \|\mathbf{x}_i - \mathbf{m}_k\|^2$), where \mathbf{m}_k is the centroid of cluster k .

Then, regularized alternating minimization (RAM) algorithm is used to solve the CALM. The RAM algorithm first separates the decision variables of the CALM into two parts, i.e., (θ, m) and δ . Next, it solves the CALM with respect to (θ, m) by fixing δ and vice versa. At each iteration, an $L_{1,1}$ the norm penalty term is added to the objective function of the CALM to ensure that the current clustering solution and the previous solution are not too different. Finally, it terminates whenever the clustering solutions from two consecutive iterations are close. The detailed RAM algorithm can be found in [31]. The number of clusters K and the parameter controlling the overlapping areas among cluster ρ are tuned for the best CALM identification accuracy.

IV. RESULTS

This section reports the tuning results of learning-based models and compares the identification performance of different models.

A. Learning-based Model Tuning

In terms of the overall identification accuracy α , model tuning is done for multiple input data time frames Δt , ranging from 0.2s to 5s.

$$\alpha = \frac{N_{AV}^{\text{correct}} + N_{HV}^{\text{correct}}}{N_{\text{all}}} \times 100\%, \quad (6)$$

where N_{AV}^{correct} is the number of successfully identified AVs, N_{HV}^{correct} is the number of successfully identified HVs, and N_{all} is the total number of observations.

1) LSTM

Different configurations are tested regarding the numbers of hidden layers N , dropout layers D , neurons N_{neuron} , batches N_{batch} , and epochs N_{epoch} and the dropout rate α_k . The best model is found for different time windows Δt . The tuning results when $\Delta t = 1$ s are reported in TABLE 2. The best result

TABLE 2
LSTM MODEL TUNING RESULTS.

Δt (s)	N	D	α_k	N_{neuron}	N_{batch}	N_{epoch}	Training α	Validation α	Testing α
1	2	3	0.5	500	2000	200	94.87%	94.66%	94.54%
1	3	3	0.5	500	2000	200	96.34%	96.18%	96.07%
1	4	3	0.5	500	2000	200	98.47%	95.47%	96.88%
1	3	5	0.5	500	2000	200	93.77%	96.54%	95.98%
1	3	3	0.7	500	2000	200	94.76%	96.72%	96.53%
1	3	3	0.2	500	2000	200	97.64%	95.16%	96.39%
1	3	3	0.5	200	2000	200	93.05%	92.87%	92.85%
1	3	3	0.5	500	1000	200	96.47%	96.24%	96.18%
1	3	3	0.5	500	2000	300	96.39%	96.37%	95.97%

3) CALM

Based on the above tuning results, LSTM produces the highest accuracy among all learning-based models. Motivated by this superiority, the CALM is developed based on the LSTM, i.e., the loss function $\ell_k(y_i, \mathbf{x}_i, \boldsymbol{\theta}_k)$ of the CALM is set as the LSTM. K and ρ are tuned for the best CALM identification accuracy, while the LSTM model parameters remain the same as the ones that produce the best results in TABLE 2, i.e., $N = D = 3$, $N_{\text{neuron}} = 500$, $N_{\text{batch}} = 1000$, $N_{\text{epoch}} = 200$, and $\alpha_k = 0.5$. The best CALM model parameters are found for different time windows Δt . The tuning results when $\Delta t = 1$ s are presented in TABLE 3. The best CALM result is yielded when $K = 2$ and $\rho = 0.1$. The accuracy is higher than the LSTM alone, indicating the effectiveness of clustering observations.

TABLE 3
CALM MODEL TUNING RESULTS.

Δt (s)	K	ρ	Training α	Validation α	Testing α
1	2	0.01	97.79%	96.37%	96.55%
1	2	0.1	97.96%	96.61%	96.82%
1	2	1	98.17%	96.77%	96.46%
1	3	0.01	96.78%	96.41%	95.94%
1	3	0.1	98.38%	96.77%	96.36%
1	3	1	97.31%	96.28%	96.42%
1	4	0.01	96.46%	95.90%	95.67%
1	4	0.1	98.36%	96.49%	96.42%
1	4	1	97.86%	95.95%	96.13%

B. Model comparison

After the above model tuning, the identification performance of different models is compared in this subsection. Different identification time window lengths are tested, ranging from 0.2s to 5s. AV and HV sample ratios vary from 0.66 to 0.69. Physics-

is observed when $N = D = 3$, $N_{\text{neuron}} = 500$, $N_{\text{batch}} = 1000$, $N_{\text{epoch}} = 200$, and $\alpha_k = 0.5$. It is noted that when N is greater or α_k is smaller, overfitting is produced. When D or α_k is greater, underfitting is yielded.

2) Other traditional learning-based models

Various parameters are tested for other traditional learning-based models. The best model is found for different time windows Δt . The best SVM result is observed (accuracy is 86.65%) when rbf kernel is used. The best KNN result is observed (accuracy is 91.76%) when $k = 15$. Overfitting exists when $k = 5$ and $k = 10$. The best FRNN result is observed (accuracy is 92.30%) when $r = 3$. The best RF result is observed (accuracy is 91.65%) when $T = 15$, $N_T = 8$. When $N_T = 15$, overfitting is produced.

based models are not investigated when $\Delta t < 1$ s due to insufficient calibration and validation data. To study the model performance of identifying AVs and HVs respectively, the identification accuracy of AVs α_{AV} (i.e., the true positive rate, recall, or sensitivity $\times 100\%$) and HVs α_{HV} (i.e., the true negative rate or specificity $\times 100\%$) are calculated as follows.

$$\alpha_{AV} = \frac{N_{AV}^{\text{correct}}}{N_{AV}} \times 100\%, \quad (7)$$

$$\alpha_{HV} = \frac{N_{HV}^{\text{correct}}}{N_{HV}} \times 100\%, \quad (8)$$

where N_{AV}^{correct} is the number of successfully identified AVs, N_{HV}^{correct} is the number of successfully identified HVs, N_{AV} is the total number of AVs, and N_{HV} is the total number of HVs.

Further, the identification precision and F1-score are computed.

$$\text{Precision} = \frac{N_{AV}^{\text{correct}}}{N_{AV}^{\text{correct}} + (N_{HV} - N_{HV}^{\text{correct}})}$$

F1 score

$$= \frac{2N_{AV}^{\text{correct}}}{2N_{AV}^{\text{correct}} + (N_{HV} - N_{HV}^{\text{correct}}) + (N_{AV} - N_{AV}^{\text{correct}})}$$

Results are shown in TABLE 4 and TABLE 5. We see that the identification accuracy of physics-based models (i.e., IDM and OVM) increases and decreases with the time window length. The same trend is observed for the precision and the F1 score. Possible reasons follow. When Δt is shorter, fewer data points result in less accurate model calibration and validation. The results are less representative in terms of car-following behavior. Thus, the identification accuracy is lower. When Δt is longer, the car-following behavior during the time window is more heterogeneous, and the explanatory power of physics-based models with only one set of parameters degrades. In this

case, RMSE_s values are greater, and the identification accuracy is lower. On average, physics-based models correctly identify over 80% AVs and over 70% HVs, indicating that AV car-following behavior likely has a smaller residual error than HVs. It also noted that the IDM's identification accuracy, precision,

and F1 score are slightly higher than the OVM across all instances. This is expected because more parameters are calibrated in the IDM, and the resulting model better represents the realistic driving behavior.

TABLE 4
MODEL COMPARISON PART 1.

Δt (s)	α_{AV}/α_{HV}	Model							
		IDM	OVM	LSTM	SVM	KNN	FRNN	RF	CALM
0.2	α_{AV}	/	/	98.14%	93.89%	93.97%	93.71%	96.12%	98.86%
	α_{HV}	/	/	94.32%	80.13%	86.54%	85.39%	83.73%	94.47%
0.5	α_{AV}	/	/	97.97%	93.45%	94.20%	94.50%	95.90%	99.59%
	α_{HV}	/	/	94.31%	79.99%	87.71%	86.48%	83.48%	95.98%
1	α_{AV}	82.03%	78.23%	98.18%	93.51%	93.79%	94.32%	96.03%	98.65%
	α_{HV}	68.91%	70.93%	93.77%	79.84%	86.99%	85.35%	84.09%	95.43%
2	α_{AV}	84.14%	85.05%	98.04%	93.74%	93.89%	93.49%	96.28%	98.37%
	α_{HV}	71.55%	72.85%	94.08%	80.51%	85.24%	87.67%	84.23%	95.76%
3	α_{AV}	86.37%	86.27%	98.12%	93.85%	93.67%	94.14%	96.12%	99.18%
	α_{HV}	74.90%	73.87%	94.25%	79.81%	86.27%	85.98%	83.18%	95.65%
4	α_{AV}	83.85%	78.90%	98.27%	93.37%	94.20%	94.07%	96.30%	98.47%
	α_{HV}	74.94%	72.71%	94.24%	80.57%	84.81%	87.55%	83.24%	95.35%
5	α_{AV}	83.19%	78.70%	98.24%	93.25%	92.96%	93.77%	96.21%	98.42%
	α_{HV}	74.13%	72.96%	94.35%	80.66%	85.06%	86.28%	84.11%	95.50%
Average α_{AV}		83.92%	81.43%	98.14%	93.58%	93.81%	94.00%	96.14%	98.79%
Average α_{HV}		72.89%	72.66%	94.19%	80.22%	86.09%	86.39%	83.72%	95.45%

TABLE 5
MODEL COMPARISON PART 2.

Δt (s)	Precision/ F1 score	Model							
		IDM	OVM	LSTM	SVM	KNN	FRNN	RF	CALM
0.2	Precision	/	/	0.929	0.783	0.842	0.830	0.818	0.932
	F1 score	/	/	0.955	0.854	0.888	0.880	0.884	0.959
0.5	Precision	/	/	0.929	0.779	0.853	0.841	0.814	0.949
	F1 score	/	/	0.953	0.850	0.895	0.890	0.881	0.972
1	Precision	0.639	0.644	0.924	0.780	0.847	0.831	0.822	0.943
	F1 score	0.719	0.706	0.952	0.851	0.890	0.884	0.886	0.964
2	Precision	0.665	0.676	0.925	0.781	0.825	0.849	0.819	0.945
	F1 score	0.743	0.754	0.952	0.852	0.878	0.890	0.885	0.964
3	Precision	0.697	0.692	0.928	0.778	0.837	0.835	0.811	0.945
	F1 score	0.772	0.768	0.954	0.850	0.884	0.885	0.880	0.968
4	Precision	0.692	0.653	0.929	0.786	0.826	0.853	0.815	0.942
	F1 score	0.758	0.715	0.955	0.854	0.880	0.895	0.883	0.963
5	Precision	0.683	0.664	0.930	0.787	0.827	0.840	0.823	0.944
	F1 score	0.750	0.720	0.956	0.854	0.875	0.886	0.887	0.964
Average	Precision	0.675	0.666	0.928	0.782	0.837	0.840	0.818	0.943
	F1 score	0.748	0.733	0.954	0.852	0.884	0.887	0.884	0.965

All learning-based models' identification accuracy, precision, and F1 score remain relatively steady across input data time window lengths. This demonstrates the model's robustness. Different results are observed across learning-based models, given their different learning capabilities. In TABLE 4, as expected, the LSTM produces the highest AV identification accuracy with an average of 98.14%, followed by the RF with an average of 96.14%, FRNN with an average of 94.00%, KNN with an average of 93.81%, and SVM with an average of 93.58%. The LSTM also produces the highest HV identification accuracy with an average of 94.19%, followed by the FRNN with an average of 86.39%, KNN with an average of 86.09%, RF with an average of 83.72%, and SVM with an average of 80.22%. In TABLE 5, the highest precision and F1 score are also produced by the LSTM. The best performance achieved by the LSTM demonstrates its superiority in learning order dependence in time series data. Enhancing LSTM with the

cluster-aware technique further improves the identification performance. In TABLE 4, the AV identification accuracy is increased from 98.14% to 98.79%, and the HV identification accuracy is increased from 94.19% to 95.45%. In TABLE 5, the precision increases from 0.928 to 0.943, and the F1 score increases from 0.954 to 0.965. Although the improvement is insignificant (the LSTM alone is already rather good in differentiating AVs and HVs), it is still valuable to real-world applications. These findings demonstrate the feasibility of identifying AVs from HVs using car-following trajectory data.

The identification accuracy of AVs is always greater than that of HVs across all models. This suggests that some HVs drive in a way quite similar to AVs, i.e., less volatile, and thus they are likely to be misidentified as AVs. Further investigations are needed to address this issue. Gaussian processes may be combined with physics-based models to reveal more insights.

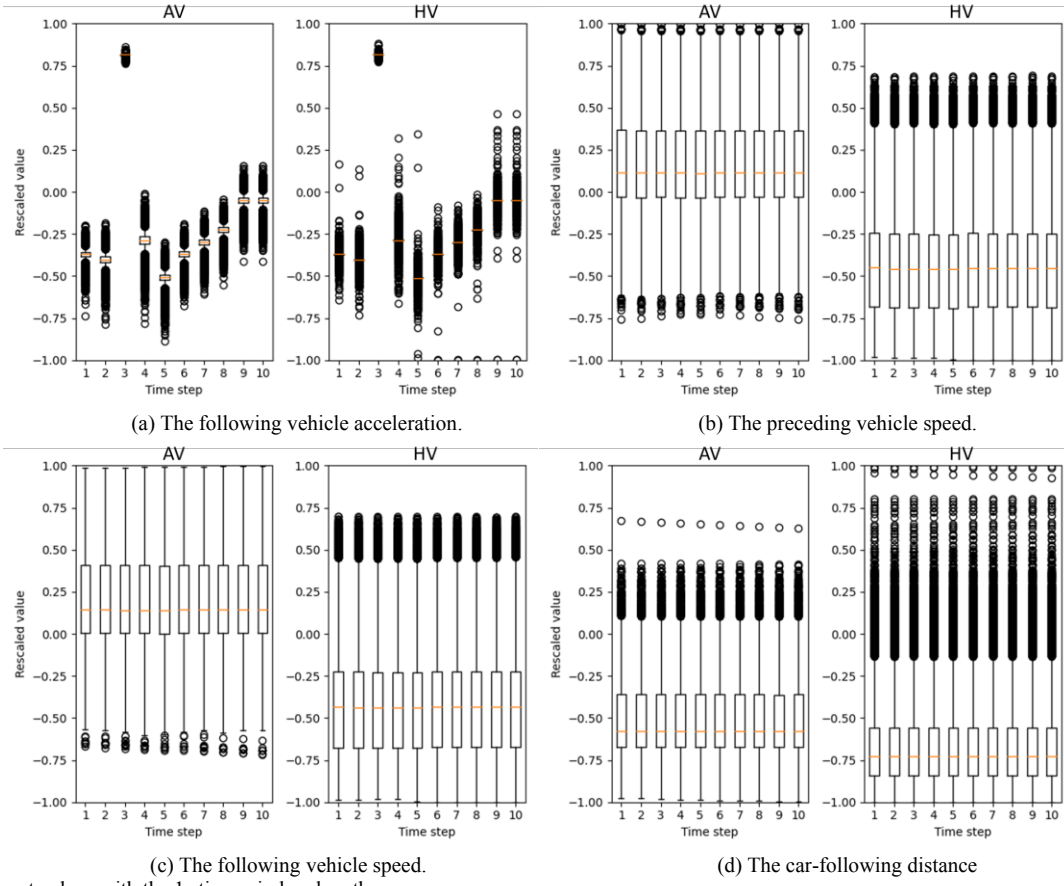


Fig. 3 Rescaled input values with the 1s time window length.

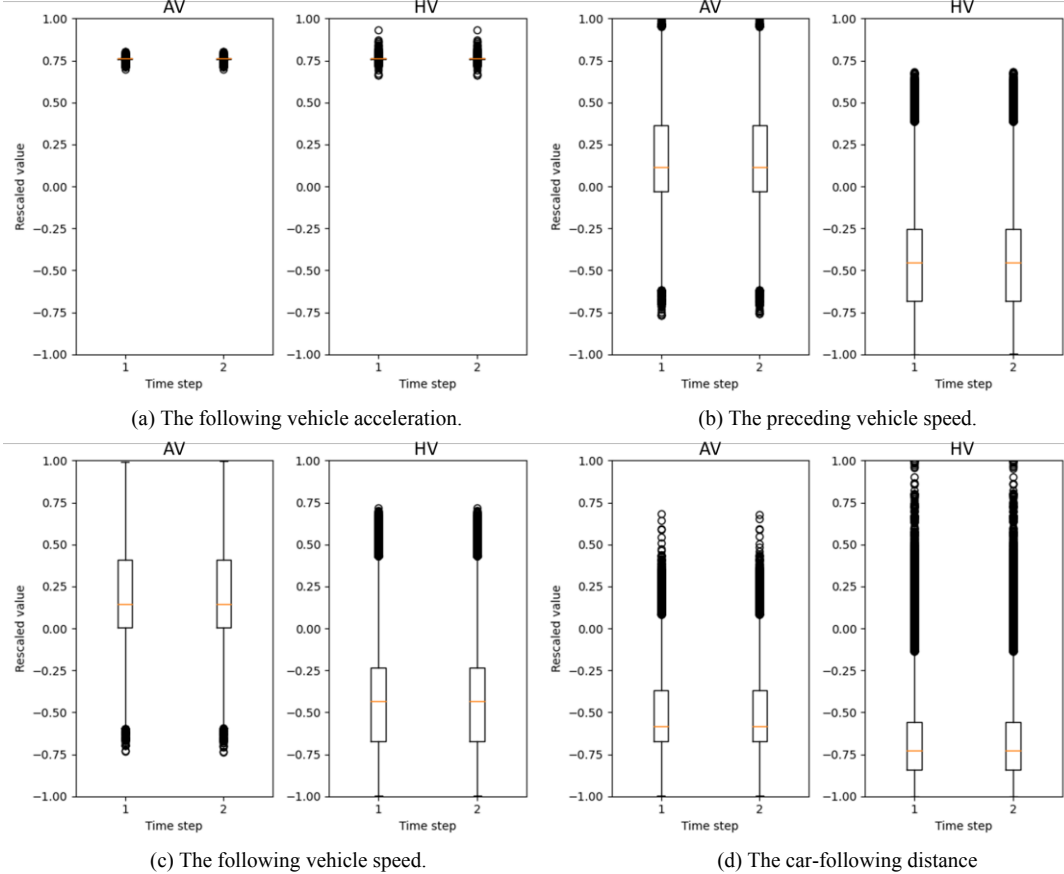


Fig. 4 Rescaled input values with the 0.2s time window length.

C. Mechanism of identification

From the above analysis, we know it is possible to identify AVs and HVs with physics and learning-based models. To better understand the mechanism of learning-based models, we take the LSTM model as an example and plot Fig. 3 and Fig. 4.

Fig. 3 shows the rescaled input values, including the following vehicle acceleration, preceding vehicle speed, following vehicle speed, and car-following distance after standardization with the 1s time window length. It is evident that HVs yield a greater range of outliers than AVs in all input values. This is because that AVs' driving behavior is more stable than HVs'. The median rescaled value of AV speed is positive, but the median rescaled value of HV speed is negative. This means AVs run at a greater speed than HVs. Also, we find that the AV's median rescaled car-following distance is greater than HV's. This is consistent with [23] that current commercial low-level AVs prefer to set a greater headway for safety. Further, the same findings could be observed in Fig. 4 with the 0.2s time window length.

V. CONCLUSION

Identifying AVs from mixed traffic streams is expected to improve traffic safety, increase roadway capacity, and promote AV management and development. Failure to recognize the presence of AVs could result in adverse consequences, e.g., crashes.

This research examines the feasibility of detecting ACC-equipped vehicles (SAE Levels 1 and 2) utilizing vehicle trajectories collected from existing infrastructure. The results reveal that learning-based models can correctly distinguish AVs from HVs with high accuracy. Integrating the cluster-aware technique into the learning-based model further improves identification accuracy. These exciting results open unparalleled data-driven potential for studying and managing mixed traffic.

Regarding future research, the first direction would be addressing the prior probability shift issue [34]. For this study, the issue happens when the AV market penetration rate changes. In this case, the identification accuracy may degrade. Second, rather than focusing solely on longitudinal vehicle behavior, it would be interesting to investigate lateral vehicle movements in the AV identification process. Third, the current dataset size is limited. It would be necessary to keep enriching the current mixed traffic data pool (e.g., more AV makes and high-level AVs) and updating the identification model. Instead of academic research, this would be an engineering practice that should be carried out when it comes to real-world implementation. Fourth, building an online learning model for real-time prediction would be more valuable. The proposed model can be implemented as a pre-trained model allowing a proper initial prediction performance. Furthermore, because AVs from different manufacturers or onboard equipment providers have distinct AV control logic and automation levels, it is critical to distinguish them. This allows for manufacturer-specific AV evaluation and the extraction of more vendor-specific insights for rapid AV technology development.

ACKNOWLEDGEMENT

The authors thank the reviewers and the editor for their constructive comments in improving the paper quality.

REFERENCES

- [1] Y. Han and S. Ahn, "Stochastic modeling of breakdown at freeway merge bottleneck and traffic control method using connected automated vehicle," *Transportation research part B: methodological*, vol. 107, pp. 146–166, 2018.
- [2] M. M. Morando, Q. Tian, L. T. Truong, and H. L. Vu, "Studying the safety impact of autonomous vehicles using simulation-based surrogate safety measures," *Journal of Advanced Transportation*, vol. 2018, 2018.
- [3] S. R. Rad, H. Farah, H. Taale, B. van Arem, and S. P. Hoogendoorn, "Design and operation of dedicated lanes for connected and automated vehicles on motorways: A conceptual framework and research agenda," *Transportation research part C: emerging technologies*, vol. 117, p. 102664, 2020.
- [4] J. Olstam, F. Johansson, A. Alessandrini, P. Sukennik, J. Lohmiller, and M. Friedrich, "An approach for handling uncertainties related to behaviour and vehicle mixes in traffic simulation experiments with automated vehicles," *Journal of advanced transportation*, vol. 2020, 2020.
- [5] J. Sun, Z. Zheng, and J. Sun, "Stability analysis methods and their applicability to car-following models in conventional and connected environments," *Transportation Research Part B: Methodological*, vol. 109, pp. 212–237, 2018.
- [6] Canlys, "10 % of new cars in the US sold with level 2 autonomy features," no. September, 2019.
- [7] M. Lavasani, X. Jin, and Y. Du, "Market penetration model for autonomous vehicles on the basis of earlier technology adoption experience," *Transportation Research Record*, vol. 2597, no. 1, pp. 67–74, 2016.
- [8] T. Litman, "Autonomous vehicle implementation predictions: Implications for transport planning," 2020.
- [9] H. Yao and X. Li, "Decentralized control of connected automated vehicle trajectories in mixed traffic at an isolated signalized intersection," *Transportation Research Part C: Emerging Technologies*, vol. 121, p. 102846, 2020.
- [10] H. Jiang, J. Hu, S. An, M. Wang, and B. B. Park, "Eco approaching at an isolated signalized intersection under partially connected and automated vehicles environment," *Transportation Research Part C: Emerging Technologies*, vol. 79, pp. 290–307, 2017.
- [11] Z. Wang, X. Zhao, Z. Xu, X. Li, and X. Qu, "Modeling and field experiments on autonomous vehicle lane changing with surrounding human-driven vehicles," *Computer-Aided Civil and Infrastructure Engineering*, 2020.
- [12] G. Gunter *et al.*, "Are commercially implemented adaptive cruise control systems string stable?," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [13] A. Deluka Tibljaš, T. Giuffrè, S. Surdonja, and S. Trubia, "Introduction of Autonomous Vehicles: Roundabouts design and safety performance evaluation," *Sustainability*, vol. 10, no. 4, p. 1060, 2018.
- [14] F. M. Favarò, N. Nader, S. O. Eurich, M. Tripp, and N. Varadaraju, "Examining accident reports involving autonomous vehicles in California," *PLoS one*, vol. 12, no. 9, p. e0184952, 2017.
- [15] X. Zhao, Z. Wang, Z. Xu, Y. Wang, X. Li, and X. Qu, "Field experiments on longitudinal characteristics of human driver behavior following an autonomous vehicle," *Transportation Research Part C: Emerging Technologies*, vol. 114, pp. 205–224, 2020.
- [16] Đ. Petrović, R. Mijailović, and D. Pešić, "Traffic Accidents with Autonomous Vehicles: Type of Collisions, Manoeuvres and Errors of Conventional Vehicles' Drivers," *Transportation research procedia*, vol. 45, pp. 161–168, 2020.
- [17] S. Zong, "How Connected Autonomous Vehicles Would Affect Our World?—A Literature Review on the Impacts of CAV on Road Capacity, Environment and Public Attitude," in *MATEC Web of Conferences*, 2019, vol. 296, p. 1007.
- [18] X. Shi and X. Li, "Empirical Study on Car Following Characteristics of Commercial Automated Vehicles with Different

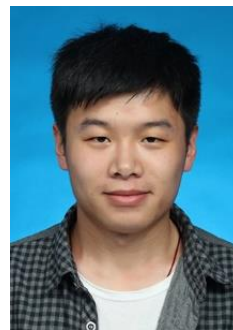
- Headway Settings,” *Preprint, Researchgate*, 2020.
- [19] M. E. Hallenbeck, O. I. Selezneva, and R. Quinley, “Verification, refinement, and applicability of long-term pavement performance vehicle classification rules,” United States. Federal Highway Administration. Office of Infrastructure ..., 2014.
- [20] T. E. Pilutti, M. Y. Rupp, R. A. Trombley, A. Waldis, and W. T. Yopp, “Autonomous vehicle identification,” Google Patents, 24-Jan-2017.
- [21] Q. Li, X. Li, H. Yao, and Z. Liang, “Automated Vehicle Identification Based on Car-following Dynamics,” *IEEE International Intelligent Transportation Systems Conference*, 2021.
- [22] H. Yao, Q. Li, and X. Li, “A study of relationships in traffic oscillation features based on field experiments,” *Transportation Research Part A: Policy and Practice*, vol. 141, pp. 339–355, 2020.
- [23] X. Shi and X. Li, “Empirical study on car-following characteristics of commercial automated vehicles with different headway settings,” *Transportation Research Part C: Emerging Technologies*, vol. 128, p. 103134, 2021.
- [24] B. Ciuffo, K. Mattas, A. Anesiadou, and M. Makridis, “Open ACC Database. European Commission, Joint Research Centre (JRC) [Dataset] PID,” 2020.
- [25] M. Treiber, A. Hennecke, and D. Helbing, “Congested traffic states in empirical observations and microscopic simulations,” *Physical review E*, vol. 62, no. 2, p. 1805, 2000.
- [26] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*, vol. 112. Springer, 2013.
- [27] M. H. Hassoun, *Fundamentals of artificial neural networks*. MIT press, 1995.
- [28] X. Zhang, J. Sun, X. Qi, and J. Sun, “Simultaneous modeling of car-following and lane-changing behaviors using deep learning,” *Transportation research part C: emerging technologies*, vol. 104, pp. 287–304, 2019.
- [29] J. L. Bentley, “Survey of techniques for fixed radius near neighbor searching,” Stanford Linear Accelerator Center, Calif.(USA), 1975.
- [30] P. Arumugam and V. Christy, “Analysis of clustering and classification methods for actionable knowledge,” *Materials Today: Proceedings*, vol. 5, no. 1, pp. 1839–1845, 2018.
- [31] S. Chen and W. Xie, “On Cluster-Aware Supervised Learning: Frameworks, Convergent Algorithms, and Applications,” *INFORMS Journal on Computing*, 2021.
- [32] F. Zhuang *et al.*, “A comprehensive survey on transfer learning,” *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2020.
- [33] L. Torrey and J. Shavlik, “Transfer learning,” in *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*, IGI global, 2010, pp. 242–264.
- [34] J. G. Moreno-Torres, T. Raeder, R. Alaiz-Rodríguez, N. V. Chawla, and F. Herrera, “A unifying view on dataset shift in classification,” *Pattern recognition*, vol. 45, no. 1, pp. 521–530, 2012.



Qianwen Li is an Assistant Professor in the School of Environmental, Civil, Agricultural and Mechanical Engineering at the University of Georgia. She obtained her Ph.D. (2022) and M.S. (2020) in transportation engineering at the University of South Florida. And she received her B.S. degree in computer science and technology from Shandong University, China, in 2018. Her main research interests are intelligent transportation systems and transportation safety.



Xiaopeng Li is currently a Professor in the Department of Civil and Environmental Engineering at the University of Wisconsin-Madison. His major research interests include automated vehicle traffic control and connected & interdependent infrastructure systems. He has published around 70 peer-reviewed journal papers. Dr. Li received a B.S. degree (2006) in civil engineering with a computer engineering minor from Tsinghua University, China, an M.S. degree (2007), and a Ph.D. (2011) degree in civil engineering along with an M.S. degree (2010) in applied mathematics from the University of Illinois at Urbana-Champaign, USA.



Handong Yao is an Assistant Professor in the School of Environmental, Civil, Agricultural and Mechanical Engineering at the University of Georgia. He received his Ph.D. (2020) and M.S. (2015) degrees in transportation engineering at the Harbin Institute of Technology, Harbin, China. His main research interests are intelligent transportation systems and transportation safety.



Zhaohui Liang is a Ph.D. student in the Department of Civil and Environmental Engineering at the University of Wisconsin-Madison. He obtained his bachelor's degree from the Harbin Institute of Technology, Weihai, China in June 2019. His main research interests are connected autonomous vehicles and battery management strategies.



Weijun Xie is an Assistant Professor in the H. Milton Stewart School of Industrial and Systems Engineering at Georgia Tech. Dr. Xie obtained his Ph.D. in Operations Research at the Georgia Institute of Technology in August 2017. His research interests lie in theory and applications of stochastic, discrete, and convex optimization. His works have received multiple awards, such as 2021 NSF CAREER Award, Winner of 2020 INFORMS Young Researchers Paper Prize, Runner-up of Dupacova-Prekopa Best Student Paper Prize in Stochastic Programming at ICSP 2019, Honorable Mention in George Nicholson Student Paper Competition at INFORMS 2017.