

Adaptive Safety Evaluation for Connected and Automated Vehicles with Sparse Control Variates

Jingxuan Yang^{ID}, Haowei Sun^{ID}, Honglin He^{ID}, Yi Zhang^{ID}, Member, IEEE,
Shuo Feng^{ID}, Member, IEEE and Henry X. Liu^{ID}, Member, IEEE

Abstract—Safety performance evaluation is critical for developing and deploying connected and automated vehicles (CAVs). One prevailing way is to design testing scenarios using prior knowledge of CAVs, test CAVs in these scenarios, and then evaluate their safety performances. However, significant differences between CAVs and prior knowledge could severely reduce the evaluation efficiency. Towards addressing this issue, most existing studies focus on the adaptive design of testing scenarios during the CAV testing process, but so far they cannot be applied to high-dimensional scenarios. In this paper, we focus on the adaptive safety performance evaluation by leveraging the testing results, after the CAV testing process. It can significantly improve the evaluation efficiency and be applied to high-dimensional scenarios. Specifically, instead of directly evaluating the unknown quantity (e.g., crash rates) of CAV safety performances, we evaluate the differences between the unknown quantity and known quantity (i.e., control variates). By leveraging the testing results, the control variates could be well designed and optimized such that the differences are close to zero, so the evaluation variance could be dramatically reduced for different CAVs. To handle the high-dimensional scenarios, we propose the sparse control variates method, where the control variates are designed only for the sparse and critical variables of scenarios. According to the number of critical variables in each scenario, the control variates are stratified into strata and optimized within each stratum using multiple linear regression techniques. We justify the proposed method’s effectiveness by rigorous theoretical analysis and empirical study of high-dimensional overtaking scenarios.

Index Terms—Adaptive safety evaluation, connected and automated vehicles, sparse control variates, high-dimensional scenarios

I. INTRODUCTION

TESTING and evaluation of safety performance are major challenges for the development and deployment of connected and automated vehicles (CAVs). One proposed way is to test CAVs in the naturalistic driving environments (NDE) through a combination of software simulation, test tracks, and public roads, observe their performances, and make statistical

This work is supported by National Key Research and Development Program under Grant 2021YFB2501200 and National Natural Science Foundation of China under Grant 62133002. (*Corresponding author: Shuo Feng*)

Jingxuan Yang and Honglin He are with the Department of Automation, Tsinghua University, Beijing 100084, China (email: {yangjx20, hehl21}@mails.tsinghua.edu.cn).

Haowei Sun and Henry X. Liu are with the Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI 48109, USA (e-mail: {haoweis, henryliu}@umich.edu).

Yi Zhang is with the Department of Automation, Beijing National Research Center for Information Science and Technology (BNRist), Tsinghua University, Beijing 100084, China (e-mail: zhyi@tsinghua.edu.cn).

Shuo Feng is with the Department of Automation, Tsinghua University, Beijing 100084, China and the University of Michigan Transportation Research Institute, Ann Arbor, MI 48109, USA (e-mail: fshuo@umich.edu).

comparisons with human drivers. Due to the rarity of safety-critical events in NDE, however, hundreds of millions of miles and sometimes hundreds of billions of miles would be required to demonstrate CAVs’ safety performance at the human-level [1], which is intolerably inefficient. To improve the efficiency and accelerate the evaluation process, the past few years have witnessed increasingly rapid advances in the field of testing scenario library generation (TSLG) [2]–[10], where safety-critical testing scenarios are usually purposely generated utilizing prior knowledge of CAVs such as surrogate models (SMs) of CAVs. However, due to the high complexity and black-box properties of CAVs, there exist significant performance dissimilarities between SMs and CAVs under test, which could severely compromise the effectiveness of the generated testing scenarios and decrease the evaluation efficiency.

Towards addressing this problem, several adaptive testing and evaluation methods have been proposed [11]–[14]. The basic idea of existing methods is to adaptively generate the testing scenarios during the testing process of CAVs. With more testing results of CAVs, more posteriori knowledge of CAVs can be obtained, and therefore the testing scenarios can be more customized and optimized for the CAVs under test. However, most existing methods can only be applied to relatively simple scenarios, and how to handle high-dimensional scenarios remains an open question. For example, Mullins *et al.* [11] proposed an adaptive sampling method that uses Gaussian process regression (GPR) and k -nearest neighbors to discover performance boundaries of the system under test and then updates the SM with new testing results obtained near the performance boundaries. Koren *et al.* [12] put forward an adaptive stress testing method that uses deep reinforcement learning to find the most-likely failure scenarios. Feng *et al.* [13] proposed an adaptive testing scenario library generation method using Bayesian optimization techniques with classification-based GPR and acquisition functions to select subsequent testing scenarios and then update the SMs with new testing results. Sun *et al.* [14] presented an adaptive design of experiments method to detect safety-critical scenarios, which uses supervised machine learning models as SMs to approximate the testing results and devises acquisition functions for updating the SMs.

The challenge for adaptively generating high-dimensional scenarios comes from the compounding effects of the “Curse of Rarity” (CoR) and the “Curse of Dimensionality” (CoD) [15]. The CoR refers to the concept that, due to rarity of safety-critical events, the amount of data needed to obtain sufficient

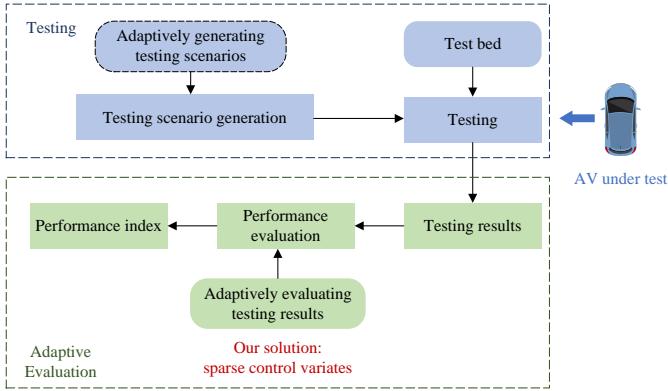


Fig. 1. Illustration of the adaptive testing and evaluation framework. The focus of this study is the adaptive evaluation method for high-dimensional scenarios, where the sparse control variates method is proposed.

information grow dramatically, while the CoD refers to the dimensionality of variables to represent realistic scenarios, which makes the computation cost increase exponentially with the growth of scenario dimensions. Most existing scenario-based testing approaches can only handle short scenario segments with limited background road users, where the decision variables are low-dimensional, which cannot represent the full complexity and variability of the real-world driving environment [16]–[20]. Towards addressing this challenge, the naturalistic and adversarial driving environment (NADE) method has been developed in our previous work [21], which can generate high-dimensional highway driving scenarios. However, the NADE did not consider the performance gap between CAVs and SMs, which could also slow down the testing process. To the best of the authors' knowledge, there is no existing work that can handle the adaptive testing and evaluation problem in high-dimensional scenarios, and the goal of this paper is to fill this gap.

In general, the adaptive testing and evaluation methods can be categorized into two types including adaptive testing scenario generation and adaptive testing result evaluation, which are complementary to each other as shown in Fig. 1. Most existing studies focus on the former one, while in this study, we focus on the latter one and propose an adaptive evaluation framework that can handle high-dimensional scenarios. We note that how to realize the former one in high-dimensional scenarios also remains unsolved, which we leave for future study. In the proposed framework, we apply the NADE method to generate high-dimensional testing scenarios, where combinations of multiple SMs are utilized to improve the robustness of the generated scenarios for different CAVs under test. Then we propose a sparse control variate (SCV) method to adjust the testing results and evaluate CAVs' performance adaptively. Essentially, the SCV method could reduce the estimation variance for the CAV under test and thus reduce the required number of tests, accelerating the evaluation process adaptively.

In the following paragraphs, we further explain the major idea of the proposed SCV method. The control variates (CV) method [22] is a popular variance reduction technique applied in research areas such as deep learning [23] and reinforcement

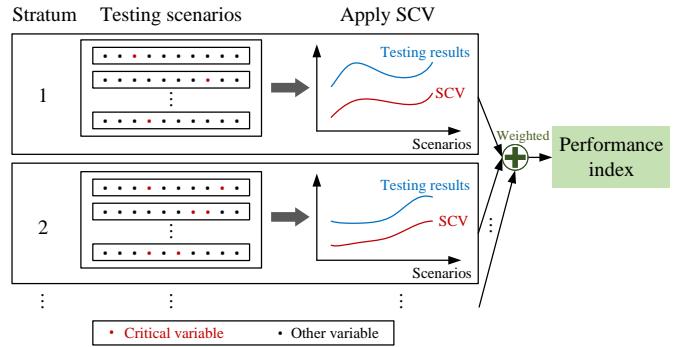


Fig. 2. Illustration of the sparse control variates method. The SCV are constructed by only considering critical variables (represented as red dots in testing scenarios). The testing results are stratified into strata according to the number of critical variables and then adjusted by SCV within each stratum. Finally, the performance index are obtained by summing up these evaluation results with proportion weights.

learning [24]. Suppose we want to estimate $\mu \triangleq \mathbb{E}_p[f(X)]$ by Monte Carlo sampling [25], where p is the probabilistic distribution of the random variable X and f is the performance index of interest. Instead of directly estimating the unknown quantity μ , the control variates method estimates the differences between the unknown quantity and known quantity as $\mu' \triangleq \mathbb{E}_p[f(X) - h(X) + \theta]$, where $h(X)$ is the control variate and $\theta \triangleq \mathbb{E}_p[h(X)]$ is a known value. Then, if $h(X)$ correlates with the performance index $f(X)$ (hence can provide some information about $f(X)$), the estimation variance of μ' will always be less than directly estimating μ [26]. For testing and evaluation of CAVs, the control variate $h(X)$ can be designed by utilizing the prior knowledge of CAVs (e.g., different SMs). $h(X)$ usually contains adjustable control parameters, which can be optimized by leveraging the testing results. In such way, the information about the CAV under test could be incorporated, which makes the adaptive evaluation possible. However, due to the CoD, the computation cost of optimal control parameters will increase exponentially with the growth of scenario dimensions, so directly applying the ordinary CV method in high-dimensional scenarios is problematic.

To address this problem, we propose the sparse control variates (SCV) method, as shown in Fig. 2. The key idea is to construct the SCV by only considering the sparse but critical variables (e.g., behaviors of principal other vehicles at critical moments), following the similar idea from [21] that handles the CoD. However, the number of critical variables varies in different testing scenarios, which cannot be handled by ordinary CV method. To address this issue, in the SCV method, we stratify the testing scenarios into strata according to the number of critical variables. Then the control parameters can be optimized by multiple linear regression (MLR) [27] within each stratum, and the final evaluation results are obtained by summing up those evaluation results in each stratum with the proportion weights. Since the number of critical variables is much less than the dimension of testing scenarios, the computation cost of optimal control parameters for SCV could be greatly reduced, overcoming the CoD challenge.

To verify the proposed method, we theoretically analyze

its accuracy, efficiency, and optimality. The theorems show that our method is unbiased, and its estimation variance is nearly proportional to the best one that all the SMs used for generating testing scenarios could have. Moreover, under certain assumptions about the SMs, our method can provide a zero-variance estimator. To validate our method, the high-dimensional overtaking scenarios with large-scale naturalistic driving data are investigated. Simulation results show that our method can further accelerate the evaluation process by about one order of magnitude for different types of CAV models, comparing with the estimation efficiency in NADE.

Compared with our previously published conference paper about SCV [28], the new contributions of this paper are listed as follows. First, we significantly extend our methodology into high-dimensional scenarios and establish the theoretical analysis for the accuracy, efficiency, and optimality of the proposed method with rigorous proofs. Second, a more realistic overtaking case study with large-scale naturalistic driving data is investigated to systematically validate the performances of our method.

The remainder of this paper is organized as follows. Section II provides preliminary knowledge for the generation of NDE and NADE. Section III formulates the adaptive testing and evaluation problem and elaborates the challenges of applying ordinary CV for adaptive safety evaluation. To address these challenges, in Section IV, the SCV method is proposed. Then Section V and VI verify and validate the accuracy and efficiency of the proposed method from the theoretical and experimental perspectives, respectively. Finally, Section VII concludes the paper and discusses future research.

II. PRELIMINARIES

A. Naturalistic Driving Environment Testing

As discussed above, the prevailing approach for CAV evaluation is to test CAVs in the naturalistic driving environments (NDE) [29], observe their performances, and make statistical comparisons with human drivers. In NDE, one of the vehicles is the automated vehicle (AV) under test and the others are background vehicles (BVs), which can be formulated as Markov games [30]. A Markov game for N agents (i.e., BVs) is defined by a set of states \mathcal{S} describing the positions and velocities of all vehicles and a collection of action (i.e., acceleration) sets $\mathcal{A}_1, \dots, \mathcal{A}_N$, one for each agent in NDE. The total action space is denoted as $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_N$. Then a scenario is defined as the time series of the states of all vehicles and the actions of all agents, i.e.,

$$x = (s_0, a_0, \dots, s_T, a_T) \in \mathcal{X}, \quad (1)$$

where x represents the scenario, \mathcal{X} is the set of all feasible scenarios, $s_t \in \mathcal{S}$ is the state of all vehicles at time t , $a_t \in \mathcal{A}$ is the action of all agents at time t , and T is the time horizon.

Let $\Omega = \mathcal{X}$ be the sample space incorporating all feasible scenarios. Consider the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where $\mathcal{F} \triangleq 2^\Omega$ is the power set of Ω and \mathbb{P} is a probability measure on \mathcal{F} . Let $X : x \mapsto x, \forall x \in \mathcal{X}$ be the random variable of scenarios. For testing and evaluation of CAVs, the crash event is usually of most interest, which can be defined as

$A = \{x \in \mathcal{X} : s_T \in \mathcal{S}_c\}$, where \mathcal{S}_c is the set of all crash states. Then the crash rate is selected as the performance index, which can be computed as

$$\mu = \mathbb{P}(A) = \mathbb{E}_p[\mathbb{I}_A(X)] = \sum_{x \in \mathcal{X}} \mathbb{P}(A|x)p(x), \quad (2)$$

where \mathbb{I}_A is the indicator function of A , and p is the naturalistic joint distribution of x . The essence of testing AV in NDE is to estimate the performance index μ by Monte Carlo simulation, i.e.,

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(A|X_i), \quad X_i \sim p. \quad (3)$$

B. Naturalistic and Adversarial Driving Environment Generation

The NDE faces the CoR, making its estimation catastrophically inefficient. To improve the estimation efficiency, the importance sampling (IS) technique [17]–[19] has been used to sample testing scenarios from the importance function q , which puts more weights on crash-prone scenarios. In IS, the performance index can be estimated as

$$\hat{\mu}_q = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)p(X_i)}{q(X_i)}, \quad X_i \sim q. \quad (4)$$

However, the IS method faces the CoD if the testing scenarios are high-dimensional [31]. To address both the CoR and the CoD, the naturalistic and adversarial driving environment (NADE) [21] has been proposed to only sample critical variables of testing scenarios from importance functions, while other variables remain their naturalistic distributions.

Denote $x = (x_c, x_{-c})$, where $x_c = \{x_{c_1}, \dots, x_{c_L}\}$ is the set of critical variables, c_1, \dots, c_L are called the critical moments, $L = 0, 1, \dots, L$ is the number of control steps (i.e., the number of critical variables in x_c), and x_{-c} is the set of other variables. Let $X_c : x \mapsto x_c$ be the random variable of critical variables and $X_{-c} : x \mapsto x_{-c}$ be the random variable of other variables, then we have $X = (X_c, X_{-c})$. The importance function can then be formulated as $q(x) = q(x_c)p(x_{-c})$, and therefore the performance index can be estimated in NADE as

$$\tilde{\mu}_q = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)p(X_{c,i})}{q(X_{c,i})}, \quad X_i \sim q, \quad (5)$$

where $X_{c,i}$ is the random variable of critical variables of X_i .

III. PROBLEM FORMULATION

A. Adaptive Testing and Evaluation

Due to the black-box property and various types of CAVs, how to adaptively test and evaluate CAVs remains a major challenge. One way of adaptive testing and evaluation is adaptively generating testing scenarios. For example, we can minimize the estimation variance by optimizing the importance function, i.e.,

$$\min_{q \in \mathcal{Q}} \text{Var}_q \left(\frac{\mathbb{P}(A|X)p(X)}{q(X)} \right), \quad (6)$$

where \mathcal{Q} is the function space of q . Better importance functions can be found by leveraging the posteriori knowledge of

CAVs obtained from testing results. Then the testing scenarios can be adaptively generated by sampling from updated importance functions.

In this paper, we focus on another way of adaptive testing and evaluation, i.e., adaptively evaluating weighted testing results. Specifically, the control variates (CV) method is adopted. This problem can be formulated as

$$\min_{h \in \mathcal{H}} \text{Var}_q \left(\frac{\mathbb{P}(A|X)p(X)}{q(X)} - h(X) \right), \quad (7)$$

where $h : \mathcal{X} \rightarrow \mathbb{R}$ is the control variate and \mathcal{H} is the function space of h . The goal is to further reduce the estimation variance by optimizing h in \mathcal{H} , leveraging the testing results.

B. Control Variates

Control variates are widely used as a basic variance reduction technique in Monte Carlo simulation. They can be usefully combined with the mixture importance sampling, where individual importance functions can serve as CV. In mixture IS, the scenarios $X_i, i = 1, \dots, n$ are sampled from the mixture importance function $q_\alpha = \sum_{j=1}^J \alpha_j q_j$, where $\alpha_j \geq 0$, $\sum_{j=1}^J \alpha_j = 1$ and the q_j are importance functions. One commonly used way is to construct CV by using the linear combination of individual importance functions as

$$h_\beta(X) = \sum_{j=1}^J \beta_j \left[\frac{q_j(X)}{q_\alpha(X)} - 1 \right], \quad (8)$$

where $\beta = (\beta_1, \dots, \beta_J)^\top$ is the control vector, $\beta_j \in \mathbb{R}$ are control parameters, and $q_j/q_\alpha - 1$ are individual control variate. Combining the control variate h_β with mixture IS gives the estimation

$$\hat{\mu}_{q_\alpha, \beta} = \frac{1}{n} \sum_{i=1}^n \left[\frac{\mathbb{P}(A|X_i)p(X_i)}{q_\alpha(X_i)} - h_\beta(X_i) \right] \quad (9)$$

for $X_i \sim q_\alpha$.

The unbiasedness of $\hat{\mu}_{q_\alpha, \beta}$ is guaranteed since

$$\mathbb{E}_{q_\alpha} [\hat{\mu}_{q_\alpha, \beta}] = \mathbb{E}_{q_\alpha} \left[\frac{\mathbb{P}(A|X)p(X)}{q_\alpha(X)} - h_\beta(X) \right] = \mu, \quad (10)$$

where the second equality is obtained from the unbiasedness of IS and $\mathbb{E}_{q_\alpha} [h_\beta(X)] = 0$. The variance of $\hat{\mu}_{q_\alpha, \beta}$ can be compared to that of IS with individual importance functions q_j . We have the following lemma.

Lemma 1: Let β^* be any minimizer over β of $\text{Var}_{q_\alpha}(\hat{\mu}_{q_\alpha, \beta})$, then

$$\text{Var}_{q_\alpha}(\hat{\mu}_{q_\alpha, \beta^*}) \leq \min_{1 \leq j \leq J} \frac{\sigma_{q_j}^2}{n\alpha_j}, \quad (11)$$

where $\sigma_{q_j}^2$ is the asymptotic variance of $\hat{\mu}_{q_j}$, i.e.,

$$\sigma_{q_j}^2 = \text{Var}_{q_j} \left(\frac{\mathbb{P}(A|X)p(X)}{q_j(X)} \right), \quad j = 1, \dots, J. \quad (12)$$

Proof: This is the Theorem 2 in [32]. \square

It can be seen from Lemma 1 that the variance of $\hat{\mu}_{q_\alpha, \beta}$ will be zero if any one of the q_j is optimal. This is a significant feature because we can nearly omit the influence of all other worse-performed importance functions. In applications, using

only one SM to test CAVs is usually under huge risk, because the performance gap between the SM and various types of CAVs may be too large to give a good estimation efficiency. Therefore, to ensure the robustness, we can combine multiple SMs to test the CAVs. However, there often exist some poor-performed SMs that will compromise the overall estimation efficiency. Using mixture IS with CV provides an effective way to ensure both good estimation efficiency and robustness to various types of CAVs.

In practice, the optimal control vector β^* is usually unknown, and its estimation $\hat{\beta}$ can be obtained by multiple linear regression (MLR). Denote the weighted testing results as $Y_i = \mathbb{P}(A|X_i)p(X_i)/q_\alpha(X_i)$, $i = 1, \dots, n$, and the individual control variate as $Z_{ij} = q_j(X_i)/q_\alpha(X_i) - 1$, $i = 1, \dots, n$, $j = 1, \dots, J - 1$. Then the $\hat{\beta}$ is given as the vector of coefficients obtained from MLR of Y_i on Z_{ij} . In essence, this process is to search for the best control variate defined in Eq. (8) in the function space spanned by individual control variate $q_j/q_\alpha - 1$. However, challenges of estimating optimal control parameters arise when the testing scenarios are high-dimensional.

C. CoD of Control Variates

Considering the Markov chain structure of scenarios with $T + 1$ time steps, the mixture importance function is given by

$$q_\alpha(x) = q_\alpha(s_0) \prod_{t=0}^T q_\alpha(a_t | s_t), \quad \forall x \in \mathcal{X}, \quad (13)$$

where $q_\alpha(s) = \sum_{j=1}^J \alpha_j q_j(s)$, $\forall s \in \mathcal{S}$, and $q_\alpha(a|s) = \sum_{j=1}^J \alpha_j q_j(a|s)$, $\forall a \in \mathcal{A}$, $s \in \mathcal{S}$. It can be found that $q_\alpha(x)$ is the product of $T + 2$ individual importance functions and thus is also the summation of J^{T+2} combinations of different importance functions at each time step. Specifically, these individual importance functions are

$$q_{j_0, \dots, j_{T+1}}(x) = q_{j_0}(s_0)q_{j_1}(a_0 | s_0) \cdots q_{j_{T+1}}(a_T | s_T), \quad (14)$$

where $j_0, \dots, j_{T+1} = 1, \dots, J$. Then the individual control variate are given by $q_{j_0, \dots, j_{T+1}}/q_\alpha - 1$.

To find the estimation of optimal control parameters, we have to conduct MLR of n weighted testing results on J^{T+2} individual control variate. The number J^{T+2} will increase exponentially with the dimension of scenarios, leading to the CoD of MLR. For example, if we have $J = 10$ individual importance functions and the testing scenarios last for 10 seconds at a frequency of 10 Hz, then the number of individual control variate will be 10^{102} . This means that a matrix with dimension 10^{102} should be inverted in MLR, which is not tractable. Moreover, the situation will get even worse if the duration of scenarios grows to several hours, which are common in daily driving yet far from being tractable. The following section aims to address this challenge.

IV. ADAPTIVE SAFETY EVALUATION WITH SPARSE CONTROL VARIATES

In this section, we will address the CoD discussed above and show how to estimate the optimal control parameters.

A. Sparse Control Variates

We propose the sparse control variates (SCV) method to address the CoD of applying CV in high-dimensional scenarios. Specifically, the SCV are constructed by only considering the importance functions of only sparse and critical variables in high-dimensional testing scenarios. The number of critical variables is usually much less than the dimension of scenarios in NADE. Therefore, the number of SCV is also much less than the number of ordinary CV, which could greatly address the CoD. However, as the number of SCV varies in different testing scenarios, we can not directly apply SCV to the weighted testing results. Towards addressing this issue, we propose to stratify the testing scenarios into strata according to the number of critical variables and then apply SCV within each stratum.

Let $\mathcal{X}_l = \{x \in \mathcal{X} : |x_c| = l\}$, $l = 0, 1, \dots, L$ be the stratum of scenarios that are controlled l steps, satisfying $\bigcup_{l=0}^L \mathcal{X}_l = \mathcal{X}$. Using mixture importance function q_α , the estimation of the performance index in NADE is

$$\tilde{\mu}_{q_\alpha} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)p(X_{c,i})}{q_\alpha(X_{c,i})}, \quad X_i \sim q_\alpha. \quad (15)$$

The performance index of scenarios in stratum \mathcal{X}_l can be written as $\mu_l \triangleq \mathbb{E}_p[\mathbb{I}_A(X)\mathbb{I}_{\mathcal{X}_l}(X)]$, $l = 0, 1, \dots, L$, then we have

$$\mu = \sum_{l=0}^L \mathbb{E}_p[\mathbb{I}_A(X)\mathbb{I}_{\mathcal{X}_l}(X)] = \sum_{l=0}^L \mu_l. \quad (16)$$

Similar to Eq. (15), the estimation of μ_l is given by

$$\tilde{\mu}_{l,q_\alpha} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)\mathbb{I}_{\mathcal{X}_l}(X_i)p(X_{c,i})}{q_\alpha(X_{c,i})}, \quad (17)$$

and then we have

$$\tilde{\mu}_{q_\alpha} = \sum_{l=0}^L \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)\mathbb{I}_{\mathcal{X}_l}(X_i)p(X_{c,i})}{q_\alpha(X_{c,i})} = \sum_{l=0}^L \tilde{\mu}_{l,q_\alpha}. \quad (18)$$

Let $q_{j_1, \dots, j_l}(x) = p(x_{-c})q_{j_1}(x_{c_1}) \cdots q_{j_l}(x_{c_l})$ be the importance functions that sample x_{-c} from p and sample x_{c_1}, \dots, x_{c_l} from q_{j_1}, \dots, q_{j_l} respectively, where $j_1, \dots, j_l = 1, \dots, J$, $l = 1, \dots, L$. Then the individual importance functions of critical variables are given by $q_{j_1, \dots, j_l}(x_c)$. Denote the linear combination of these individual importance functions as

$$\tilde{h}_l(x) \triangleq \sum_{j_1, \dots, j_l} \beta_{l,j_1, \dots, j_l} q_{j_1, \dots, j_l}(x), \quad l = 1, \dots, L, \quad (19)$$

where $\beta_{l,j_1, \dots, j_l} \in \mathbb{R}$ are associated control parameters. Then the SCV are given by

$$h_l(x_c) = \frac{\tilde{h}_l(x_c)\mathbb{I}_{\mathcal{X}_l}(x_c)}{q_\alpha(x_c)} - \theta_l, \quad l = 1, \dots, L, \quad (20)$$

where $\theta_l \triangleq \mathbb{E}_{q_\alpha}[\tilde{h}_l(X)\mathbb{I}_{\mathcal{X}_l}(X)/q_\alpha(X)]$. Therefore, the estimation $\tilde{\mu}_{l,q_\alpha}$ in Eq. (17) can be evaluated with SCV as

$$\begin{aligned} \tilde{\mu}_{l,q_\alpha, \beta_l} &= \frac{1}{n} \sum_{i=1}^n \left[\frac{\mathbb{P}(A|X_i)\mathbb{I}_{\mathcal{X}_l}(X_i)p(X_{c,i})}{q_\alpha(X_{c,i})} - h_l(X_{c,i}) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(A|X_i)p(X_{c,i}) - \tilde{h}_l(X_{c,i})}{q_\alpha(X_{c,i})} \mathbb{I}_{\mathcal{X}_l}(X_i) + \theta_l \end{aligned} \quad (21)$$

for $l = 1, \dots, L$, where $\beta_l = \text{vec}(\beta_{l,j_1, \dots, j_l})$ is the vector of control parameters, and $\text{vec}(\cdot)$ is the vectorization operator that flattens a tensor into a long vector. Note that there is no critical variable for $l = 0$, and thus we set $\beta_0 \triangleq 0$. In summary, the performance index estimated by the proposed SCV method is given by

$$\tilde{\mu}_{q_\alpha, \beta} = \sum_{l=0}^L \tilde{\mu}_{l,q_\alpha, \beta_l}, \quad (22)$$

where $\beta = \{\beta_l\}_{l=0}^L$ is the set of all control vectors.

B. Optimal Control Parameters

To estimate the optimal control parameters that minimize the estimation variance, multiple linear regression (MLR) technique is applied in each stratum. Let $\mathbb{X}_l \triangleq \{X_i | X_i \in \mathcal{X}_l, i = 1, \dots, n\}$ be the set of sampled scenarios with l controlled steps, $n_l \triangleq \sum_{i=1}^n \mathbb{I}_{\mathcal{X}_l}(X_i)$ be the number of tests with l controlled steps and $d_l \triangleq J^l$ be the number of SCV, $l = 1, \dots, L$. Denote the vector of testing results as

$$Y_l \triangleq \left[\frac{\mathbb{P}(A|X_i)p(X_i)}{q_\alpha(X_i)} \text{ for } X_i \in \mathbb{X}_l \right] \in \mathbb{R}^{n_l}, \quad (23)$$

the individual SCV as

$$h'_{j_1, \dots, j_l}(x_c) = \frac{q_{j_1, \dots, j_l}(x_c)}{q_\alpha(x_c)} - \sum_{x_c \in \mathcal{X}_l} q_{j_1, \dots, j_l}(x_c), \quad (24)$$

for $l = 1, \dots, L$. Then the matrix of individual SCV can be formulated as

$$H_l \triangleq [\text{vec}(h'_{j_1, \dots, j_l}(X_{c,i})) \text{ for } X_i \in \mathbb{X}_l] \in \mathbb{R}^{n_l \times d_l}, \quad (25)$$

for $l = 1, \dots, L$. Then the regression formula is given by $Y_l \approx \eta_l + H_l \beta_l$. The MLR of Y_l on H_l is to find the optimal solution of the following optimization problem, i.e.,

$$\min_{\eta_l, \beta_l} f(\eta_l, \beta_l) = \|Y_l - \eta_l - H_l \beta_l\|_2^2. \quad (26)$$

Letting the partial derivatives of f with respect to η_l and β_l both equal zero, we have $\hat{\eta}_l = 1^\top Y_l / n_l$ and $\hat{\beta}_l = (H_l^\top H_l)^{-1} H_l^\top Y_l$, assuming that the control matrix $M_l \triangleq H_l^\top H_l \in \mathbb{R}^{d_l \times d_l}$ is invertible. Then the estimated performance index is $\hat{\mu}_l = n_l \hat{\eta}_l / n$. In practice the control matrix may often not be invertible, then we use singular value decomposition (SVD) [33] to compute the regression coefficients $\hat{\beta}_l$, and the rank of the control matrix is

$$\text{rank}(M_l) = \text{rank}(H_l) \leqslant \min\{n_l, d_l\}. \quad (27)$$

If $n_l < d_l$, then the control matrix M_l will be singular and has utmost n_l nonzero singular values. As the number of tests n_l in \mathbb{X}_l will not grow exponentially with the number

of control steps l , the rank of the control matrix will also not, albeit the dimension $d_l = J^l$ of the control matrix increases exponentially with l . In conclusion, solving the optimal control parameters for SCV is tractable and will not face the CoD challenge. We will further demonstrate this in Subsection VI-E.

V. THEORETICAL ANALYSIS

This section theoretically justifies the accuracy, efficiency and optimality of the proposed SCV method.

A. Accuracy Analysis

We first prove that the estimation is unbiased.

Theorem 1: Let $\tilde{\mu}_{q_\alpha, \beta}$ be given by Eq. (22) where $q_\alpha > 0$ whenever $\mathbb{P}(A|x)p(x) > 0$, then $\mathbb{E}_{q_\alpha}[\tilde{\mu}_{q_\alpha, \beta}] = \mu$.

Proof: To establish unbiasedness, write

$$\begin{aligned}\mathbb{E}_{q_\alpha}[\tilde{\mu}_{q_\alpha, \beta}] &= \mathbb{E}_{q_\alpha}\left[\sum_{l=0}^L \tilde{\mu}_{l, q_\alpha, \beta_l}\right] \\ &= \sum_{l=0}^L \mathbb{E}_{q_\alpha}\left[\tilde{\mu}_{l, q_\alpha} - \frac{\tilde{h}_l(X)}{q_\alpha(X)} \mathbb{I}_{\mathcal{X}_l}(X) + \theta_l\right] \quad (28) \\ &= \sum_{l=0}^L (\mu_l - \theta_l + \theta_l) = \mu.\end{aligned}$$

□

Remark 1: This theorem indicates that the estimation is unbiased if the control parameters β are independent of the sample data. It's worth noting that in practice the control parameters are usually estimated by the sample data, which would bring a bias. However, that bias is ordinarily negligible (please see Section 8.9 in [26] for more discussions).

B. Efficiency Analysis

Next, we evaluate the efficiency of the SCV method. The variance of the estimation $\tilde{\mu}_{q_\alpha, \beta}$ is $\text{Var}_{q_\alpha}(\tilde{\mu}_{q_\alpha, \beta}) = \sigma_{q_\alpha, \beta}^2/n$, where $\sigma_{q_\alpha, \beta}^2$ is the asymptotic variance of $\tilde{\mu}_{q_\alpha, \beta}$, i.e.,

$$\sigma_{q_\alpha, \beta}^2 = \text{Var}_{q_\alpha}\left(\sum_{l=0}^L \frac{\mathbb{P}(A|X)p(X) - \tilde{h}_l(X)}{q_\alpha(X)} \mathbb{I}_{\mathcal{X}_l}(X)\right) \quad (29)$$

for $X \sim q_\alpha$. Denote

$$Z_l \triangleq \frac{\mathbb{P}(A|X)p(X) - \tilde{h}_l(X)}{q_\alpha(X)} \mathbb{I}_{\mathcal{X}_l}(X), \quad l = 0, \dots, L, \quad (30)$$

then the asymptotic variance $\sigma_{q_\alpha, \beta}^2$ can be expressed as

$$\sigma_{q_\alpha, \beta}^2 = \text{Var}_{q_\alpha}\left(\sum_{l=0}^L Z_l\right) = \mathbb{E}_{q_\alpha}\left[\left(\sum_{l=0}^L [Z_l - \mathbb{E}_{q_\alpha}[Z_l]]\right)^2\right]. \quad (31)$$

Let $L' = L + 1$, then by convexity of quadratic function and Jensen's inequality, we have

$$\begin{aligned}\sigma_{q_\alpha, \beta}^2 &\leq \mathbb{E}_{q_\alpha}\left[L' \sum_{l=0}^L (Z_l - \mathbb{E}_{q_\alpha}[Z_l])^2\right] \\ &= L' \sum_{l=0}^L \text{Var}_{q_\alpha}(Z_l).\end{aligned} \quad (32)$$

Denote $\sigma_{l, q_\alpha, \beta_l}^2 \triangleq \text{Var}_{q_\alpha}(Z_l)$ and the asymptotic variance of $\tilde{\mu}_{l, q}$ over \mathcal{X}_l as $\sigma_{l, q}^2$, i.e.,

$$\sigma_{l, q}^2 \triangleq \sum_{x \in \mathcal{X}_l} \left(\frac{\mathbb{P}(A|x)p(x)}{q(x)} - \mu_l \right)^2 q(x), \quad l = 1, \dots, L, \quad (33)$$

then we have the following theorem.

Theorem 2: If β^* is any minimizer of $\sigma_{q_\alpha, \beta}^2$, then

$$\begin{aligned}\sigma_{q_\alpha, \beta^*}^2 &\leq L' \sigma_{0, p, \beta_0}^2 \\ &+ L' \sum_{l=1}^L \min_{j_1, \dots, j_l} \left\{ \frac{\sigma_{l, q_{j_1}, \dots, j_l}^2}{\prod_{\ell=1}^l \alpha_{j_\ell}} + 3 \left(\frac{\mu_l}{\prod_{\ell=1}^l \alpha_{j_\ell}} \right)^2 \right\}.\end{aligned} \quad (34)$$

Proof: Take $\sigma_{1, q_\alpha, \beta_1}^2$ as an example. Following the proof in [32], we consider the particular vector β_1 having $\beta_{1,1} = 0$ and $\beta_{1,j} = -\mu_1 \alpha_j / \alpha_1$ for $j > 1$. Let $r_1(x) \triangleq [\mathbb{P}(A|x)p(x) - \mu_1 q_1(x)] \mathbb{I}_{\mathcal{X}_1}(x)$, then we have $\sum_{x \in \mathcal{X}} r_1(x) = \mu_1(1 - \xi_1)$, where $\xi_1 \triangleq \sum_{x \in \mathcal{X}_1} q_1(x)$, $\xi_1 \in [0, 1]$. Substituting these values, we find that for this β_1 ,

$$\begin{aligned}Z_1 &= \frac{\mathbb{P}(A|X)p(X) - \tilde{h}_1(X)}{q_\alpha(X)} \mathbb{I}_{\mathcal{X}_1}(X) \\ &= \frac{\mathbb{P}(A|X)p(X) - \mu_1 q_1 + \mu_1 q_1 - \tilde{h}_1(X)}{q_\alpha(X)} \mathbb{I}_{\mathcal{X}_1}(X) \quad (35) \\ &= \frac{r_1(X)}{q_\alpha(X)} + \frac{\mu_1}{\alpha_1} \mathbb{I}_{\mathcal{X}_1}(X),\end{aligned}$$

and $\mathbb{E}_{q_\alpha}[Z_1] = \mu_1 \alpha_{1,1} / \alpha_1$, where $\alpha_{1,1} \triangleq \alpha_1 + \sum_{j=2}^J \alpha_j \sum_{x \in \mathcal{X}_1} q_j(x)$, $\alpha_{1,1} \in [0, 1]$. Therefore, we have

$$\begin{aligned}\sigma_{1, q_\alpha, \beta_1}^2 &= \mathbb{E}_{q_\alpha}\left[\left(Z_1 - \mathbb{E}_{q_\alpha}[Z_1]\right)^2\right] \\ &= \sum_{x \in \mathcal{X}} \left[\frac{r_1(x)}{q_\alpha(x)} + \frac{\mu_1}{\alpha_1} (\mathbb{I}_{\mathcal{X}_1}(x) - \alpha_{1,1}) \right]^2 q_\alpha(x) \quad (36) \\ &\triangleq V_{1,1} + V_{1,2} + V_{1,3},\end{aligned}$$

where

$$\begin{aligned}V_{1,1} &\triangleq \sum_{x \in \mathcal{X}} \frac{r_1^2(x)}{q_\alpha(x)} = \sum_{x \in \mathcal{X}} \frac{[\mathbb{P}(A|x)p(x) - \mu_1 q_1(x)]^2}{q_\alpha(x)} \mathbb{I}_{\mathcal{X}_1}(x) \\ &\leq \sum_{x \in \mathcal{X}_1} \frac{[\mathbb{P}(A|x)p(x) - \mu_1 q_1(x)]^2}{\alpha_1 q_1(x)} = \frac{\sigma_{1, q_1}^2}{\alpha_1},\end{aligned} \quad (37)$$

$$\begin{aligned}V_{1,2} &\triangleq \sum_{x \in \mathcal{X}} \frac{2\mu_1 r_1(x)(\mathbb{I}_{\mathcal{X}_1}(x) - \alpha_{1,1})}{\alpha_1} \\ &= \frac{2\mu_1^2(1 - \xi_1)(1 - \alpha_{1,1})}{\alpha_1} \leq 2 \left(\frac{\mu_1}{\alpha_1} \right)^2,\end{aligned} \quad (38)$$

and

$$\begin{aligned}V_{1,3} &\triangleq \sum_{x \in \mathcal{X}} \left[\frac{\mu_1(\mathbb{I}_{\mathcal{X}_1}(x) - \alpha_{1,1})}{\alpha_1} \right]^2 q_\alpha(x) \\ &\leq \sum_{x \in \mathcal{X}} \left(\frac{\mu_1}{\alpha_1} \right)^2 q_\alpha(x) = \left(\frac{\mu_1}{\alpha_1} \right)^2.\end{aligned} \quad (39)$$

Therefore, we conclude that

$$\sigma_{1, q_\alpha, \beta_1^*}^2 \leq \sigma_{1, q_\alpha, \beta_1}^2 \leq \frac{\sigma_{1, q_1}^2}{\alpha_1} + 3 \left(\frac{\mu_1}{\alpha_1} \right)^2. \quad (40)$$

By making similar arguments for $j = 2, \dots, J$, we have

$$\sigma_{1,q_\alpha,\beta^*}^2 \leq \min_j \left\{ \frac{\sigma_{1,q_j}^2}{\alpha_j} + 3 \left(\frac{\mu_1}{\alpha_j} \right)^2 \right\}. \quad (41)$$

It's straightforward to extend the proof for $l = 2, \dots, L$, then Eq. (34) is established. \square

Remark 2: For $l = 1$, we expect to get approximately $n_1 \alpha_j$ scenarios in \mathcal{X}_1 from the importance function q_j . The quantity $\sigma_{1,q_j}^2 / \alpha_j$ in Eq. (41) is the variance we would obtain from $n_1 \alpha_j$ such scenarios alone. It is hard to imagine that we could do better in general, because when $\sigma_{1,q_j}^2 = \infty$ for all but one of the mixture components it is guaranteed that those bad components do not make the estimation worse than what we would have had from the one good importance function. Moreover, if there exists an optimal importance function in q_j , then the minimum value of $\sigma_{1,q_j}^2 / \alpha_j$ will be zero, which will greatly reduce the estimation variance. It should be noted that the upper bound for variance in Eq. (41) contains a residual term $3(\mu_1/\alpha_j)^2$, which is the cost for stratifying the scenarios.

C. Optimality Analysis

Under the following assumptions, the estimation variance of the SCV method can be zero.

Assumption 1: The scenarios in \mathcal{X}_0 will not be sampled by q_α , i.e., $q_\alpha(x) = 0, \forall x \in \mathcal{X}_0$.

Assumption 2: The control policy satisfies $|x_c| = 1$, i.e., the number of critical variable of all sampled scenarios is 1.

Assumption 3: There exists an optimal control policy such that $\mathbb{P}(A|x_c) = \mathbb{P}(A|x)$, which means that the critical variable x_c can totally dominate the crash probability.

Assumption 4: There exists an optimal importance function among q_j . Without loss of generality, let q_1 be the optimal importance function, i.e., $q_1(x_c) \triangleq \mathbb{P}(A|x_c)p(x_c)/\mu$.

Theorem 3: Under Assumptions 1, 2, 3 and 4, if β^* is any minimizer of $\sigma_{q_\alpha,\beta}^2$, then $\sigma_{q_\alpha,\beta^*}^2 = 0$.

Proof: From Assumptions 1 and 2, we know that all sampled scenarios will only be controlled once, i.e., $\mathcal{X} = \mathcal{X}_1$ and $\mu = \mu_1$, then

$$Z_1 = \frac{r_1(X)}{q_\alpha(X)} + \frac{\mu_1}{\alpha_1} \mathbb{I}_{\mathcal{X}_1}(X) = \frac{r_1(X)}{q_\alpha(X)} + \frac{\mu_1}{\alpha_1}, \quad (42)$$

and $\mathbb{E}_{q_\alpha}[Z_1] = \mu_1 \alpha_{1,1} / \alpha_1 = \mu_1 / \alpha_1$. Therefore, the asymptotic variance $\sigma_{1,q_\alpha,\beta_1}^2$ is

$$\begin{aligned} \sigma_{1,q_\alpha,\beta_1}^2 &= \mathbb{E}_{q_\alpha} \left[\left(Z_1 - \mathbb{E}_{q_\alpha}[Z_1] \right)^2 \right] \\ &= \sum_{x \in \mathcal{X}} \frac{r_1^2(x)}{q_\alpha(x)} \leq \frac{\sigma_{1,q_1}^2}{\alpha_1}. \end{aligned} \quad (43)$$

By Assumptions 3 and 4, we have $\mathbb{P}(A|x_c) = \mathbb{P}(A|x)$ and $q_1(x_c) = \mathbb{P}(A|x_c)p(x_c)/\mu$, then

$$\begin{aligned} \sigma_{1,q_1}^2 &= \sum_{x \in \mathcal{X}_1} \left(\frac{\mathbb{P}(A|x)p(x)}{q_1(x)} - \mu_1 \right)^2 q_1(x) \\ &= \sum_{x \in \mathcal{X}} \left(\frac{\mathbb{P}(A|x_c)p(x_c)}{q_1(x_c)} - \mu \right)^2 q_1(x) = 0. \end{aligned} \quad (44)$$

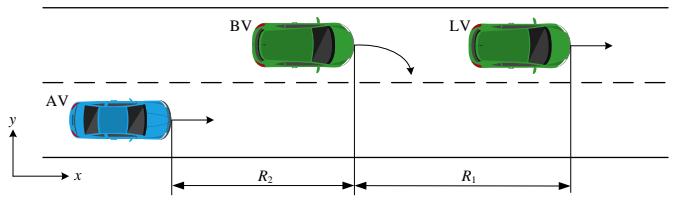


Fig. 3. Illustration of the overtaking scenarios.

Therefore, we conclude that $\sigma_{q_\alpha,\beta^*}^2 = \sigma_{1,q_\alpha,\beta_1}^2 = 0$. \square

Remark 3: Assumption 1 suggests that the scenarios in \mathcal{X}_0 should not be sampled. Since there are no crash in these scenarios, they can not make any contribution to the estimation. Assumption 2 requires that the number of critical variable is 1, because stratifying scenarios into different strata leads to some residual terms (e.g., $3(\mu_1/\alpha_j)^2$ in Eq. (41)) in estimation variance that can not be eliminated. Assumption 3 indicates that the critical variables should dominate the crash probability, since otherwise we may lose some critical information about the scenarios and obtain the suboptimal testing results. Assumption 4 requires that one of the importance functions should be optimal, together with Assumption 3 further reducing the asymptotic variances to zero. Although in practice these assumptions may not be fully satisfied, they could provide useful guidance for us to implement the SCV method.

Remark 4: The theorems in this section hold regardless of the specifics of SMs, which may be constructed by traditional traffic models or by neural networks.

VI. OVERTAKING CASE STUDY

A. Overtaking Scenarios

The overtaking scenarios are shown in Fig. 3, where the leading vehicle (LV) runs at the left lane, the background vehicle (BV) follows LV and the automated vehicle (AV) runs at the right lane. If BV cuts in to the right lane, then AV will follow BV and may rear-end BV, resulting in a crash. The state of the overtaking scenarios can be formulated as

$$s \triangleq (v_{BV}, R_1, \dot{R}_1, R_2, \dot{R}_2), \quad (45)$$

where $R_1 \triangleq x_{LV} - x_{BV}$, $\dot{R}_1 \triangleq v_{LV} - v_{BV}$, $R_2 \triangleq x_{BV} - x_{AV}$, and $\dot{R}_2 \triangleq v_{BV} - v_{AV}$. The x_{BV} , x_{LV} , x_{AV} are the positions and v_{BV} , v_{LV} , v_{AV} are the velocities of BV, LV and AV, respectively. The action of the overtaking scenario is defined as the actions of LV and BV, i.e., $a \triangleq (a_{LV}, a_{BV})$. We note that the overtaking scenarios are more stochastic and complicated than simple scenarios such as cut-in scenarios and car-following scenarios, since the BV in overtaking scenarios may have many chances to cut in, resulting in different cut-in scenarios and car-following scenarios between BV and AV. This is the reason why overtaking scenarios are always much more high-dimensional than cut-in scenarios.

B. Generation of NDE

The essence of NDE is to provide a driving environment where all BVs travel like humans. To generate NDE, the

probability distributions of the behaviors of all BVs should be consistent with the naturalistic driving data (NDD) [29]. In this paper, the probability distributions of free-driving, car-following, and cut-in behaviors are extracted from the NDD of the Safety Pilot Model Deployment (SPMD) [34] program and Integrated Vehicle-Based Safety System (IVBSS) [35] at the University of Michigan, Ann Arbor. The initial state is set as

$$s_0 = [v_{BV,0}, R_{1,0}, \dot{R}_{1,0}, R_{2,0}, \dot{R}_{2,0}], \quad (46)$$

where $v_{BV,0}$, $R_{1,0}$, $\dot{R}_{1,0}$ are sampled from the naturalistic distributions of car-following scenarios, $R_{2,0} \sim \mathcal{U}(20 \text{ m}, 100 \text{ m})$, $\dot{R}_{2,0} \sim \mathcal{U}(-5 \text{ m/s}, -10 \text{ m/s})$, where \mathcal{U} is the uniform distribution. After sampling the initial state, all vehicles select actions independently and simultaneously for each time step (0.1 s). The cut-in maneuver of BV is set completed within one time step. The car-following maneuver of AV is controlled by the intelligent driver model (IDM) [36]. The simulation continues until AV rear-ends BV or maximum simulation time (20 s) reached. Typically, the dimension of overtaking scenarios will exceed 1400 (201 time steps, each with 5 state variables and 2 action variables), leading to the high-dimensionality challenge.

C. Generation of NADE

The goal of NADE is to generate high-dimensional testing scenarios where the behaviors of BVs are adjusted only at critical moments, while keeping naturalistic distributions as in NDE at other time steps [21]. To construct the importance function, the maneuver criticality of BV is evaluated at each time step, which is defined as the multiplication of the exposure frequency and the maneuver challenge. The exposure frequency represents the probability of each action given current state in NDE. The maneuver challenge measures the probability of crash between AV and BV given current state and action. Since the AV models are usually black-boxes, the surrogate models (SMs) are adopted to approximate the maneuver challenge. In this paper, we use IDM and full velocity difference model (FVDM) [36] as SMs with different parameters: (1) IDM, denoted as SM-I; (2) FVDM with $a_{\min} = -1 \text{ m/s}^2$, denoted as SM-II; (3) FVDM with $a_{\min} = -6 \text{ m/s}^2$, denoted as SM-III. Then the importance functions can be obtained from the maneuver criticalities estimated by these SMs. Readers can find more technical details in [21].

D. Application of SCV

As shown in Algorithm 1, the SCV method can be applied to adjust the testing results and reduce estimation variance after testing AV in NADE. The key is to use importance functions of only sparse and critical variables to construct SCV, and then apply MLR of weighted testing results on SCV in each stratum. Finally, the estimated performance index is given by the summation of weighted intercepts obtained from MLR in all strata.

Algorithm 1: Adaptive safety evaluation with sparse control variates by multiple linear regression

```

Input:  $p$ ,  $q_\alpha$ ,  $X_{c,i}$ , and  $\mathbb{P}(A|X_i)$ ,  $i = 1, \dots, n$ 
Output:  $\tilde{\mu}_{q_\alpha, \hat{\beta}}$ ,  $\text{Var}_{q_\alpha}(\tilde{\mu}_{q_\alpha, \hat{\beta}})$ 
1 initialize  $Y_l$  and  $H_l$  as empty arrays,  $l = 0, \dots, L$ ;
2 initialize  $n_l = 0$ ,  $l = 0, \dots, L$ ;
3 for  $i \leftarrow 1$  to  $n$  do
4    $l \leftarrow$  number of control steps of  $X_{c,i}$ ;
5    $n_l \leftarrow n_l + 1$ ;
6   if  $l = 0$  then
7     append  $Y_l$  with  $\mathbb{P}(A|X_i)$ ;
8     append  $H_l$  with 0;
9   else
10    append  $Y_l$  with  $\mathbb{P}(A|X_i)p(X_{c,i})/q_\alpha(X_{c,i})$ ;
11    append  $H_l$  with  $\text{vec}(q_{j_1, \dots, j_l}(X_{c,i})/q_\alpha(X_{c,i}))$ ,
12       $j_1, \dots, j_l = 1, \dots, J - 1$ ;
13  end
14 end
15 for  $l \leftarrow 0$  to  $L$  do
16    $H_l \leftarrow H_l - \text{average}(H_l)$ ;
17   MLR  $\leftarrow$  multiple linear regression of  $Y_l$  on  $H_l$ ;
18    $\hat{\beta}_l \leftarrow$  estimated coefficients from MLR;
19    $\hat{\eta}_l \leftarrow$  estimated intercept from MLR;
20    $\tilde{\mu}_{l, q_\alpha, \hat{\beta}_l} \leftarrow n_l \hat{\eta}_l / n$ ,  $Z_l \leftarrow Y_l - H_l \hat{\beta}_l$ ;
21 end
22  $Z \leftarrow [Z_0, \dots, Z_L]$ ;
23  $\tilde{\mu}_{q_\alpha, \hat{\beta}} \leftarrow \sum_{l=0}^L \tilde{\mu}_{l, q_\alpha, \hat{\beta}_l}$ ,  $\text{Var}_{q_\alpha}(\tilde{\mu}_{q_\alpha, \hat{\beta}}) \leftarrow \text{var}(Z)$ ;
24 return  $\tilde{\mu}_{q_\alpha, \hat{\beta}}$ ,  $\text{Var}_{q_\alpha}(\tilde{\mu}_{q_\alpha, \hat{\beta}})$ ;

```

E. Evaluation Results

We validate the accuracy and efficiency of AV evaluation in NDE and NADE by the simulation of overtaking scenarios. The simulation is parallel conducted using 100 threads on a computer equipped with AMD® EPYC™ 7742 CPU and 512 GB RAM. Fig. 4 shows the crash rates of AV in NDE and NADE, respectively. The crash rate in NDE is presented as the black line in Fig. 4, with the bottom x -axis as its number of tests. The blue line in Fig. 4 represents the crash rate in NADE, and the top x -axis is the number of tests. The light shadow gives the 90% confidence interval. It can be seen that the crash rates in NDE and NADE converge to the same value, while NADE requires a much smaller number of tests. To measure the estimation precision of the crash rate, the relative half-width (RHW) [19] is adopted as the metric. The threshold of RHW is set to 0.3. To reach this threshold, NADE requires 6.76×10^6 number of tests, while NDE requires 1.21×10^8 number of tests, as shown in Fig. 5. It can be found that NADE can accelerate the evaluation by about 17.90 times compared with NDE. We note that the acceleration ratio is smaller than that in [21], because combinations of multiple various SMs are applied in this paper, which improves the robustness yet decreases the efficiency. The goal of the adaptive evaluation is to improve the efficiency while keeping the robustness.

To investigate the performance of the SCV method, the accuracy and efficiency of AV evaluation in NADE with and

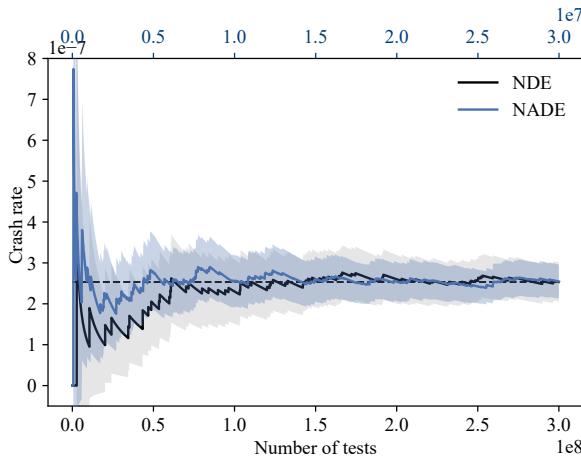


Fig. 4. Crash rates of AV in NDE and NADE, where the dashed line is the crash rate estimated by NDE.

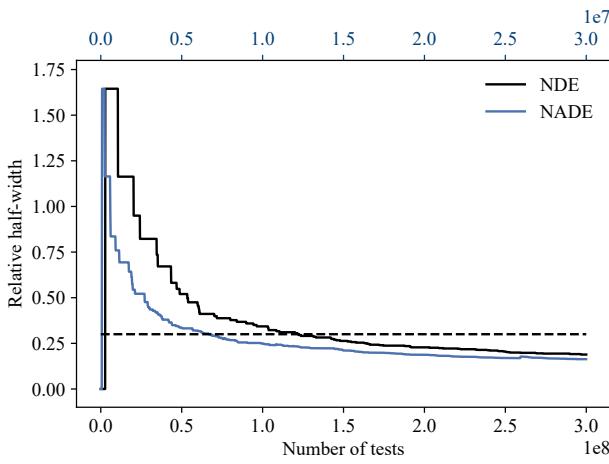


Fig. 5. RHW of AV evaluation in NDE and NADE, where the dashed line represents the RHW threshold (0.3).

without SCV are compared. It can be seen in Fig. 6 (a)-(e) that the crash rates of NADE and SCV converge to the same value for different number of tests. Fig. 6 (f) shows that the required numbers of tests of NADE and SCV for reaching the RHW threshold are 6.76×10^6 and 5.92×10^5 , respectively, resulting in a further acceleration ratio of 11.42. The weighted testing results before and after being adjusted by SCV with different number of control steps are compared in Fig. 7 (a)-(i), and Fig. 7 (j) shows the total 10^7 adjusted testing results. It can be seen that the SCV method is able to adjust the testing results into a much narrower interval, especially for relatively large number of control steps (e.g., $l \geq 4$), resulting in a considerable reduction of the estimation variance.

The detailed regression processes of the SCV method are also investigated. Fig. 8 shows the number of tests, the number of SCV and the maximum rank of the control matrices for the number of control steps $l = 1, \dots, 9$, respectively. Note that for $l \geq 10$, we only use the first 9 control steps to construct the SCV. It can be seen that the maximum number of tests appears at $l = 6$ and then the number of tests decreases to a relatively low level. As shown in Eq. (27), the maximum rank

TABLE I
AARs of SCV WHERE AV ADMITS IDMS WITH DIFFERENT α VALUES,
AND THE RIGHTMOST COLUMN CORRESPONDS TO THE VT-IDM.

α	0.5	1.0	1.5	2.0	2.5	3.0	VT-IDM
AAR	11.52	9.02	7.87	6.76	7.73	10.90	7.30
α	3.5	4.0	4.5	5.0	5.5	6.0	
AAR	13.44	11.95	11.12	10.61	10.45	10.05	

of the control matrices is the minimum value between the number of tests and the number of SCVs, and hence will not grow exponentially with the number of control steps, although the number of SCVs will do. Therefore, the SVD of control matrices is always tractable in each stratum and the optimal control parameters can be found to minimize the estimation variance.

Since the scenario generation processes are stochastic, the testing and evaluation results are usually not the same in different experiments. Therefore, to find the average performances, we shuffle the testing results 200 times to bootstrap them and obtain the frequency distributions of the required number of tests (RNoT) in NDE and NADE. The average RNoT of NDE and NADE are 1.20×10^8 and 8.71×10^6 , respectively. Therefore, the average acceleration ratio (AAR) of NADE with respect to NDE is 13.78. The testing results of SCV are also bootstrapped by 200 times. For cases with maximum RHW below 0.3, we use the RNoT when the maximum RHW is reached. The average RNoT of SCV is 1.29×10^6 , resulting in an AAR of 6.76 times compared with NADE.

F. Generalizability Analysis

In the above experiments, we have set the AV model the same as SM-I, i.e., they are both IDMs with same parameters. To investigate the generalizability of the SCV method for different AV models, the IDMs with a series of parameters $\alpha = 0.5, 1.0, \dots, 6.0$ are chosen as AV models. The AARs of SCV compared with NADE are shown in Table I. The testing results of all AV models are shuffled 200 times to obtain the AARs. It can be seen that the minimum AAR appears at $\alpha = 2.0$, where the AV model is the same as SM-I, while the maximum AAR appears at $\alpha = 3.5$. The mean AAR for different AV models is 10.12. Therefore, the SCV method can further accelerate the evaluation process by about one order of magnitude for various types of AV models. Moreover, the AARs of SCV with AV models different from SM-I are always greater than that of AV model the same as SM-I. The reason is that although using AV models different from SM-I will do harm to both the estimation efficiency of NADE and SCV, the damage to NADE is more than to SCV.

In addition, we also select the calibrated IDM in [37] (denoted as VT-IDM) as the AV model to further validate the generalization performance of the SCV method. The testing results shuffled 200 times give an AAR of 7.30 for SCV compared with NADE, which is shown at the rightmost column in Table I. Therefore, the SCV method can also increase the evaluation efficiency considerably for AV model with completely different calibrated parameters. This is not a

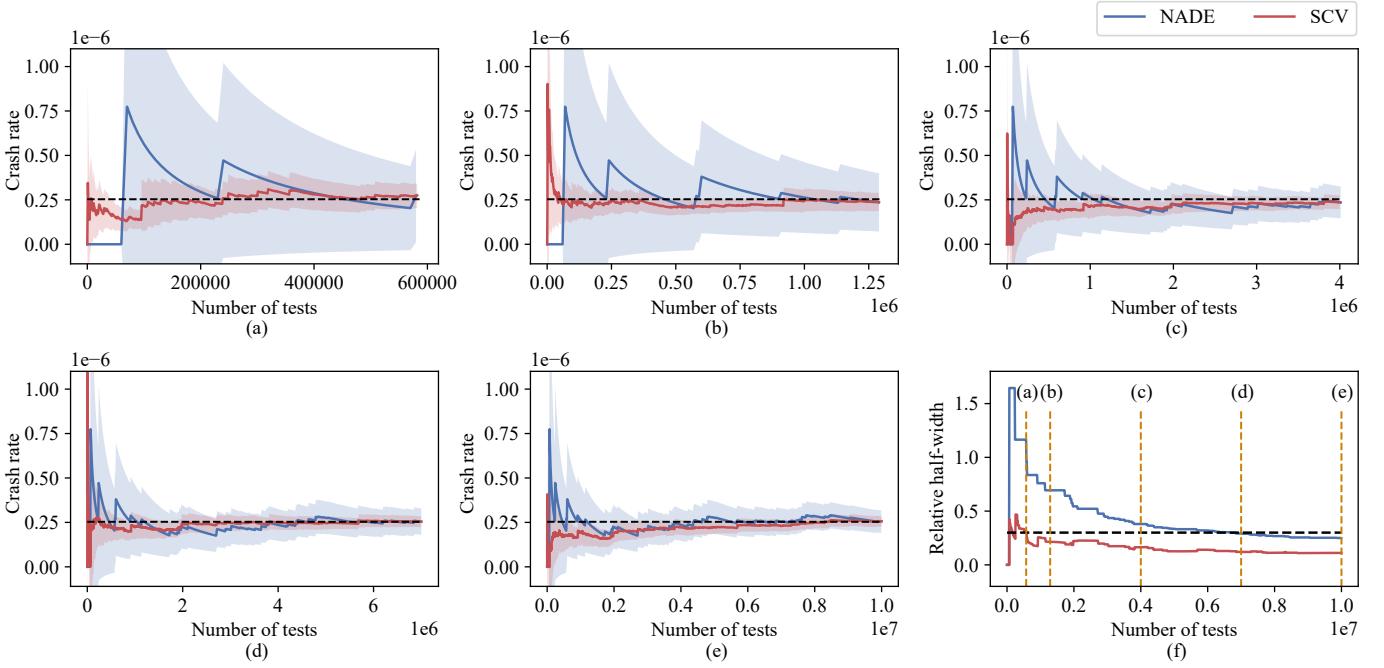


Fig. 6. Crash rate of AV using NADE and SCV for (a) $n = 5.92 \times 10^5$, (b) $n = 1.29 \times 10^6$, (c) $n = 4 \times 10^6$, (d) $n = 7 \times 10^6$ and (e) $n = 1 \times 10^7$, where n is the total number of tests and the dashed line is the crash rate estimated by NDE; (f) RHW of AV evaluation using NADE and SCV, where the dashed line in black represents the RHW threshold (0.3) and 5 dashed lines in orange correspond to (a)-(e).

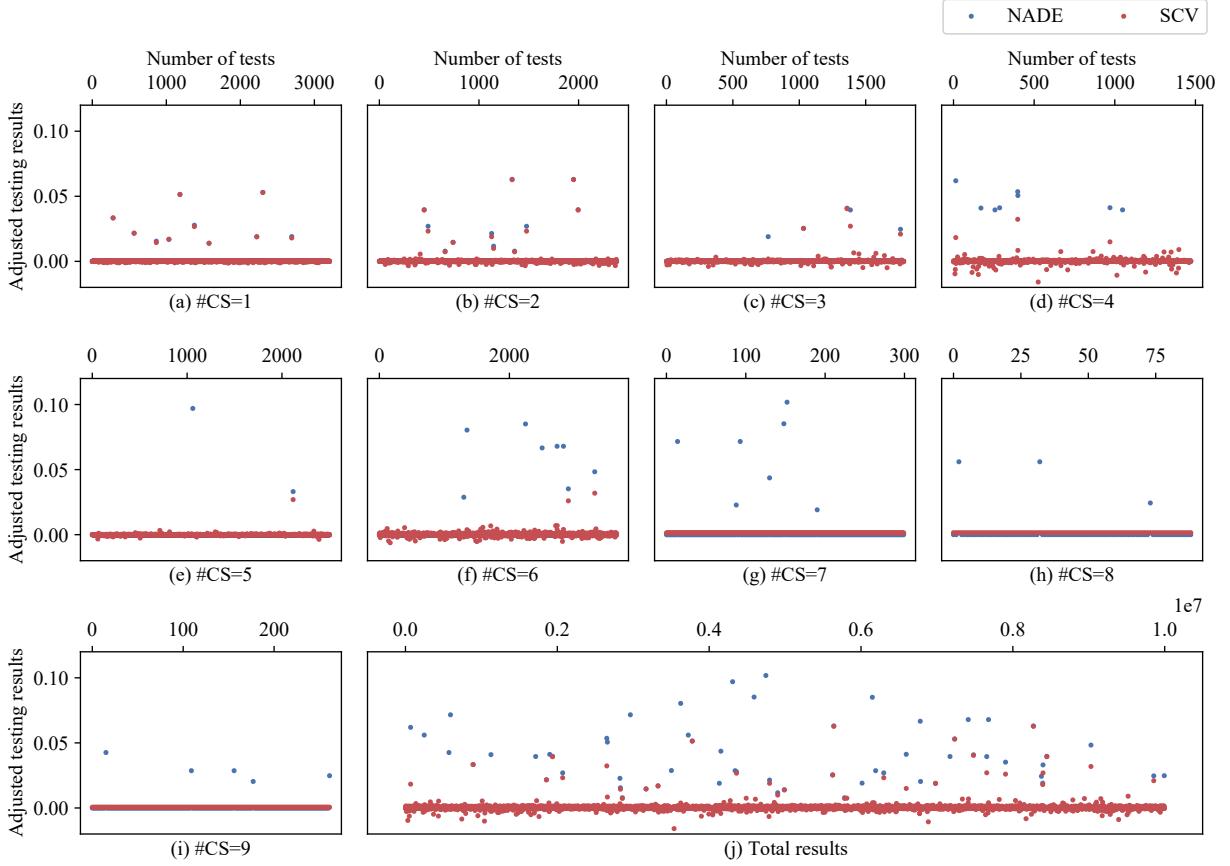


Fig. 7. Adjusted testing results by NADE and SCV for (a)-(i) the number of control steps (#CS) from 1 to 9 and (j) total 10^7 testing results.

surprising result because the only requirement for the SCV method to work is that the SMs and the AV model have

some correlation, and more correlation contributes to more variance reduction. Although the VT-IDM and IDM have

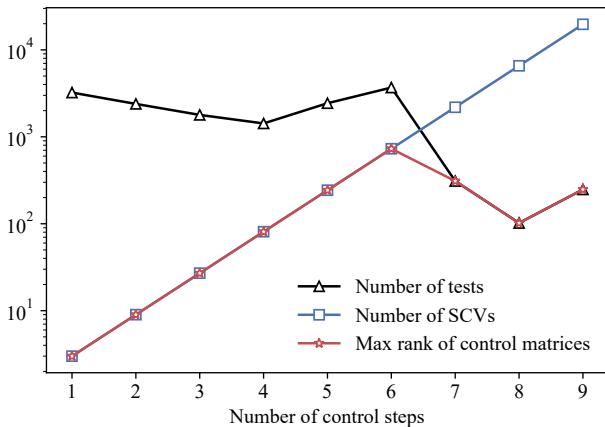


Fig. 8. Number of tests, number of SCV and maximum rank of control matrices for different number of control steps.

totally different parameters, they are still correlated to some extent.

VII. CONCLUSION

In this paper, we propose an adaptive safety evaluation framework for CAVs in high-dimensional scenarios with a newly developed sparse control variates (SCV) method. To address the CoD, the SCV are constructed by only considering the sparse and critical variables of testing scenarios and stratified into strata accordingly. By optimizing the SCV leveraging the testing results within each stratum, the estimation variance is significantly reduced for different CAVs adaptively, accelerating the evaluation process. The accuracy, efficiency and optimality of the proposed method are verified and validated by both theoretical analysis and empirical studies. Comparing with the evaluation efficiency in NDE and NADE, our method is always more efficient particularly for CAVs that are different from SMs. It has been noted that adaptive testing scenario generation and adaptive testing result evaluation are two complementary approaches for adaptive testing and evaluation of CAVs. How to develop the former in high-dimensional scenarios deserves further investigation.

REFERENCES

- [1] N. Kalra and S. M. Paddock, “Driving to safety: How many miles of driving would it take to demonstrate autonomous vehicle reliability?” *Transportation Research Part A: Policy and Practice*, vol. 94, pp. 182–193, 2016.
- [2] A. Li, S. Chen, L. Sun, N. Zheng, M. Tomizuka, and W. Zhan, “Scgene: Bio-inspired traffic scenario generation for autonomous driving testing,” *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [3] J. Wang, A. Pun, J. Tu, S. Manivasagam, A. Sadat, S. Casas, M. Ren, and R. Urtasun, “Advsim: Generating safety-critical scenarios for self-driving vehicles,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9909–9918.
- [4] T. Menzel, G. Bagschik, and M. Maurer, “Scenarios for development, test and validation of automated vehicles,” in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1821–1827.
- [5] Y. Tian, K. Pei, S. Jana, and B. Ray, “Deepertest: Automated testing of deep-neural-network-driven autonomous cars,” in *Proceedings of the 40th International Conference on Software Engineering*, 2018, pp. 303–314.
- [6] D. Rempe, J. Phlion, L. J. Guibas, S. Fidler, and O. Litany, “Generating useful accident-prone driving scenarios via a learned traffic prior,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17305–17315.
- [7] L. Li, W.-L. Huang, Y. Liu, N.-N. Zheng, and F.-Y. Wang, “Intelligence testing for autonomous vehicles: A new approach,” *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 2, pp. 158–166, 2016.
- [8] L. Li, Y.-L. Lin, N.-N. Zheng, F.-Y. Wang, Y. Liu, D. Cao, K. Wang, and W.-L. Huang, “Artificial intelligence test: A case study of intelligent vehicles,” *Artificial Intelligence Review*, vol. 50, no. 3, pp. 441–465, 2018.
- [9] L. Li, X. Wang, K. Wang, Y. Lin, J. Xin, L. Chen, L. Xu, B. Tian, Y. Ai, J. Wang et al., “Parallel testing of vehicle intelligence via virtual-real interaction,” *Science Robotics*, 2019.
- [10] S. Riedmaier, T. Ponn, D. Ludwig, B. Schick, and F. Diermeyer, “Survey on scenario-based safety assessment of automated vehicles,” *IEEE access*, vol. 8, pp. 87456–87477, 2020.
- [11] G. E. Mullins, P. G. Stankiewicz, R. C. Hawthorne, and S. K. Gupta, “Adaptive generation of challenging scenarios for testing and evaluation of autonomous vehicles,” *Journal of Systems and Software*, vol. 137, pp. 197–215, 2018.
- [12] M. Koren, S. Alsaif, R. Lee, and M. J. Kochenderfer, “Adaptive stress testing for autonomous vehicles,” in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1–7.
- [13] S. Feng, Y. Feng, H. Sun, Y. Zhang, and H. X. Liu, “Testing scenario library generation for connected and automated vehicles: an adaptive framework,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 2, pp. 1213–1222, 2022.
- [14] J. Sun, H. Zhou, H. Xi, H. Zhang, and Y. Tian, “Adaptive design of experiments for safety evaluation of automated vehicles,” *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [15] H. X. Liu and S. Feng, “‘curse of rarity’ for autonomous vehicles,” *arXiv preprint arXiv:2207.02749*, 2022.
- [16] S. Feng, Y. Feng, X. Yan, S. Shen, S. Xu, and H. X. Liu, “Safety assessment of highly automated driving systems in test tracks: A new framework,” *Accident Analysis & Prevention*, vol. 144, p. 105664, 2020.
- [17] S. Feng, Y. Feng, C. Yu, Y. Zhang, and H. X. Liu, “Testing scenario library generation for connected and automated vehicles, part i: Methodology,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 3, pp. 1573–1582, 2021.
- [18] S. Feng, Y. Feng, H. Sun, S. Bao, Y. Zhang, and H. X. Liu, “Testing scenario library generation for connected and automated vehicles, part ii: Case studies,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 9, pp. 5635–5647, 2021.
- [19] D. Zhao, H. Lam, H. Peng, S. Bao, D. J. LeBlanc, K. Nobukawa, and C. S. Pan, “Accelerated evaluation of automated vehicles safety in lane-change scenarios based on importance sampling techniques,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 3, pp. 595–607, 2016.
- [20] D. Zhao, X. Huang, H. Peng, H. Lam, and D. J. LeBlanc, “Accelerated evaluation of automated vehicles in car-following maneuvers,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 3, pp. 733–744, 2017.
- [21] S. Feng, X. Yan, H. Sun, Y. Feng, and H. X. Liu, “Intelligent driving intelligence test for autonomous vehicles with naturalistic and adversarial environment,” *Nature Communications*, vol. 12, no. 1, pp. 1–14, 2021.
- [22] R. Y. Rubinstein and R. Marcus, “Efficiency of multivariate control variates in monte carlo simulation,” *Operations Research*, vol. 33, no. 3, pp. 661–677, 1985.
- [23] W. Grathwohl, D. Choi, Y. Wu, G. Roeder, and D. Duvenaud, “Back-propagation through the void: Optimizing control variates for black-box gradient estimation,” in *International Conference on Learning Representations*, 2018.
- [24] C.-A. Cheng, X. Yan, and B. Boots, “Trajectory-wise control variates for variance reduction in policy gradient methods,” in *Conference on Robot Learning*. PMLR, 2020, pp. 1379–1394.
- [25] A. Shapiro, “Monte carlo sampling methods,” *Handbooks in operations research and management science*, vol. 10, pp. 353–425, 2003.
- [26] A. B. Owen, *Monte Carlo theory, methods and examples*. Stanford, 2013.
- [27] D. J. Olive, “Multiple linear regression,” in *Linear Regression*. Springer, 2017, pp. 17–83.
- [28] J. Yang, H. He, Y. Zhang, S. Feng, and H. X. Liu, “Adaptive testing for connected and automated vehicles with sparse control variates in overtaking scenarios,” in *IEEE International Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2022.
- [29] X. Yan, S. Feng, H. Sun, and H. X. Liu, “Distributionally consistent simulation of naturalistic driving environment for autonomous vehicle testing,” *arXiv preprint arXiv:2101.02828*, 2021.

- [30] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” *Advances in neural information processing systems*, vol. 30, 2017.
- [31] S.-K. Au and J. Beck, “Important sampling in high dimensions,” *Structural safety*, vol. 25, no. 2, pp. 139–163, 2003.
- [32] A. Owen and Y. Zhou, “Safe and effective importance sampling,” *Journal of the American Statistical Association*, vol. 95, no. 449, pp. 135–143, 2000.
- [33] M. E. Wall, A. Rechtsteiner, and L. M. Rocha, “Singular value decomposition and principal component analysis,” in *A practical approach to microarray data analysis*. Springer, 2003, pp. 91–109.
- [34] D. Bezzina and J. Sayer, “Safety pilot model deployment: Test conductor team report,” *Report No. DOT HS*, vol. 812, no. 171, p. 18, 2014.
- [35] J. Sayer, D. LeBlanc, S. Bogard, D. Funkhouser, S. Bao, M. L. Buonarosa, A. Blankspoor *et al.*, “Integrated vehicle-based safety systems field operational test: Final program report,” United States. Joint Program Office for Intelligent Transportation Systems, Tech. Rep., 2011.
- [36] J. W. Ro, P. S. Roop, A. Malik, and P. Ranjitkar, “A formal approach for modeling and simulation of human car-following behavior,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 2, pp. 639–648, 2017.
- [37] J. Sangster, H. Rakha, and J. Du, “Application of naturalistic driving data to modeling of driver car-following behavior,” *Transportation research record*, vol. 2390, no. 1, pp. 20–33, 2013.