

An Interaction-aware Evaluation Method for Highly Automated Vehicles

Xinpeng Wang¹, Songan Zhang¹, Kuan-Hui Lee², Hwei Peng¹

Abstract—It is important to build a rigorous verification and validation (V&V) process to evaluate the safety of highly automated vehicles (HAVs) before their wide deployment on public roads. In this paper, we propose an interaction-aware framework for HAV safety evaluation which is suitable for some highly-interactive driving scenarios including highway merging, roundabout entering, etc. Contrary to existing approaches where the primary other vehicle (POV) takes predetermined maneuvers, we model the POV as a game-theoretic agent. To capture a wide variety of interactions between the POV and the vehicle under test (VUT), we characterize the interactive behavior using level- k game theory and social value orientation and train a diverse set of POVs using reinforcement learning. Moreover, we propose an adaptive test case sampling scheme based on the Gaussian process regression technique to generate customized and diverse challenging cases. The highway merging is used as the example scenario. We found the proposed method is able to capture a wide range of POV behaviors and achieve better coverage of the failure modes of the VUT compared with other evaluation approaches.

I. INTRODUCTION

Highly automated vehicles (HAVs) are under rapid development all over the world. They have the potential to transform ground transportation by liberating people from tedious driving tasks and improving road safety by avoiding human errors. It's crucial to conduct verification and validation on their safety efficacy before their wide deployments.

Safety evaluation of HAVs has been conducted at multiple traffic scenarios including unprotected left-turn, cut-in, and pedestrian crossing scenarios [1]–[3], etc. They can be characterized as reactive tests, where the vehicle under test (VUT) will be challenged by the primary other vehicle (POV) "in a surprise". The test case is fully defined by the initial condition of the challenge. Due to the short duration of the challenge, the POV is typically not programmed to interact with the VUT, and a predetermined trajectory is assumed.

For SAE level 3 and above automated vehicles [4], their operational design domain (ODD) can include dynamic and complex scenarios, including highway merging, roundabout entering, turning at unsignalized intersections, etc. For these scenarios, the existing methods show their limitations. For example, in highway merging, the VUT attempts to merge from the ramp onto the main road when another vehicle is present, which serves as the POV. The merging ramp gives enough time for the two vehicles to interact, and thus the assumption of "no interaction between POV and

VUT" becomes unrealistic. Moreover, in the role of the POV, different human drivers may exhibit different behaviors under the same initial conditions, including coasting, accelerating to pull ahead, decelerating to yield, etc. These diverse behaviors pose a novel challenge for the motion prediction and decision-making module of the VUT, and should be incorporated into the evaluation framework.

We propose an interaction-aware evaluation methodology in this paper. It consists of two parts: first, we create a test case pool, in which we model a set of interactive POVs using level- k game theory and social value orientation (SVO). Second, we propose an adaptive sampling scheme based on Gaussian process regression to generate challenging test cases for a given VUT. In this paper, we focus on the highway merging scenario, but this methodology can be applied to other scenarios including roundabout entering, turning at unsignalized intersections, etc. This is the first effort to comprehensively identify the failure modes of a VUT in an interactive scenario.

The paper is organized as follows: Section 2 introduces related works; Sections 3 to 5 introduce the proposed method by first formulating the evaluation problem, then introduces the POV library construction and finally the adaptive test case generation procedure. Section 6 discusses the implementation details for the highway merging scenario; Section 7 shows the simulated testing results and the comparison with other sampling methods; finally concluding remarks are made in Section 8.

II. RELATED WORK

Evaluation of the safety of HAVs has been an active area in recent years. Many test procedures have been proposed. Test matrix has been used to evaluate advanced driver assistance system (ADAS) [5]. However, the VUT can be tuned to pass the predefined test cases, but may fail under broader conditions in real-world driving tasks. Worst-case evaluation methods attempt to generate adversarial situations or POV inputs to create edge cases. [6] and [7] used reachability analysis to find test cases where the VUT has a minimal solution space. [8] applied simulation-based falsification to find failure cases for a given VUT. However, the assumption of adversarial POVs may not be reasonable. [9] applied reinforcement learning to create adversarial yet socially acceptable POV behaviors in a highway driving scenario, while the diversity of the challenging scenarios has not been discussed. On the other hand, Monte-Carlo sampling-based evaluation methods have been proposed to generate test cases to estimate the real-world performance of the VUT. Some research uses importance sampling [1] [2] or subset

¹Xinpeng Wang (xinpengw@umich.edu), Songan Zhang, Hwei Peng are with the Department of Mechanical Engineering, the University of Michigan, Ann Arbor, MI 48109, U.S.

²Kuan-Hui Lee is with Toyota Research Institute, Los Altos, CA 94022, U.S.

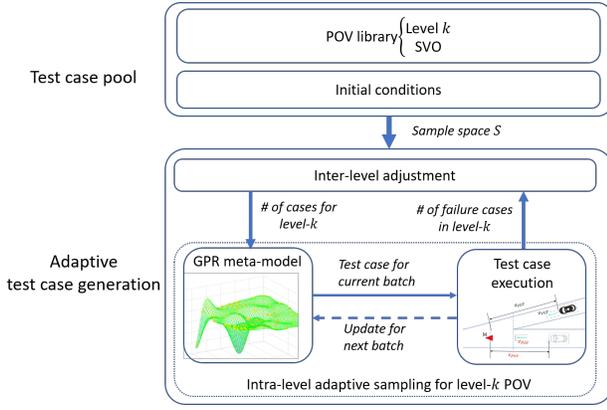


Fig. 1. Pipeline of the interaction-aware evaluation method.

simulation [10] to efficiently estimate the crash rate of the VUT, while other works customize test cases to identify the failure modes of the VUT using adaptive sampling method [11]–[13]. The interactions between POV and VUT have not been considered in these works.

To model the interactive nature of human driving behavior, game theory has been widely applied, in which humans are modeled as utility-maximizing rational agents. Nash [14] or Stackelberg [15], [16] equilibrium models have been applied to model human driving behaviors. They rely on the assumption that each agent has an infinite level of rationality, which could be too strict considering that human drivers have to make quick decisions in a complex and dynamic environment. Therefore, other researchers assumed bounded rationality of human drivers and applied level- k game theory [17], quantal response [18] or cumulative prospect theory [19] to model human driving behaviors. On the other hand, [14] and [20] considered the altruism of human driving behaviors in a game-theoretic setting. Despite the richness of game-theoretic models, they have yet to be comprehensively considered for HAV evaluations. Filling this gap is the focus of this work.

III. PROBLEM FORMULATION

We aim to systematically generate test cases for a given VUT in interactive scenarios. The tasks are two-fold: firstly, we will create a test case pool for the target scenario; secondly, we proposed a mechanism to sample test cases from the test case pool. The test cases can be characterized by two sets of attributes: the first set defines the initial condition of the scenario; the second set describes the interactive and behavioral properties of the POV, which determines its driving policy. The test case sampling procedure aims to evaluate the safety performance of a black-box VUT by finding the failure modes of it through efficient sampling schemes. The overall concept of the proposed interaction-aware evaluation method is shown in Figure 1.

IV. POV LIBRARY CONSTRUCTION

The POV library needs to capture a diverse set of POV behaviors. On the one hand, the POV model should ap-

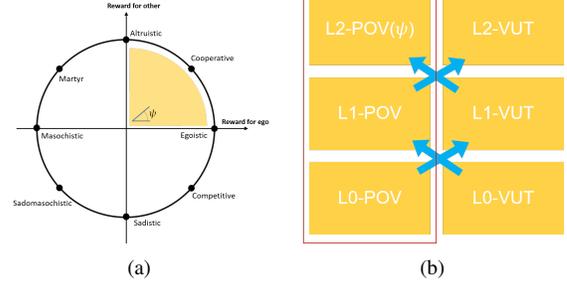


Fig. 2. (a) The SVO ring: we will focus on $\psi \in [0, \pi/2)$. (b) The hierarchy of level- k agents. The level-2 POV has an extra parameter ψ characterizing its SVO angle.

proximate the decision-making procedure of human drivers. Therefore, we assume that the POVs are game-theoretic agents, which take the (near) optimal action according to its utility function and assumptions of the opponents. On the other hand, the modeling framework should have the flexibility to describe a wide range of possible driving behaviors. We model the POVs as agents that possess different assumptions on the VUT and different utility functions for their own behavior. Specifically, we adopt the idea of level- k game theory and social value orientation to describe the diversified POVs.

A. Level- k game formulation

The level- k game theory model is based on the idea that intelligent agents (such as human drivers) have finite level of rationality. The model first assumes a known level-0 agent, which is a naive agent that behaves non-cooperatively. Then, a level- k agent ($k > 0$) will assume that all the opponents are level- $(k - 1)$ and will behave optimally according to this assumption. Using the level-0 policy as the starting point, the optimal policy for a level- k agent can be generated sequentially. According to an experimental study in economics [21], human decision-makers are usually as high as level-2 thinkers. Therefore, we only consider agents that are up to level-2 in this research to illustrate the concept.

B. Social value orientation

Social value orientation (SVO) is a concept from social psychology literature, which quantifies the agent’s degree of selfishness [22]. It can be represented as an orientation angle ψ indicating the agent’s preference on the outcome for itself versus for others, as shown in Figure 2(a), where different ψ represent personalities including egoistic, pro-social, altruistic, competitive, etc. In a common game-theoretic setting, an agent is egoistic and will solely optimize for its own utility function, i.e. $\psi = 0$. However, when the variable SVO is combined with a game-theoretic driver model, as shown in [14], it can significantly improve the accuracy of trajectory prediction for human drivers, thus better explain human driving behaviors. Moreover, agents with different SVO could represent a continuous spectrum of human drivers, which complements the level- k framework

where humans only have discrete types. In this work, we combine the SVO with the level- k game theory to capture richer POV behaviors.

C. POV library construction using reinforcement learning

Based on the level- k game theory and SVO, we create a library with the following POV agents: level-0 POV, level-1 POV and level-2 POV with varying social value orientation. Due to the non-competitive nature of driving tasks, we only consider the SVO angle in the 1st quadrant for simplicity, i.e. $0 \leq \psi < \pi/2$. The reasons that we do not consider SVO for lower-level POVs are that: a level-0 POV is non-cooperative, thus SVO cannot be defined; a level-1 POV assumes its opponent is level-0 and non-cooperative, thus SVO is not defined either.

To construct the POV library, we first design the policy for a level-0 POV as a baseline. Next, to generate the policy for a level- k POV ($k > 0$), a level- $(k-1)$ VUT is needed in advance. Therefore, we start with a level-0 VUT policy, and then generate higher-level POV and VUT sequentially in a double-helix structure, as shown in Figure 2(b). Although the targets are level- k POVs, we still need to compute level- k VUTs as the stepping stones to obtain higher-level POVs. In simulated and real tests, the VUTs we evaluate are not these model VUTs.

For level-0 POV and VUT, they behave non-cooperatively with fixed speed profiles, which capture the behavior of inattentive drivers. For level- k POV and VUT ($k > 0$), we use reinforcement learning (RL) to compute their driving policies. To train a level- k POV, we model it as an agent operating in an environment of level- $(k-1)$ VUT. The same procedure applies to VUT. To incorporate the factor of SVO, we consider the SVO angle as an extra state of the model when the level-2 POV is trained to generate a continuum of level-2 POVs.

1) *Reinforcement learning basics*: Computing a rational agent can be modelled as an Markov Decision Process (MDP) problem, which is defined by $\mathcal{M} = (\mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma)$, with the state space $\mathcal{X} \subseteq \mathbb{R}^n$, the action space $\mathcal{U} \subseteq \mathbb{R}^m$, the transition dynamics of the environment $\mathcal{P} : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$, the reward function $\mathcal{R} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$, and the discount factor $\gamma \in [0, 1)$.

At each state x_i , an agent tries to compute a best action u_i from the state-action mapping, i.e. the policy $\pi(x_i) = u_i$, that maximizes the expected cumulative reward, written as $E_\pi[\sum_{t=0}^{t_1} \gamma^t r(t)]$, where t_1 is the end time. To learn the optimal policy π^* , we use the Q-learning technique. We first define the action-value function Q :

$$Q(x, u|\pi) = E_\pi \left[\sum_{t=0}^{t_1} \gamma^t r_t | x_0 = x, u_0 = u \right] \quad (1)$$

Then π^* is learned by training the agent to learn the optimal Q function, i.e. $Q^*(x, u|\pi^*)$, which satisfies the Bellman equation. For details please refer to [23].

2) *Reinforcement learning formulation*: For a level- k VUT, the state space includes all continuous physical states of the POV and the VUT, denoted as X ($\mathcal{X} = X$). For POVs, the SVO angle is considered in the states space, which is held constant in each episode, i.e. $\mathcal{X} = X \times [0, \pi/2)$. The action space is a discrete set of acceleration or steering input.

The reward function reflects the goal of driving for each agent. We assume that the reward function can be represented as:

$$r(x, u) = W^T \Phi(x, u) = \sum_{i=1}^k w_i \phi_i(x, u) \quad (2)$$

which is a linear combination of multiple terms, each represents a different attribute for driving. There are three categories:

- 1) Ego reward for POV: $r_{POVe} = W_{POVe}^T \Phi_{POVe}$.
- 2) Ego reward for VUT: $r_{VUTE} = W_{VUTE}^T \Phi_{VUTE}$.
- 3) Safety reward for both: $r_{safe} = W_{safe}^T \Phi_{safe}$.

The final reward function for VUT is:

$$r_{VUT} = r_{safe} + r_{VUTE} \quad (3)$$

For a POV with SVO angle ψ , the reward function is:

$$r_{POV} = r_{safe} + r_{POVe} \cos(\psi) + r_{VUTE} \sin(\psi) \quad (4)$$

where ψ modulates the rewards between POV and VUT. For a level-1 POV, $\psi \equiv 0$.

3) *Training POV & VUT agents using DDQN*: In this work, since the state space is continuous, we use an artificial neural network as the function approximator for the optimal action-value function Q^* . The reinforcement learning algorithm we use is Double Deep-Q network (DDQN) [24]. DDQN is based on the Deep-Q network (DQN) [25] method. It addresses the problem of overestimating future return of DQN by decoupling the action evaluation and action selection into max operations in two different Q-networks. For the MDP with discrete action space and low dimensional state space, other advanced RL methods are not necessarily better than the DDQN approach. For other applications, DDQN can be replaced by other appropriate RL method.

V. ADAPTIVE TEST CASE GENERATION

A. Problem formulation for adaptive testing

From the previous section, we systematically generate the interactive POV library, which is characterized by the SVO ψ and rationality level L . Combined with the initial condition of the scenario x_0 , we can build the test cases pool, denoted as \mathcal{S} , where each case is $s = [x_0^T, \psi, L]^T$. Then, we need a mechanism to pick a set of N test cases $s = [s_1, \dots, s_N]$ from the pool to identify the failure modes of the VUT. The main challenge is that different VUTs may have different performance profiles and weaknesses, and thus the failure modes are unknown at the beginning of the V&V process. Therefore, we need a sampling scheme that can select new cases based on past test results to adaptively search for the weaknesses of each VUT as the testing proceeds. The goals of the test case generation process are two-fold:

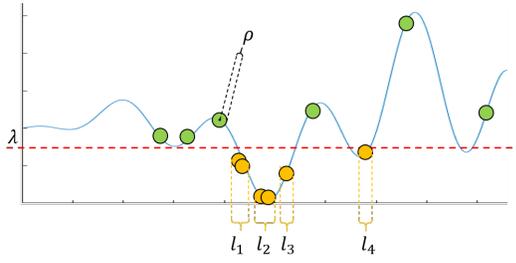


Fig. 3. Measuring the FMC of test samples: the 1-D illustration. The blue curve represents the performance score $P(s)$, the red dashed line shows the performance threshold λ , and the regions of curve below the threshold are the failure modes. The FMC is computed as $M(s, \rho, \lambda) = l_1 + l_2 + l_3 + l_4$.

- 1) Challenge: find cases where the VUT performs poorly (i.e. identify its weakness).
- 2) Coverage: identify (possibly disassociated) regions of weak performance.

For a test run with case s , the performance of a VUT can be evaluated by function P , which takes the VUT trajectory $\tau = [x(0), u(0), x(1), u(1) \dots x(t_1 - 1), u(t_1 - 1), x(t_1)]$ as input, and computes the performance score. It is written as:

$$P(s) = f(\tau) = \mu_1 I_{crash} + \mu_2 P_{safety} + \mu_3 P_{task} \quad (5)$$

where I_{crash} is the indicator function for collision; P_{safety} is the safety score; P_{task} is the score on task accomplishment (success highway merge, smooth acceleration, etc); μ_1, μ_2, μ_3 are weighting factors.

To describe the aforementioned two goals, we propose the criterion of failure mode coverage (FMC) for evaluating the quality of test samples:

$$M(s, \rho, \lambda) = \int_{\bigcup B(\rho, s_{I(\lambda)})} \mathbf{1} dv \quad (6)$$

where $B(\rho, s)$ is a hyper-ball centered around case s with radius ρ ; $s_{I(\lambda)}$ is a subset of s , such that $\forall s \in s_{I(\lambda)}, P(s) < \lambda$. The FMC evaluates the volume of the union of hyper-balls centered around test cases for which the VUT behaves poorly ($P(s) < \lambda$), which characterizes the coverage of failure modes for the VUT. Here, all the dimensions are normalized between [0,1]. Figure 3 is a graphic illustration of the FMC in 1-D.

B. Adaptive testing method overview

To meet the goals of adaptive testing, we apply an adaptive sampling method. We generate N test cases in batches sequentially, with batch size n . For each case in S , the last attribute L is a categorical variable, while all others are continuous variables, i.e. $S = S \times \{0, 1, 2\}$. Therefore, we separate the sampling scheme for each batch into two stages, as shown in the lower part of Figure 1. In the 1st stage, we allocate the number of samples into different POV levels, i.e. assign n_k^i cases to be tested with level- k POV at batch i . In the 2nd stage, we generate new test cases within each POV level from S using Gaussian process regression (GPR). We will elaborate on the two stages later in this section.

C. Intra-level adaptive sampling

Within each POV level, we conduct adaptive sampling using the Gaussian process regression (GPR). GPR is a non-parametric probabilistic model [26]. The key idea is to maintain and update a GPR based meta-model based on existing samples, and use the meta-model to guide the generation of a new batch of samples.

1) *Gaussian process regression*: Gaussian process (GP) is a stochastic process, for which the joint distribution of every finite collection of random variables follows a multivariate Gaussian distribution. A GP, as shown in (7), is characterized by its mean function $m(x)$ and a covariance function $k(x, x')$ (kernel).

$$f(x) \sim GP(m(x), k(x, x')) \quad (7)$$

In this work, we use GP to model the performance surface of each VUT, as shown in (8).

$P(s) = \epsilon + f(s)$, where

$$f(s) \sim GP(0, k(s, s'|\theta)), \epsilon \sim N(\beta, \sigma^2) \quad (8)$$

where (β, σ, θ) are the parameters of the model. In this work, we use a zero mean function and a square-exponential kernel function for the GPR model. Model parameters are optimized using maximum likelihood estimation. The procedure of adaptive sampling is illustrated in Algorithm 1, and some details are explained below.

Algorithm 1 Intra-level adaptive sampling

Input: batches number i ; batch size n_k^i ; previous GPR model \hat{P}_k^{i-1} ; exploration factor ϵ_0 .

Output: test cases with level- k POV s_k^i and test results y_k^i ; updated GPR model \hat{P}_k^i .

- 1: **if** $i = 1$ **then**
 - 2: Sample initial test batch s_k^1 uniformly from S .
 - 3: **else**
 - 4: Uniformly sample p queries \check{s} from S ($p \gg n_k^i$).
 - 5: $\epsilon = \epsilon_0 \alpha^{i-1}$.
 - 6: Pick $(1 - \epsilon)n_k^i$ queries from \check{s} according to $q_{exploit}(s)$ as $s_{exploit}$.
 - 7: Pick ϵn_k^i queries according to $q_{explore}(s)$ as $s_{explore}$.
 - 8: $s_k^i = [s_{exploit}, s_{explore}]$.
 - 9: **end if**
 - 10: Execute test cases s_k^i , acquire results y_k^i .
 - 11: Fit/Update the GPR model: $y = \hat{P}_k^i(s) = \hat{P}(s | s_k^{1:i}, y_k^{1:i})$.
-

2) *Balancing between exploration and exploitation*: To achieve good coverage of the failure modes, we need to balance exploration and exploitation. On the one hand, it is desirable to explore regions with high uncertainty and pick more informative samples for more accurate meta-model, which helps on coverage. On the other hand, samples with low predicted $\hat{P}(s)$ represents more challenging cases, which are preferred for the goal of challenge. We attempt to solve this dilemma in an ϵ -greedy way: on lines 6,7 of the Algorithm 1, we evaluate the queries with two sets of query

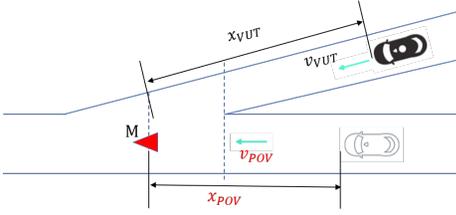


Fig. 4. The configuration of highway merging scenario.

quality metric according to the GPR model, $q_{exploit}(s)$ and $q_{explore}(s)$:

$$q_{exploit}(s) = \hat{\mu}^{z_1}(s) \hat{\sigma}^{z_2}(s) \quad (9)$$

$$q_{explore}(s) = \hat{\mu}^{z_3}(s) \hat{\sigma}^{z_4}(s) \quad (10)$$

where $\hat{\mu}(s) = \mathbb{E}[\hat{P}(s)]$, $\hat{\sigma}(s) = \mathbb{V}ar[\hat{P}(s)]$. $q_{exploit}(s)$ and $q_{explore}(s)$ have different parameters: the former prefers exploitation ($z_1 > z_2$), and the latter focuses on exploration ($z_3 < z_4$). For each batch, We pick the cases which maximize one of these metrics. The portion of cases for exploration and exploitation are determined by the parameter ϵ . It which will gradually decrease across batches at the rate of α ($\alpha \in (0.9, 1)$), such that the procedure starts with more exploration, and bias towards exploitation as more data are collected and a better meta-model is built.

D. Inter-level ratio adjustment

In this section, we will consider all the POV levels together by distributing cases into each level based on the results from the previous batch. Targeting on maximizing the expected coverage of the failure region, the strategy is to invest more samples in better-performing POV levels, while we keep exploring the other options. Specifically, we implement the following softmax decision rule on batch allocation:

$$n_k^{i+1} = \pi_k^{i+1} n = \frac{\exp(\xi U(i, k))}{\sum_{j=0}^2 \exp(\xi U(i, j))} n \quad (11)$$

where $U(i, k) = \frac{\# \text{ of cases with } P(s) < \text{threshold}}{n_k^i}$, $\sum_{j=0}^2 \pi_j^i = 1$. For the first batch, we distribute cases equally to all POV levels. After that, the cases are distributed according to the ratio of challenging cases found within that level in the previous batch. The parameter ξ controls how "greedy" the decision rule is.

VI. IMPLEMENTATION ON HIGHWAY MERGING SCENARIO

A. Scenario model

The highway merging scenario is the focus of this paper, for which the configuration is illustrated in Figure 4. The VUT attempts to merge onto the highway, while the POV is driving on the main lane of the road. We make the following assumptions for simplification:

- 1) The POV and the VUT do see and interact with each other through the simulation horizon.

- 2) The POV is not able to change lanes to yield to the VUT; the VUT can only merge at the merge point M , which is the origin for the lane-fixed coordinates for both the ramp and the main lane.
- 3) There is only one POV on the main lane and there is no vehicle in front of the VUT on the ramp.
- 4) The scenario ends when the VUT reaches point M .

We model both vehicles as double integrators and they only move longitudinally in their own lane. The equations of motion are:

$$\begin{cases} x_{POV}(t+1) = x_{POV}(t) + v_{POV}(t)\delta t \\ v_{POV}(t+1) = v_{POV}(t) + a_{POV}(t)\delta t \\ x_{VUT}(t+1) = x_{VUT}(t) + v_{VUT}(t)\delta t \\ v_{VUT}(t+1) = v_{VUT}(t) + a_{VUT}(t)\delta t \end{cases} \quad (12)$$

where x_{POV}, x_{VUT} are the longitudinal position, and v_{POV}, v_{VUT} are the longitudinal speed of POV and VUT in their lanes. The input for each vehicle is the longitudinal acceleration, which ranges between $[a_{min}, a_{max}]$. The initial condition is characterized by $(x_{POV}^0, v_{POV}^0, x_{VUT}^0, v_{VUT}^0)$. Without loss of generality, we assume x_{VUT}^0 is fixed. Moreover, v_{VUT}^0 is observed rather than determined by the test conductor. Therefore, the initial condition to sample from is $x_0 = [x_{POV}^0, v_{POV}^0]^T$.

B. Level-0 policy

For the highway merging scenario, a level-0 POV is assumed to keeps a constant speed, regardless of the VUT. A level-0 VUT will accelerate with constant acceleration ($1m/s^2$) until the assumed highway speed ($28m/s$).

C. Training RL agents at the highway merging scenario

When applied to the highway merging scenario, the physical state space of the MDP is $X = \mathbb{R}^4$, where each state is $x = [x_{POV}, v_{POV}, x_{VUT}, v_{VUT}]$. The transition dynamics are illustrated in (12), with the opponent's action governed by the level- $(k-1)$ policy. The actions are discrete acceleration choices within a_{min} and a_{max} for both POV and VUT, i.e. $u = a_{POV}/a_{VUT} \in U = \{-4, -3, \dots, 0, +1, +2\}(m/s^2)$. Each episode terminates when the VUT reaches the merge point $x_{VUT}(t_1) = 0$.

The detailed definitions of the three categories of reward mentioned in section IV-C for the highway merging scenario are as follow:

$\Phi_{POVe} = [\phi_{acc}, \phi_{v_{HW}}]^T$, where ϕ_{acc} penalizes acceleration action; $\phi_{v_{HW}}$ penalizes speed exceeding the highway speed limits (either v_{HWmin} or v_{HWmax}). The parameter values are shown in Table I.

$\Phi_{VUTe} = [\phi_{acc}, \phi_{v_{min}}, \phi_{v_{end}}]^T$, where ϕ_{acc} is the same as in Φ_{POVe} ; $\phi_{v_{min}}$ penalizes speed lower than a minimum speed v_{min} during the episode; $\phi_{v_{end}}$ penalizes final merging speed of the VUT that is faster or slower than highway speed limits.

TABLE I
PARAMETERS FOR REWARD DESIGN

v_{HWmax}	35.0 m/s	v_{HWmin}	24.6 m/s
v_{min}	12.0 m/s	TTC_{min}	7.0 s
Δx_{crash}	6 m	$\Delta x_{critical}$	15 m

$\Phi_{safe} = [\phi_{TTC}, \phi_{\Delta x}, \phi_{crash}]^T$ are the safety terms evaluated at the end of the episode t_1 . We define:

$$\begin{aligned} \Delta x_1 &= x_{POV}(t_1) - x_{VUT}(t_1) \\ \Delta v_1 &= v_{POV}(t_1) - v_{VUT}(t_1) \\ TTC &= \begin{cases} \frac{\Delta x_1}{-\Delta v_1} & \text{when } \Delta x_1 \Delta v_1 < 0 \\ \infty & \text{otherwise} \end{cases} \end{aligned}$$

where ϕ_{TTC} gives penalty when $TTC < TTC_{min}$; $\phi_{\Delta x}$ rewards large $|\Delta x|$, and gives penalty when $|\Delta x| < \Delta x_{critical}$; ϕ_{crash} gives heavy penalty when $|\Delta x| < \Delta x_{crash}$.

Finally, the DDQN algorithm for training the level- k POVs and VUTs is implemented using the MATLAB reinforcement learning toolbox and Simulink.

VII. SIMULATION RESULTS

We conduct interactive-aware testing to several baseline VUTs in simulations to validate the performance and benefits of the proposed method.

A. Baseline VUT

For the highway-merge scenario, we design a rule-based algorithm for the merging vehicle (the VUT). Its decision-making has 3 stages:

- 1) The VUT starts by following the speed profile of a level-0 VUT policy π_{VUT}^0 . Go to stage 2 when it is x_1^{rb} close to the merge point M .
- 2) The VUT predicts Δx relative to the POV when arriving at M , assuming the POV keeps a constant speed, and VUT follows π_{VUT}^0 . If too close, switch to coast; else, keep following π_{VUT}^0 . Go to stage 3 when VUT is x_2^{rb} close to M ($x_2^{rb} < x_1^{rb}$).
- 3) The VUT predicts Δx with the POV when arriving at M , assuming the POV maintains a constant speed, and VUT follows π_{VUT}^0 . If too close, switch to PID-control on acceleration; if not, follows π_{VUT}^0 .

By adjusting the parameters, we can manipulate the VUT to have different failure modes.

B. Various interactive test cases

In this section, we present exemplar test cases with different interactions between the POV and the VUT. The road geometry of the highway merging scenario is based on an entrance ramp on US 23 North near exit 41. The VUT started at $x_{VUT}^0 = -182[m]$. In Figure 5, we present three test cases with different POVs and VUTs. In all cases, the initial conditions are the same. In the 1st case, as shown in Figure 5(a), 5(b), the level-0 VUT is accelerating non-cooperatively, while the POV yields by reducing its speed to let the VUT merge first. In the 2nd scenario (Figure 5(c), 5(d)), the

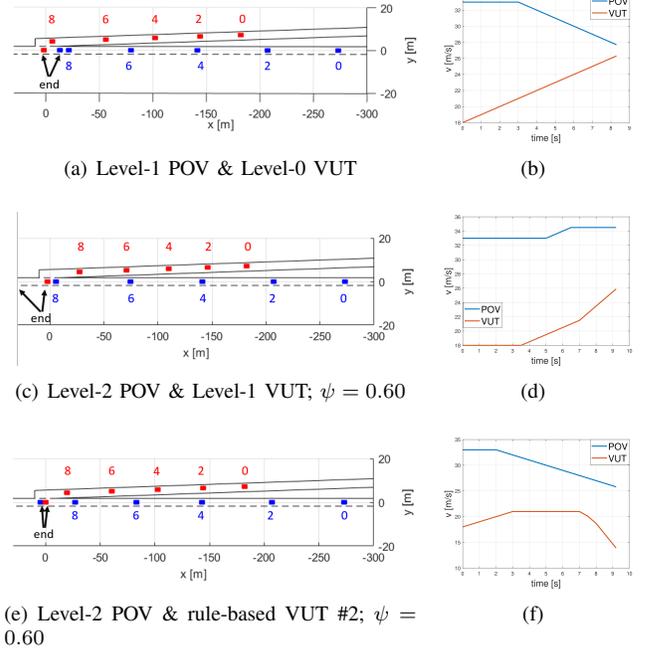


Fig. 5. Results with initial condition: $x_{POV}^0 = -273m$, $v_{POV}^0 = 33m/s$, $v_{VUT}^0 = 18m/s$; blue for VUT, red for POV; the numbers show time lapses in seconds.

level-1 VUT yields by starting its accelerating phase later, while the level-2 POV with a cooperative SVO accelerates to leave room for the VUT to merge behind. In the 3rd scenario, (Figure 5(e), 5(f)), the same level-2 POV yields to let the VUT enter first. However, the rule-based VUT fails to understand the POV's intention. It starts to accelerate, then coasts and even decelerates hard before it crashes with the VUT. This last case shows a "stalemate" situation, when both agents try to yield to each other and create an inefficient and dangerous scene. These three test cases capture different interactions, which makes the evaluation scenarios diverse.

C. Test results comparison

1) *Results with a single POV level:* We first show the results of simulated testing with a fixed POV level. Each VUT is put through $N = 400$ test cases. A case with a score $P(s) < -500$ means a collision has occurred, thus it is deemed as a failure case. We compare the proposed GPR-based adaptive sampling scheme to other test case generation schemes, including uniform sampling, simulated annealing [27], and subset simulation [10]. The FMC M is the criterion for comparing their capability of discovering failure cases. For the VUT, two rule-based algorithm designs are selected, denoted as design #1 and design #2, where design #2 has a faster response and is deemed "smarter". We test the two designs against Level-0 to Level-2 POVs. All methods are compared against the ground truth, which is generated by 10000 samples for level-0,1 and 20000 samples for level-2 POV using uniform sampling from the test case pool. The quantitative results comparison is shown in table II. For all the three combinations of POV and VUT, the GPR-based

TABLE II
ADAPTIVE TESTING RESULTS COMPARISON

Methods	FMC $M(s, 0.05, -500)$		
	L-0 POV #1 VUT	L-1 POV #2 VUT	L-2 POV #2 VUT
Ground truth	0.1296	0.0849	0.0073
Uniform sampling	0.0223	0.0144	0
Simulated annealing	0.0497	0.0184	0.0006
Subset simulation	0.0830	0.0445	0.0019
GPR-based sampling	0.1012	0.0585	0.0032

adaptive sampling achieves the highest FMC among all the methods, and is also closest to the ground truth using only 4% of the cases.

Specifically, Figure 6 compares the results for testing VUT design #2 against level-1 POV, which has three disjoint failure regions according to the ground truth. Uniform sampling can locate one failure region with very few failure cases; simulated annealing can find only one failure region, with many test cases concentrated around one local minimum; both subset simulation and GPR-based sampling can identify all three failure regions with only 4% of the samples compared to the ground truth, while the GPR-based method achieves higher FMC and reconstructs the shape of the failure regions better. In Figure 7, the progression of the GPR meta-model is displayed, where it finds more accurate failure modes as the batch number grows and test results accumulate. Figure 8 compares the results of testing with a level-2 POV. While the GPR-based method can find the two failure modes far away from each other, subset simulation can only identify one of them with same number of tests.

2) *Results with multiple POV levels:* Finally, we simulate the adaptive testing procedure with all the POV levels for VUT design #2. The goal is to identify failure modes in all three levels within $N = 800$ cases. Figure 9 shows the change of sample allocation across different POVs. The sample sizes start evenly, but since more failure cases were found with level-1, and none for level-0 POV, The sample size grows for level-1 in later batches while shrinks for level-0. It is demonstrated that the proposed method is able to focus on the more promising interactive POV level for efficient identification of challenging test cases, while keeping exploring the under-performed ones.

VIII. CONCLUSIONS

In this paper, we study the evaluation problem for black-box HAVs in scenarios with significant human interactions. We apply two game-theoretic methodologies, level- k game theory and social value orientation, to model the interactive POV driving policies and incorporate them into our test case pool design. Then, we design an adaptive test case sampling scheme based on Gaussian process regression and propose a metric to assess failure mode coverage (FMC) to measure the test sample quality. We verify the proposed method by running simulated testing on several baseline VUTs. The POV library is able to emulate a wide variety of interactive behaviors and the sampling method can customize test cases to discover the failure modes of the VUTs by using only

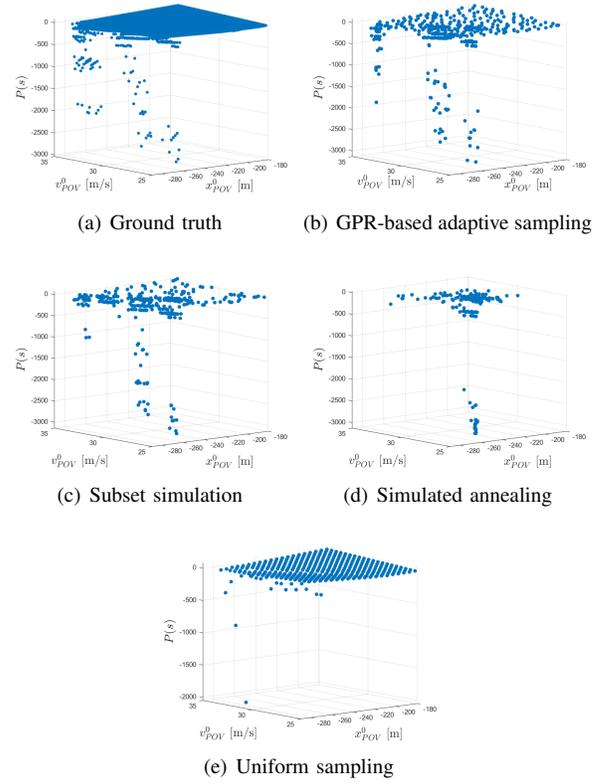


Fig. 6. Simulated testing results with level-1 POV on VUT design #2; for (b)-(e), $k = 20$, $n = 20$.

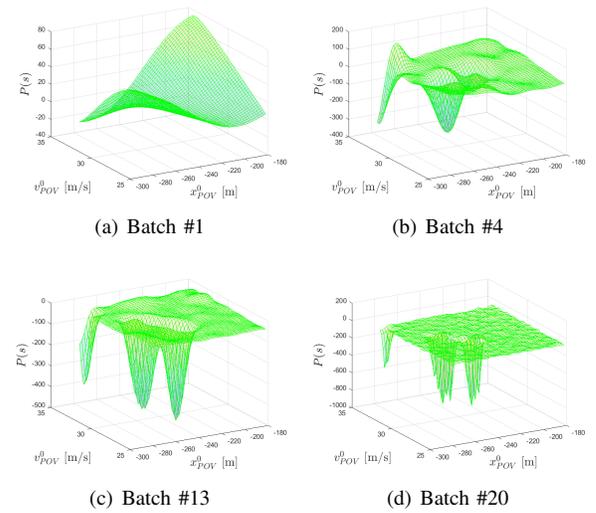


Fig. 7. The evolution of the GPR meta-model, with level-1 POV on VUT design #2; $k = 20$, $n = 20$.

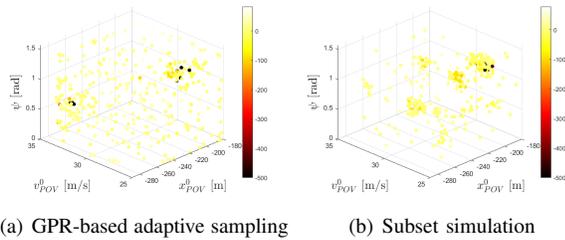


Fig. 8. Simulated testing results with level-2 POV on VUT design #2; $k = 20$, $n = 20$.

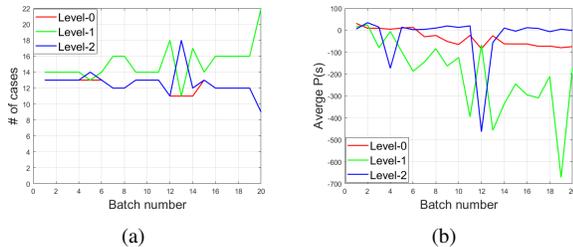


Fig. 9. Simulated testing results with the full POV library on VUT design #2; $k = 20$, $n = 40$. (a) Batch sample allocation across POV levels in all batches. (b) Average performance score for each POV level in all batches.

a fraction of the number of cases compared to the ground truth. It out-performs other sampling methods according to the FMC metric.

ACKNOWLEDGMENT

Toyota Research Institute (TRI) provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

We thank Shaobing Xu, Geunseob Oh, Yuanxin Zhong for their insightful suggestions and help.

REFERENCES

- [1] 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, 3 2017.
- [2] X. Wang, H. Peng, and D. Zhao, "Combining reachability analysis and importance sampling for accelerated evaluation of highly automated vehicles at pedestrian crossing," in *ASME 2019 Dynamic Systems and Control Conference, DSCC 2019*, vol. 3. American Society of Mechanical Engineers, 10 2019.
- [3] X. Wang, Y. Dong, S. Xu, H. Peng, F. Wang, and Z. Liu, "Behavioral Competence Tests for Highly Automated Vehicles," in *Accepted by IEEE Intelligent Vehicles Symposium*, 2020.
- [4] "Automated Vehicles for Safety | NHTSA." [Online]. Available: <https://www.nhtsa.gov/technology-innovation/automated-vehicles-safety>
- [5] NCAP, "European New Car Assessment Programme - TEST PROTOCOL – AEB systems," Tech. Rep., 2015.
- [6] G. Chou, Y. E. Sahin, L. Yang, K. J. Rutledge, P. Nilsson, and N. Ozay, "Using control synthesis to generate corner cases: A case study on autonomous driving," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2906–2917, 2018.
- [7] M. Althoff and S. Lutz, "Automatic Generation of Safety-Critical Test Scenarios for Collision Avoidance of Road Vehicles," in *2018 IEEE Intelligent Vehicles Symposium (IV)*, vol. 2018-June. IEEE, 6 2018, pp. 1326–1333.

- [8] C. E. Tuncali, T. P. Pavlic, and G. Fainekos, "Utilizing S-TaLiRo as an automatic test generation framework for autonomous vehicles," *IEEE Conference on Intelligent Transportation Systems, Proceedings, ITSC*, no. ii, pp. 1470–1475, 2016.
- [9] S. Zhang, H. Peng, S. Nageshrao, and H. E. Tseng, "Generating socially acceptable perturbations for efficient evaluation of autonomous vehicles," *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, vol. 2020-June, pp. 1341–1347, 2020.
- [10] S. Zhang, H. Peng, D. Zhao, and H. E. Tseng, "Accelerated Evaluation of Autonomous Vehicles in the Lane Change Scenario Based on Subset Simulation Technique," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 11 2018, pp. 3935–3940.
- [11] G. E. Mullins, P. G. Stankiewicz, and S. K. Gupta, "Automated generation of diverse and challenging scenarios for test and evaluation of autonomous vehicles," *Proceedings - IEEE International Conference on Robotics and Automation*, pp. 1443–1450, 2017.
- [12] Z. Huang, H. Lam, and D. Zhao, "Towards Affordable On-track Testing for Autonomous Vehicle - A Kriging-based Statistical Approach," *IEEE Conference on Intelligent Transportation Systems, Proceedings, ITSC*, vol. 2018-March, pp. 1–6, 7 2017.
- [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*, pp. 1–10, 3 2020.
- [14] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, "Social behavior for autonomous vehicles," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 116, no. 50, pp. 2492–24978, 2019.
- [15] J. F. Fisac, E. Bronstein, E. Stefansson, D. Sadigh, S. S. Sastry, and A. D. Dragan, "Hierarchical game-theoretic planning for autonomous vehicles," in *Proceedings - IEEE International Conference on Robotics and Automation*, vol. 2019-May, 2019, pp. 9590–9596.
- [16] J. H. Yoo and R. Langari, "A stackelberg game theoretic driver model for merging," *ASME 2013 Dynamic Systems and Control Conference, DSCC 2013*, vol. 2, pp. 1–8, 2013.
- [17] N. Li, D. W. Oyler, M. Zhang, Y. Yildiz, I. Kolmanovsky, and A. R. Girard, "Game theoretic modeling of driver and vehicle interactions for verification and validation of autonomous vehicle control systems," *IEEE Transactions on Control Systems Technology*, vol. 26, no. 5, pp. 1782–1797, 2018.
- [18] A. Sarkar and K. Czamecki, "A behavior driven approach for sampling rare event situations for autonomous vehicles," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 11 2019, pp. 6407–6414.
- [19] L. Sun, W. Zhan, Y. Hu, and M. Tomizuka, "Interpretable Modelling of Driving Behaviors in Interactive Driving Scenarios based on Cumulative Prospect Theory," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 10 2019, pp. 4329–4335.
- [20] L. Sun, W. Zhan, M. Tomizuka, and A. D. Dragan, "Courteous Autonomous Cars," *IEEE International Conference on Intelligent Robots and Systems*, pp. 663–670, 2018.
- [21] M. A. Costa-Gomes, V. P. Crawford, and N. Iriberry, "Comparing Models of Strategic Thinking in Van Huyck, Battalio, and Beil's Coordination Games," *Journal of the European Economic Association*, vol. 7, no. 2-3, pp. 365–376, 4 2009.
- [22] C. G. McClintock and S. T. Allison, "Social Value Orientation and Helping Behavior," *Journal of Applied Social Psychology*, vol. 19, no. 4, pp. 353–362, 3 1989.
- [23] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [24] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-Learning," in *30th AAAI Conference on Artificial Intelligence, AAAI 2016*, 2016, pp. 2094–2100.
- [25] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [26] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. Cambridge, Massachusetts: MIT Press, 2006.
- [27] L. C. W. Dixon and G. P. Szegő, *Towards global optimisation*. North-Holland Amsterdam, 1978, vol. 2.