

# Behaviorally Diverse Traffic Simulation via Reinforcement Learning

Shinya Shiroshita<sup>1</sup>, Shirou Maruyama<sup>1</sup>, Daisuke Nishiyama<sup>1</sup>, Mario Ynocente Castro<sup>1</sup>,  
Karim Hamzaoui<sup>1</sup>, Guy Rosman<sup>2</sup>, Jonathan DeCastro<sup>2</sup>, Kuan-Hui Lee<sup>2</sup> and Adrien Gaidon<sup>2</sup>

**Abstract**—Traffic simulators are important tools in autonomous driving development. While continuous progress has been made to provide developers more options for modeling various traffic participants, tuning these models to increase their behavioral diversity while maintaining quality is often very challenging. This paper introduces an easily-tunable policy generation algorithm for autonomous driving agents. The proposed algorithm balances diversity and driving skills by leveraging the representation and exploration abilities of deep reinforcement learning via a distinct policy set selector. Moreover, we present an algorithm utilizing intrinsic rewards to widen behavioral differences in the training. To provide quantitative assessments, we develop two trajectory-based evaluation metrics which measure the differences among policies and behavioral coverage. We experimentally show the effectiveness of our methods on several challenging intersection scenes.

## I. INTRODUCTION

Modeling driving policies is an important and challenging problem in autonomous vehicles research, with a growing range of promising applications. One such example is the use in behavior prediction modules, which aim to help automated vehicles estimate future trajectories of other vehicles and take necessary safety measures. Another application of interest is traffic simulation, which we explore in this work.

Modeling human drivers is an important step towards a better understanding of the complex interactions among all traffic participants. One difficulty is reconciling driving skills and diversity at the same time. Human drivers are generally skilled, considering the complexity of the driving task. Conversely, they also display a variety of driving characteristics, including aggressive, conservative, or careless types, which make the driving behavior sub-optimal. Therefore, it is important to strike a balance between expanding driving behavioral diversity while maintaining high driving skills.

As shown in Fig.1a, one way to impose diversity in a simulator is to use random policies, but the resulting vehicles will have non-sensical trajectories and become useless for simulating realistic traffic situations. Another way is to use a cost function that expresses driving skills when modeling a policy, which will result in the intended driving type. However, trying to produce various types of driving behaviors using such a method usually results only in minor deviations from optimal behavior. The goal of our research is to acquire

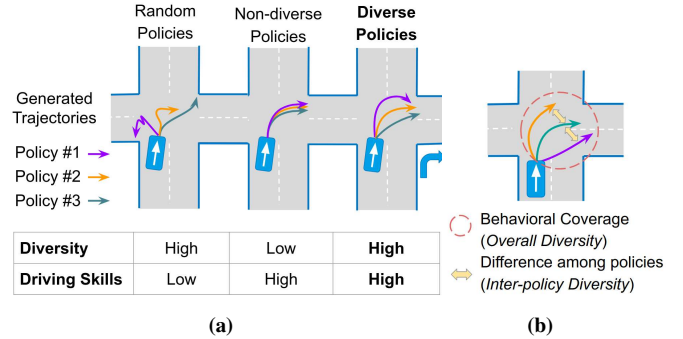


Fig. 1: (a) Comparison of random policies, non-diverse policies, and diverse policies and their driving skills. (b) Inter-policy diversity and overall diversity.

behaviorally diverse driving policies while maintaining high driving skills.

Various approaches have been proposed for diverse behavior modeling. SUMO [1] provides hand-tuned rule-based policies limited to scripted trajectories in complex scenes such as intersections. Data-driven approaches (see [2]–[4] and references therein) have difficulty in representing novel behaviors. Previous reward shaping methods (e.g. Hu et al. [5]) made use of reinforcement learning (RL) as a way of generating diversity. However, manually adjusting the reward function to produce various driving policies is difficult for complex traffic scenes such as interactions. In this work, our main contribution is an RL algorithm using intrinsic rewards and a highly expressive model to generate behaviorally diverse driving policies (Fig. 1b) without requiring additional tuning per target behavior.

To define and reason about the diversity of a set of driving policies, we propose as our second contribution two evaluation metrics based on generated trajectories, reflecting both higher-order decisions and lower-order controls (Fig. 1b). Our first metric, Inter-Policy Diversity, measures how different two policies are from each other (Fig. 1b). Our second metric, Overall Diversity, represents the overall behavioral coverage (Fig. 1b). Since the difference in trajectories largely depends on the traffic scene, including the behavior of other vehicles, both metrics calculate an average of the trajectory distances across multiple traffic scene simulations.

To produce diverse policies with high driving skills, our method first generates various snapshots leveraging the high exploration ability of reinforcement learning, then uses our inter-policy diversity metric to select snapshots where policies express the most distinct behavior. For the snapshot se-

<sup>1</sup>Preferred Networks, Inc., Japan., {shiroshita, maruyama, dnishiyama, marioyc, karim}@preferred.jp

<sup>2</sup>Toyota Research Institute, U.S., {guy.rosman, jonathan.decastro, kuan.lee, adrien.gaidon}@tri.global

lection, we filter out those of lower success rates and choose policies by a method similar to Farthest Point Sampling (FPS) [6]. Since diversity from policy selection relies on candidate policies, it is crucial to generate diverse candidates. To acquire them, we use large scale exploration and a variant of Diversity-Driven Exploration (DDE) [7] to enlarge the differences among training policies.

We evaluate our approach on an intersection simulation environment to confirm that it efficiently selects diverse policies. We also qualitatively check that it produces a variety of cooperation models.

Our contributions can be summarized as follows:

- We propose a method of acquiring policies that can balance driving skills and diversity.
- We develop trajectory-based diversity evaluation measures for a set of driving policies.
- We show the effectiveness of our method for diverse policy acquisition through experiments using a simulation platform which can model vehicle-to-vehicle interactions at intersections.

One application of this work is conducted within another work of ours on planner testing [8], where diverse policies are used in surrounding vehicles aiming to detect more diverse failure cases.

The remaining parts are organized as follows. Section II and III mention related works and preliminaries, respectively. Our main contribution is described in the following three sections: Section IV for the definition of diversity evaluation metrics, Section V for the explanation of our policy generation algorithms optimizing the metrics, and Section VI for the simulator to evaluate our approach’s effectiveness. Section VII summarizes the experimental setting and results. Section VIII concludes our work.

## II. RELATED WORKS

### A. Driving policies in existing traffic simulators

Many simulators have been developed focusing on various aspects such as photo-realistic views with a variety of traffic participants (CARLA [9]), bird-view traffic flow (SUMO [1]), light-weighted simulation environment (FLUIDS [10]), multi-agent framework (CoInCar-Sim [11]), and benchmark scenarios (CommonRoad [12]).

In terms of driving policies, however, the behavioral diversity of other vehicle policies is not a major focus. Existing policies consist of simple rule-based policies [1], [9]–[12] while some make use of reinforcement learning and imitation learning [9]. However, in all of them, behavioral diversity can be obtained by specifying it manually, and it is not shown how to select a diverse policy set. Our research purpose is to establish algorithms that can generate policies for these simulators, which can sustain diversity and good driving skills without having to manually specify one-by-one.

### B. Diversity acquisition for driving policies

Hu et al. [5] and Bacchiani et al. [13] generate diverse driving policies by controlling the environmental settings

(especially rewards) in merging and roundabout scenes, respectively. However, these methods depend on specifying such settings manually, and they do not provide a quantitative comparison of their obtained diversity. In contrast, our work puts more focus on diverse policy acquisition efficiency and its quantitative evaluation.

In the RL field, diversity appears as a policy exploration tool. Diversity-Driven Exploration (DDE) [7] adds a bonus reward when a policy takes different actions from those taken by previous policies. Stein Variational Policy Gradient (SVPG) [14] maximizes the entropy of policy parameters to make policies behave differently. Diversity is All You Need (DIAYN) [15] combines a policy with a discriminator network to distribute states to each action type. See also Aubret et al. [16] for a survey about intrinsic rewards. To the best of our knowledge, our work is the first to investigate RL with intrinsic rewards for traffic simulation tasks. These RL algorithms generally focus on improving benchmark scores instead of providing a diverse simulation platform.

### C. Diversity evaluation metrics

In the field of behavior prediction, diversity is measured through the prediction error. Social GAN [2] and DRO-GON [17], for example, evaluated the minimum prediction error among  $n$  samples. However, this metric requires diverse expert trajectories and corresponding maps in the simulator, which is not applicable to any traffic scenes. Also, the quality of diversity depends on experts and is not directly measured. Our metrics, on the other hand, evaluate diversity and behavioral coverage separately. We also show an approximation of diverse trajectory sets to show coverage without experts.

## III. PRELIMINARIES

Reinforcement learning is a method to learn policies by interacting with an environment. It assumes the environments as a *Markov decision process (MDP)*. Let  $\mathcal{S}$ ,  $\mathcal{A}$ ,  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ ,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  be the state space, the action space, the transition probability, and the reward function, respectively. Then, the goal is to find a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  to maximize accumulated rewards with a discount factor  $0 < \gamma < 1$ :

$$J = \sum_{t=0}^{\infty} \gamma^t R_{t,\pi} \quad (1)$$

where  $R_{t,\pi}$  is the expected reward with the policy  $\pi$  at the step  $t$ . In traffic environments, the observation  $o$  of the ego vehicle is usually not equal to the (global) state since some information such as occluded objects or other vehicle’s policy is invisible to the ego vehicle. Therefore, the policy  $\pi$  becomes a mapping from the set of observations to the set of actions. Section VI shows the details of the observation.

## IV. DIVERSITY EVALUATION METRICS

We firstly define objective functions for diversity acquisition algorithms shown in the next section. For the sake of evaluating the discriminability and the coverage of policies, we consider two types of diversity measures: *inter-policy*

*diversity* and *overall diversity*. The first metric, inter-policy diversity, focuses on the difference between the trajectories generated by pairs of policies. The second metric, overall diversity, measures the difference between the trajectories generated by a target policy and some particular set of reference trajectories.

Our metrics consider the distances of generated trajectories for each *scenario*, a setting of the surrounding environment, such as a map, initial positions, and other vehicles' policies. A *trajectory*  $\tau = ((x_t, y_t))_{t=1, \dots, |\tau|}$  is a sequence of two dimensional coordinates. We employed the average Euclidean distance for two trajectories  $\tau, \tau'$  like

$$d(\tau, \tau') = \frac{1}{T} \sum_{t=1}^T \|\tau(t) - \tau'(t)\|_2 \quad (2)$$

where  $T = \min\{|\tau|, |\tau'|\}$ . For a scenario  $s$  and a policy  $\pi$ , let  $\tau_s(\pi)$  be a trajectory generated by  $\pi$  in  $s$ . We assume that each scenario also fixes random seeds to ensure deterministic behavior, and thus there is a mapping from the ego policy to the resulted trajectory.

#### A. Inter-policy diversity

Let  $S, \Pi$  be a set of scenarios and a policy set, respectively. Let  $S_\pi \subseteq S$  be a set of scenarios where a policy  $\pi$  succeeds the mission to reach a given goal within a time limit without collision (see Section VII-A for more details), and similarly let  $\Pi_s \subseteq \Pi$ . Note that for any scenario  $s$  and any policy pair  $\pi, \pi'$ , we assume that  $\Pi_s \neq \emptyset$  and  $S_\pi \cap S_{\pi'} \neq \emptyset$ . Exceptions rarely happen since policies with low success ratios are removed at the policy selection step.

We define the inter-policy diversity  $D_{IP}$  of a policy set  $\Pi$  by the average distance of two policy pairs like

$$D_{IP}(\Pi) = \frac{1}{|\Pi|(|\Pi| - 1)} \sum_{\pi \in \Pi} \sum_{\pi' \in \Pi \setminus \{\pi\}} D_{IP}(\pi, \pi') \quad (3)$$

where

$$D_{IP}(\pi, \pi') = \frac{1}{|S_\pi \cap S_{\pi'}|} \sum_{s \in S_\pi \cap S_{\pi'}} d(\tau_s(\pi), \tau_s(\pi')) \quad (4)$$

means an average trajectory distance of two policies  $\pi, \pi'$ . Our work utilizes the diversity for discriminability evaluation and a basis for the policy set selection.

#### B. Overall diversity

Overall diversity calculates the distance between the obtained trajectory and the expected trajectories called *reference trajectories*  $\mathcal{T}$ . To capture both the density and the spatial similarity of trajectories, we employ the Wasserstein-1 distance as a distance measure of distributions. For each scenario  $s$ , let  $\tau_s(\Pi) = \{\tau_s(\pi) | \pi \in \Pi_s\}$  be a set of successful trajectories of  $\Pi$  in  $s$ . Then, we define the overall diversity  $D_{OA}^{\mathcal{T}, s}$  of a policy set  $\Pi$  as

$$D_{OA}^{\mathcal{T}, s}(\Pi) = \inf_{\gamma \in \Gamma(\tau_s(\Pi), \mathcal{T})} \mathbb{E}_{(\tau, \tau') \sim \gamma} [d(\tau, \tau')] \quad (5)$$

where  $\Gamma(\tau_s(\Pi), \mathcal{T})$  is a set of all possible joint distributions of  $\tau_s(\Pi)$  and  $\mathcal{T}$ . We finally obtain the overall diversity  $D_{OA}^{\mathcal{T}}$  of the policy set by the average over scenarios like

$$D_{OA}^{\mathcal{T}}(\Pi) = \frac{1}{|S|} \sum_{s \in S} D_{OA}^{\mathcal{T}, s}(\Pi_s). \quad (6)$$

Note that a lower value indicates better performance.

Since counting all possible trajectories is infeasible, we need to set a well-approximated  $\mathcal{T}$  covering behaviors as diverse as possible. We modeled  $\mathcal{T}$  as a set of *Brownian bridges* modified to take into consideration vehicle dynamics and surrounding vehicles. The process of generating the bridges consists of the following steps.

- 1) Define an expected right-turn trajectory manually. A generator draws diverse trajectories along with that trajectory.
- 2) Draw *target trajectories*. These trajectories are created by perturbing the expected trajectory with longitudinal and lateral movements, respectively. The longitudinal movement is for speed variation and the lateral movement is for steering variation, respectively. Each movement is generated by Brownian bridges, which are stochastic processes whose initial and end points are fixed, to restrict the destination and arrival time.
- 3) Run a Proportional control (P-control) vehicle trying to mimic the target trajectories. Let  $\tau^{\text{target}}$  and  $\tau$  denote the target and the converted trajectory, respectively. The vehicle at time  $t$  will try to move from point  $\tau(t)$ , in  $\nu$  seconds, to the point  $\tau^{\text{target}}(t + \nu)$  which is the point in the original trajectory  $\nu$  seconds later.

We run the above procedure in the intersection environment shown in Section VI. If the agent reaches the goal without collision, the trace of the vehicle's center is added to the output. Fig. 7a in the experimental section shows an example of the generated set of trajectories in one scenario.

### V. DIVERSITY ACQUISITION ALGORITHMS

This section shows the diversity acquisition modules that generate a set of policies with better inter-policy and overall diversity values defined in the previous section. Since there exists a trade-off between behavioral diversity and driving quality, we put a restriction on policies where their mission success rates must not be worse than a threshold  $\delta$ .

Our strategy is as follows: we train a set of  $n (\geq 2)$  policies  $\Pi$  simultaneously and select snapshots (not only the latest ones but also those during training) in order to maximize the inter-policy diversity. Details are explained in Section V-A. In the case of environments where random policies tend to converge into the same policy, we also propose an algorithm which uses intrinsic rewards in Section V-B.

#### A. Diverse policy selection

Reinforcement learning produces a variety of policies through the exploration process guided by initial parameters and randomness during training. The basic strategy is to first create many snapshots of the policies during their training process. Let  $\Pi_{\text{all}}$  be a set of candidate policies that consists

---

**Algorithm 1** Diverse policy selection
 

---

**Input:** Policies  $\Pi_{all}$ , scenarios  $S$ , driving score threshold  $\delta$ , number of selected policies  $k$

**Output:** Set of  $k$  policies  $Q_k$

$\Pi \leftarrow \{\pi \in \Pi_{all} \mid \pi\text{'s driving score is not less than } \delta\}$

$Q_1 \leftarrow \{\pi_1\}$  where  $\pi_1$  is randomly selected from  $\Pi$

**for**  $i = 2, \dots, k$  **do**

$\pi_i \leftarrow \operatorname{argmax}_{\pi \in \Pi \setminus Q_{i-1}} \min_{\pi' \in Q_{i-1}} D_{IP}(\pi, \pi')$

$Q_i \leftarrow Q_{i-1} \cup \{\pi_i\}$

**end for**

**return**  $Q_k$

---

of snapshots generated during the training process. Firstly, we remove policies from  $\Pi_{all}$  with a driving score (e.g. success ratio) is less than the threshold  $\delta$  to filter out the policies with low driving skills. Let  $\Pi \subseteq \Pi_{all}$  be the filtered set, and  $k$  be the number of agents to be selected. Then, our target is to find a subset  $Q_k^* \in \mathcal{P}_k(\Pi) = \{Q \mid Q \in 2^\Pi \text{ and } |Q| = k\}$  such that

$$Q_k^* = \operatorname{argmax}_{Q \in \mathcal{P}_k(\Pi)} (D_{IP}(Q)) \quad (7)$$

which maximizes the inter-policy diversity.

However, finding an optimal subset from all possible choices is computationally infeasible since  $|\mathcal{P}_k(\Pi)|$  exponentially increases according to  $k$ . Therefore, we propose a greedy method based on farthest point sampling [6] to find an approximation of  $Q_k^*$ .

Our method first selects one policy at random and repeatedly adds a new one where the average inter-policy diversity with the selected policies is the highest among the remainder. Algorithm 1 shows the precise procedure. Although this algorithm considers only inter-policy diversity, our experimental results suggest that the generated set of policies is also diverse in terms of overall diversity due to the spatially separated selection process.

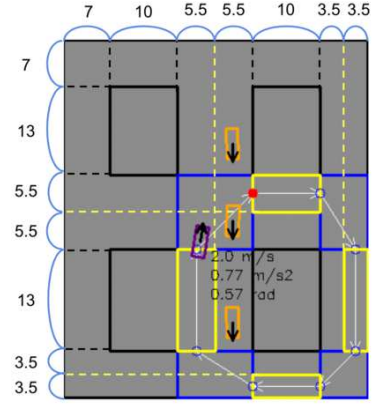
### B. Diverse policy training with intrinsic rewards

Another way of enhancing diversity is by utilizing intrinsic rewards on different actions to enlarge policy differences, especially for reward settings where policies with different initial parameters are likely to converge into similar ones. Our approach is a modification of DDE [7]. To put a repulsive force among training policies, our method runs  $m$  training threads in parallel and set intrinsic rewards as

$$r^{(env)} + \mathbb{E}_{\pi' \in \hat{\Pi} \setminus \{\pi\}} [\alpha d_{KL}(\pi, \pi')] \quad (8)$$

where  $r^{(env)}$  is the environmental reward,  $\pi$  is a current policy,  $\hat{\Pi}$  is the current policy set beyond all threads (while the original DDE used a set of previous policies),  $\alpha$  is the weight of the intrinsic reward, and  $d_{KL}(\pi, \pi')$  is the KL divergence of action distributions of  $\pi$  and  $\pi'$ , respectively.

Although the original DDE [7] used a tunable  $\alpha$ , we fixed  $\alpha$  in the experiments.



**Fig. 2:** The scales of the simulation map. All units are meters. The width and length of each vehicle are 1.8 and 4.5 meters, respectively.

## VI. TRAFFIC SIMULATION

To verify our proposed evaluation and method, we have implemented a simulation environment with intersection-based traffic scenes. This simulator can perform simulation of multiple agents. The goals of each agent include: (i) moving along with the given navigation route as fast as possible, and (ii) performing collision-free movements to walls and other vehicles.

### A. Map

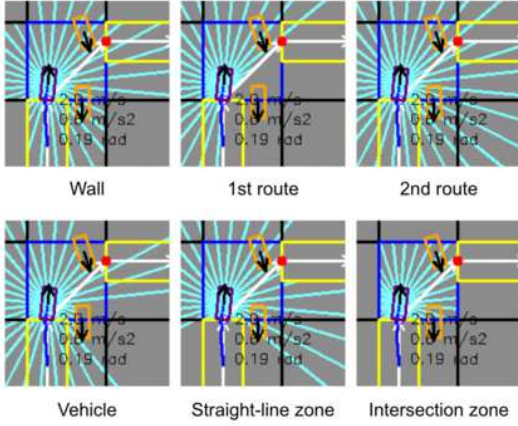
The map we prepared follows the Japanese traffic rule, that is, vehicles on the road drive on the left side of the driving lanes. Fig. 2 shows the scales of the map. The purple rectangle and orange ones represent the ego vehicle and the other vehicles, respectively. The black arrow on a vehicle represents its heading. Each vehicle is given a predefined navigation route consisting of a set of target arrows (white arrows) the vehicle is expected to follow. Each target arrow is combined with a rectangular zone to define rewards for the vehicle. The zone has two types: straight-line zones (yellow rectangles) and intersection zones (blue rectangles) with the different reward calculation formula. The map also contains walls (black lines) which vehicles cannot pass beyond.

### B. Episode termination conditions

In the evaluation phase, the episode is stopped if a collision occurs, if the ego vehicle reaches a predefined goal area, or if the episode time limit expires. In training, there is no fixed goal, and the trainee tries to go around the navigation route forever unless the time is over or it has more than 150 collisions in an episode.

### C. Vehicle Dynamics

The simulator regards each vehicle as a rectangle running under the kinematic bicycle model [18]. The length from the gravity center to each wheel is equal to the half of the vehicle length. The control input for each vehicle at time  $t$  is specified as  $u_t = (\phi_t, a_t)$ , where  $\phi_t \in [-0.785, 0.785]$  is the steering angle (rad) and  $a_t \in [-1.0, 1.0]$  is the acceleration ( $\text{m/s}^2$ ), respectively. According to the control input, the



**Fig. 3:** Six types of distances for observation. Each light blue line emitted from the purple rectangle represents the distance of the captioned parameter in the corresponding direction.

**TABLE I:** Correspondence of the action types.

Action	$\gamma_1$	$\gamma_2$	Action	$\gamma_1$	$\gamma_2$
Forward	0	2.5	Right-forward	0.628	2.5
Backward	0	-2.5	Left-forward	-0.628	2.5
Right	0.628	0	Right-backward	0.628	-2.5
Left	-0.628	0	Left-backward	-0.628	-2.5
Holding	0	0			

vehicle state transition from  $c_t = (x_t, y_t, \theta_t, v_t)$  to  $c_{t+\Delta t}$  is computed where  $x_t$  and  $y_t$  represent the two dimensional coordinates of the vehicle center (m),  $\theta_t$  is the heading of the vehicle (rad),  $v_t$  is the velocity of the vehicle (m/s), and  $\Delta t$  is time difference between two adjacent frames. The velocity range is limited to  $v_t \in [0.0, 2.0]$ , which is sufficient for driving at an intersection.

#### D. Observation and Action Space

In terms of observation, the agent acquires six types of distances for each of 32 directions  $\mu_w, \mu_{r1}, \mu_{r2}, \mu_v, \mu_{z^s}, \mu_{z^i} \in \mathbb{R}^{32}$ . Each value, visualized in Fig. 3, is defined as, the shortest distance (50 meters is maximum) to walls, the closest navigation route, the second closest one, other vehicles, the straight-line zone boundaries, and the intersection zone ones, respectively. Combining with the last three ego vehicle states  $h_{t-\Delta t}, \dots, h_{t-3\Delta t}$  where  $h_t = (v_t, a_t, \phi_t) \in \mathbb{R}^3$ , the overall observation becomes

$$o_t = [\mu_w, \mu_{r1}, \mu_{r2}, \mu_v, \mu_{z^s}, \mu_{z^i}, h_{t-\Delta t}, h_{t-2\Delta t}, h_{t-3\Delta t}] \in \mathbb{R}^{201}. \quad (9)$$

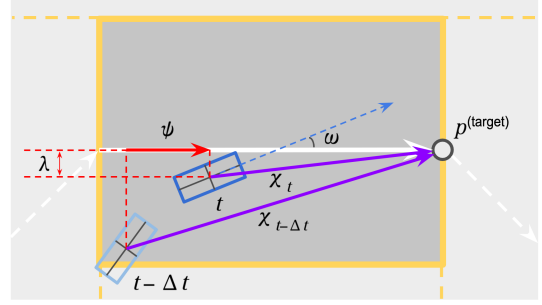
Actions, on the other hand, consist of nine discrete pairs of  $(\gamma_1, \gamma_2)$  shown in Table I. The control input at time  $t$  is determined as  $u_t = (\gamma_1 \Delta t + \phi_{t-\Delta t}, \gamma_2 \Delta t + a_{t-\Delta t})$ .

#### E. Reward Function

The environmental reward  $r^{(env)}$  consists of four types of rewards  $r^{(move)}$ ,  $r^{(collision)}$ ,  $r^{(angle)}$ ,  $r^{(center)}$  as

$$r^{(env)} = r^{(move)} + r^{(collision)} + r^{(angle)} + r^{(center)}. \quad (10)$$

These rewards are parametrized by non-negative weight coefficients  $w^{(move)}$ ,  $w^{(collision)}$ ,  $w^{(angle)}$ , and  $w^{(center)}$ . Fig. 4 summarizes variables used for each reward's description.



**Fig. 4:** Variables used for reward calculation. The yellow rectangle represents a zone and the dark blue one represents the ego vehicle at time  $t$ .

1) *Movement Reward:*  $r^{(move)}$  is specified based on the progress toward the desired direction, which depends on the zone the ego vehicle is in. In the intersection zones,

$$r^{(move)} = w^{(move)} \cdot \max\{0, \chi_{t-\Delta t} - \chi_t\} \quad (11)$$

where  $\chi_{t-\Delta t}, \chi_t$  are the Euclidean distances from the end point  $p^{(target)}$  of the target arrow to  $(x_{t-\Delta t}, y_{t-\Delta t}), (x_t, y_t)$ , respectively. In the straight-line zones,

$$r^{(move)} = w^{(move)} \cdot \psi \quad (12)$$

where  $\psi$  is the travel distance from frame  $t-\Delta t$  to  $t$  projected on the target arrow ( $\psi$  can be negative if the ego vehicle recedes). If the ego vehicle exists in neither,  $r^{(move)} = 0$ .

2) *Collision Reward:*  $r^{(collision)}$  is defined as

$$r^{(collision)} = \begin{cases} -w^{(collision)} & \text{if a collision occurs} \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

3) *Angle Reward:*  $r^{(angle)}$  is based on the difference between angles of the ego vehicle and the expected route. When the agent is in a straight-line zone and  $r^{(move)} \geq 0$ ,

$$r^{(angle)} = w^{(angle)} \cdot (0.5 - (\omega/\pi)^2). \quad (14)$$

where  $\omega$  is the difference (rad) of  $\phi_t$  and the target arrow's angle.  $r^{(angle)} = 0$  otherwise.

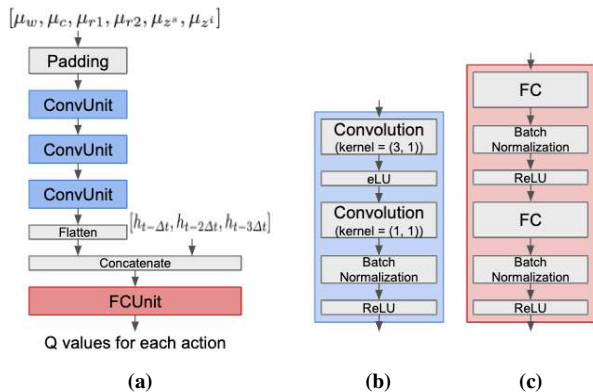
4) *Center-line Reward:*  $r^{(center)}$  represents how close the ego vehicle is from the route. When the agent is in a straight-line zone and  $r^{(move)} \geq 0$ ,

$$r^{(center)} = w^{(center)} \cdot (5 \cdot \exp(-8 \cdot \lambda^2) - 0.5) \quad (15)$$

where  $\lambda$  is the shortest distance from  $(x_t, y_t)$  to the target arrow.  $r^{(center)} = 0$  otherwise.

#### F. Reinforcement Learning and Neural Network

We used the double PAL as the reinforcement learning algorithm and the distances convolutional network model proposed in the previous work [19]. The differences are the number of input/output dimensions of the network and the absence of the time series of distance inputs. The details of the architecture are shown in Fig. 5. We also use the Adam optimizer [20] and the epsilon-greedy method for exploration.



**Fig. 5:** (a) The overall architecture. The padding layer appends the first few repeats inputs to distances. (b) The *ConvUnit*. Any convolution layer has 50 hidden channels and strides (1, 1). (c) The *FCUnit*. Each fully-connected (FC) layer has 600 hidden channels.

## VII. EXPERIMENTS

In order to check the effectiveness of our method, we conducted experiments in both single-agent and multi-agent settings. Firstly, we show the common experimental settings for both cases. Then, we show some quantitative results for the diversity acquisition in the single-agent setting, where the surrounding vehicles ignore the ego vehicle. Finally, we show qualitative results to confirm that the generated policies can perform various kinds of interactions among multiple agents.

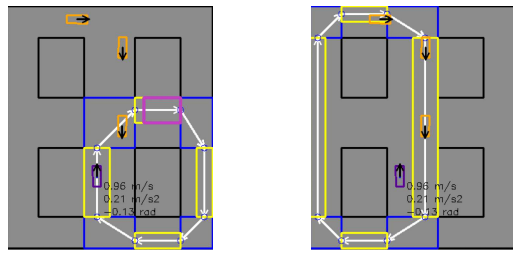
Although the following subsections discuss the right-turn traffic scene (Fig. 6), we also ran experiments on the straight-straight intersection scene and the 4-directional intersection scene. The results can be found in the supplemental video.

### A. Experimental Setting

Details of our experiments settings are listed in this subsection. Fig. 6 shows the navigation route of the target evaluated vehicle making a right turn (dark purple rectangle) and that of the other vehicles going straight down the opposite lane (orange rectangles). The mission of the ego vehicle is to reach to the goal area (purple rectangle in Fig. 6a) within 25 seconds without colliding with vehicles nor into walls. We set  $\Delta t = 0.1$  (thus, the maximum episode length is 250). The driving score of a policy  $\pi$  is defined as the proportion of successful scenarios  $S_\pi$ , where the agent successfully reaches its goal, to all scenarios  $S$ . In all settings, we set the driving score threshold as  $\delta = 90\%$ .

Weight coefficients and hyperparameters are as follows:

- $w^{(\text{move})} = 100$  in all experiments.
- $(w^{(\text{angle})}, w^{(\text{center})}) = (15, 5)$  in the multi-agent settings and  $(0, 0)$  otherwise.
- $w^{(\text{collision})}$  was linearly increased from 0 to 300 during the first 300,000 steps in all settings as in Miyashita et al. [19].
- $lr = 0.001$  in learning rate of Adam.
- $\epsilon$  of epsilon-greedy method was linearly decreased from 1.0 to 0.1 during the first 100,000 steps.
- $\alpha = 0.01$  which is the coefficient in DDE.



(a) Route for the ego vehicle. (b) Route for the other vehicles.

**Fig. 6:** An example of initial positions of vehicles and their navigation routes for the right-turn traffic.

**TABLE II:** Comparison of diversity scores for 50 policies.  $|\Pi_{\geq 90\%}|$  is the number of candidates policies with a success rate 90% or higher. Suc., O.A., and I.P. are the average success rate, overall diversity, and inter-policy diversity for selected 50 policies, respectively. Since *PolicySelect+DDE300* had only 25 filtered candidates, we selected all 25 policies instead.

Method	$ \Pi_{\geq 90\%} $	Suc. (%)	O.A. (ave.)	I.P. (ave.)
<i>RandomTrajectories</i>	N/A	N/A	1.49	3.44
<i>RandomSelect300</i>	79	93.20	3.31	1.32
<i>PolicySelect300</i>	79	93.40	3.12	1.66
<i>PolicySelect+DDE300</i>	25	93.20	2.84	2.11
<i>RandomSelect1200</i>	338	93.00	3.22	1.46
<i>PolicySelect1200</i>	338	94.32	2.72	2.24
<i>PolicySelect+DDE1200</i>	111	93.20	2.64	2.50
<i>RandomSelect14400</i>	3926	93.56	3.25	1.48
<i>PolicySelect14400</i>	3926	92.88	<b>2.39</b>	2.90
<i>PolicySelect+DDE14400</i>	1236	93.23	2.46	<b>3.18</b>

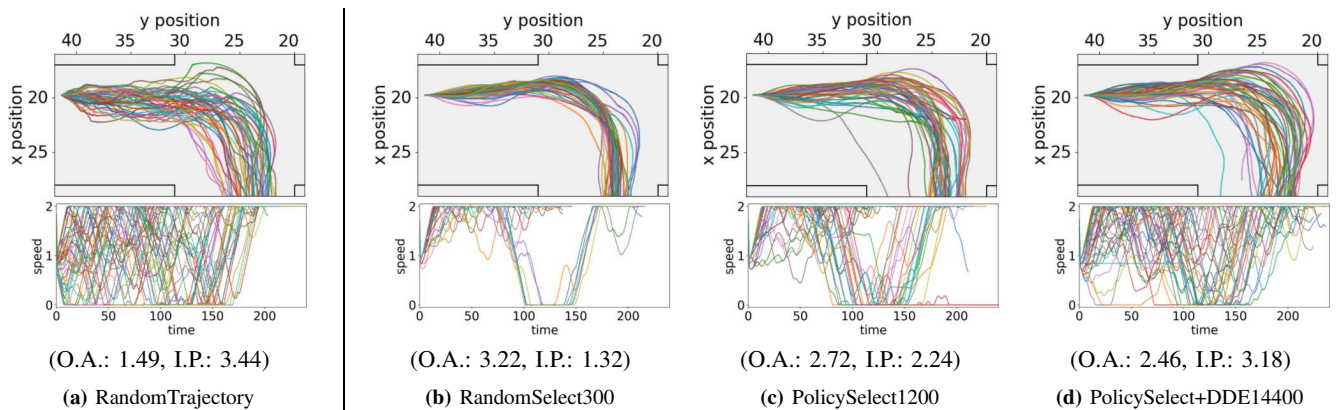
During training, the initial position of the cars is perturbed by a random number for each episode. In order to compare all methods under the same scenarios, we prepared 50 scenarios for evaluation, in which random numbers for the initial positions are fixed. To generate the candidates, we have run several training sessions in parallel. Each session trains an agent to create 150 snapshots of every 20,000 steps out of 3 million training steps. Consequently, the number of sessions for each candidate set is equal to the number of candidates divided by 150.

We did all the implementation in Python 3. We also used the implementation of Chainer [21] and ChainerRL [22] for the RL algorithm.

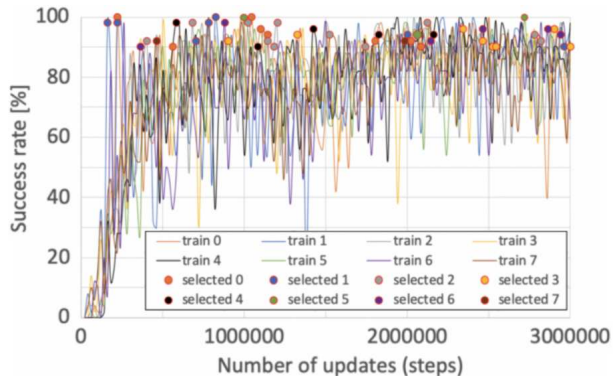
### B. Experiments for acquiring diversity

In this subsection, we focus on quantitatively evaluating how much diversity the proposed method can provide to a single agent. For this purpose, the behavior of peripheral agents other than the targeted agent is based on a fixed policy prepared by pre-training. Besides, the peripheral agents run in a mode where they cannot observe the evaluated target agent not to interact with it. Thus, the trajectories of other vehicles are fixed for each scenario.

The evaluation results are shown in Table II. As previously described, larger inter-policy diversity means better, while smaller overall diversity means better. As a reference, we added *RandomTrajectory*, a set of 50 trajectories generated with the Brownian bridge described in Section IV-B for each scenario. Here, the parameter  $\nu$  required for P-control is



**Fig. 7:** Visualizations of trajectories over 50 policies for a right-turn. The top represents positional changes (rotated clockwise by 90 degrees), and the bottom represents speed changes.



**Fig. 8:** Changes in the success rate of 8 training sessions for *PolicySelect1200*. The circle markers correspond to 50 policies selected in diverse policy selection. Differences in marker colors indicate that policies were generated from different training sessions.

set to 2.0 [s]. The names beginning with *RandomSelect*, *PolicySelect*, and *PolicySelect+DDE* are the results of extracted 50 policies by using the uniformly at random selector, our selector, and our selector whose candidates are trained with our modified DDE, respectively. Each number appended to each method name represents the number of candidate policies. Table II shows that, as the number of candidate policies increases, *PolicySelect* and *PolicySelect+DDE* get better scores than *RandomSelect* in both diversities. Although *PolicySelect14400* gets the best overall diversity, incorporating DDE gives a better score for the inter-policy diversity and competitive scores for the overall diversity whereas  $|\Pi_{\geq 90\%}|$  is smaller than ones without DDE. One possible reason is that the repulsive force induced by the intrinsic rewards helps to search spatially distinct trajectories, which leads to higher inter-policy diversity.

To show that our diversity measures correctly evaluate diversity, we also visualized the trajectories generated from trained 50 policies for a scenario in Fig. 7. There seems to be a positive correlation between better diversities and trajectory coverage. Especially, this shows that there is a great difference in the variety of speed changes.

We also visualized an example (*PolicySelect1200*) to

check the distribution of selected 50 policies generated from multiple training sessions. As shown in Fig. 8, the policy selector employed a variety of training processes. The success rate reached around 80% within 200,000 steps then kept fluctuating strongly until the training completed. One hypothesis is that the space of successful trajectories is multi-modal and the RL trainer finds various failures during the exploration. Therefore, the policies in different training steps have different strategies, and actually Fig. 7c indicates diverse behaviors. That would be one possible reason why using past snapshots helps diversity acquisition.

### C. Experiments with multi-agent settings

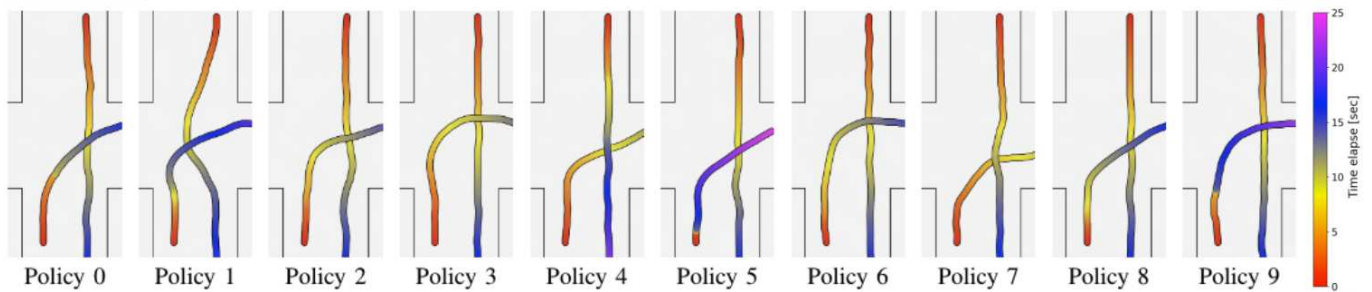
One last experiment we conducted with the purpose of checking if our RL-based policies can express a variety of interactive behaviors. In this experiment, each policy learns behaviors of both the right-turn and oncoming vehicles. Each vehicle is controlled in a decentralized manner by using a copy of the same policy. In training, we used transitions of one vehicle whose index is switched for each episode. We selected ten policies from 24,000 candidates by using our policy selector.

We examined the behavioral differences in the same scenario. As shown in Fig. 9, the selected policies display a variety of interactions:

- Policies 0, 1, 2, 5, 8, 9 seem to yield the other vehicle. The yielding position differs in each case. Policies 5, 9 had long stops, while Policy 2 decelerated gradually.
- Policies 3, 6 gave way to the upcoming vehicle by taking a large detour to the left side.
- Policies 4, 7 took no yielding behavior. Especially, Policy 4 made the upcoming vehicle wait for the right-turn vehicle.

## VIII. CONCLUSION

We presented a method for acquiring diverse driving policies that leverages RL’s exploration abilities and intrinsic rewards. Our experiments showed how we were able to acquire driving policies that optimize their behavioral diversity while maintaining good driving skills. Our proposed



**Fig. 9:** Behavioral differences for each policy in the same scenario. The color transition means the time change. Two points of the same color show the positions of two vehicles at the same time.

trajectory-based metrics to evaluate diversity, are also shown to be very helpful in building an efficient policy selection agnostic to how candidate policies are generated. We do expect that our approach could generate even more diverse policies by combining various policy generation methods. Such diverse driving modeling could also open new possibilities in building stronger in-car behavior prediction modules in the future. Although this work utilized driving skills as the filter, it would be possible to generate human-like behaviors by comparing them with real-world driving data. Since the parameters of reference trajectories are dependent on each simulation environment, an automatic parameter selector could help the diversity evaluation for other traffic scenes.

#### ACKNOWLEDGMENT

This work was supported by Toyota Research Institute - Advanced Development, Inc. We would like to thank Shin-ichi Maeda, Masanori Koyama, Mitsuru Kusumoto, Yi Ouyang, Daisuke Okanohara and anonymous reviewers for their helpful feedback and suggestions.

#### REFERENCES

- [1] P. Alvarez Lopez, M. Behrisch, L. Bieker-Walz, J. Erdmann, Y.-P. Flötteröd, R. Hilbrich, L. Lücken, J. Rummel, P. Wagner, and E. Wießner, “Microscopic traffic simulation using SUMO,” in *Proc. of ITSC*, 2018, pp. 2575–2582.
- [2] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, “Social gan: Socially acceptable trajectories with generative adversarial networks,” in *Proc of CVPR*, no. CONF, 2018.
- [3] N. Deo and M. M. Trivedi, “Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms,” in *Proc. of IV. IEEE*, 2018, pp. 1179–1184.
- [4] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone, “Trajectory++: Multi-agent generative trajectory forecasting with heterogeneous data for control,” *arXiv preprint arXiv:2001.03093*, 2020.
- [5] Y. Hu, A. Nakhaei, M. Tomizuka, and K. Fujimura, “Interaction-aware decision making with adaptive strategies under merging scenarios,” in *Proc. of IROS*, 2019, pp. 151–158.
- [6] T. F. Gonzalez, *Clustering to minimize the maximum intercluster distance*. Theoretical Computer Science, 1985, vol. 38, pp. 293–306.
- [7] Z.-W. Hong, T.-Y. Shann, S.-Y. Su, Y.-H. Chang, and C.-Y. Lee, “Diversity-driven exploration strategy for deep reinforcement learning,” in *Advances in Neural Information Processing Systems 31*, 2018, pp. 10 489–10 500.
- [8] D. Nishiyama, M. Y. Castro, S. Maruyama, S. Shiroshita, K. Hamzaoui, Y. Ouyang, G. Rosman, J. DeCastro, K.-H. Lee, and A. Gaidon, “Discovering avoidable planner failures of autonomous vehicles using counterfactual analysis in behaviorally diverse simulation,” in *Proc. of ITSC*, 2020, accepted.
- [9] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “CARLA: An open urban driving simulator,” in *Proc of CoRL*, 2017, pp. 1–16.
- [10] H. Zhao, A. Cui, S. A. Cullen, B. Paden, M. Laskey, and K. Goldberg, “Fluids: A first-order local urban intersection driving simulator,” in *Proc. of CASE*, 2018.
- [11] M. Naumann, F. Poggenhans, M. Lauer, and C. Stiller, “CoInCar-Sim: An open-source simulation framework for cooperatively interacting automobiles,” in *Proc. of IV*, Changshu, China, June 2018, pp. 1879–1884.
- [12] M. Althoff, M. Koschi, and S. Manzingler, “Commonroad: Composable benchmarks for motion planning on roads,” in *Proc. of IV*, 2017, pp. 719 – 726.
- [13] G. Bacchiani, D. Molinari, and M. Patander, “Microscopic traffic simulation by cooperative multi-agent deep reinforcement learning,” in *Proc. of AAMAS*, 2019.
- [14] Y. Liu, P. Ramachandran, Q. Liu, and P. Jian, “Stein variational policy gradient,” in *Proc. of UAI*, 2017.
- [15] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, “Diversity is all you need: Learning diverse skills without a reward function,” 2018, arXiv preprint arXiv:1802.06070.
- [16] A. Aubret, L. Matignon, and S. Hassas, “A survey on intrinsic motivation in reinforcement learning,” 2019, arXiv preprint arXiv:1908.06976.
- [17] C. Choi, A. Patil, and S. Malla, “DROGON: A causal reasoning framework for future trajectory forecast,” 2019, arXiv preprint arXiv:1908.00024.
- [18] P. Polack, F. Althché, B. d’Andréa Novel, and A. de La Fortelle, “The kinematic bicycle model: A consistent model for planning feasible trajectories for autonomous vehicles?” in *Proc. of IV*, 06 2017, pp. 812–818.
- [19] M. Miyashita, S. Maruyama, Y. Fujita, M. Kusumoto, T. Pfeiffer, E. Matsumoto, R. Okuta, and D. Okanohara, “Toward onboard control system for mobile robots via deep reinforcement learning,” in *Deep RL Workshop at Neurips*, 2018.
- [20] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proc. of ICLR*, 2015.
- [21] S. Tokui, R. Okuta, T. Akiba, Y. Niitani, T. Ogawa, S. Saito, S. Suzuki, K. Uenishi, B. Vogel, and H. Yamazaki Vincent, “Chainer: A deep learning framework for accelerating the research cycle,” in *Proc. of SIGKDD*, 2019, pp. 2002–2011.
- [22] Y. Fujita, T. Kataoka, P. Nagarajan, and T. Ishikawa, “Chainerrl: A deep reinforcement learning library,” in *Deep RL Workshop at Neurips*, December 2019.