

# Learning a Decision Module by Imitating Driver’s Control Behaviors

Junning Huang<sup>\*3</sup>, Sirui Xie<sup>\*4</sup>, Jiankai Sun<sup>1</sup>, Qiurui Ma<sup>5</sup>, Chunxiao Liu<sup>2</sup>, Dahua Lin<sup>1</sup>, Bolei Zhou<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>SenseTime Research

<sup>3</sup>Technische Universität Darmstadt

<sup>4</sup>University of California, Los Angeles

<sup>5</sup>Hong Kong University of Science and Technology

**Abstract:** Autonomous driving systems are usually considered as a pipeline of perception, decision, planning, and control. The decision module processes information from the perception module and directs the execution of downstream planning and control modules. Under the recent data-driven paradigm, end-to-end neural control policy has been proposed to the previous rule-based decision module. But one drawback is the difficulty of enforcement of physical or logical constraints to ensure safety and stability. In this work, we propose a hybrid framework for learning a decision module, agnostic to the mechanisms of perception, planning, and control modules. Instead of directly learning from human demonstration of decisions, our method learns high-level driving decisions by imitating drivers’ low-level control behaviors, a more reliable source of supervision. We show in the simulation experiments that our driving agent with a learned decision module can be generalized to various complex driving scenarios where the rule-based approach fails, while it generates smoother and safer driving trajectories than the end-to-end neural policy.

## 1 Introduction

An autonomous driving system is essentially a decision-making system that takes a stream of on-board sensory data as input, processes it with prior knowledge about driving scenarios to make a *decision*, and then outputs *control* signals to steer the vehicle [1]. It is often modularized into perception, decision, planning, and control. While the perception and its auxiliary, map reconstruction, tend to be more adaptive to machine learning methods [2], the downstream modules such as *decision*, *planning* and *control* remain as program-based to be configurable to ensure safe interaction with the physical world. In other words, this safety guarantee is built upon a human-inspectable and interruptible basis. However, in complex driving environments, building a complete rule-based system for decision making is infeasible. A corner case that is inconsistent with the existing rules may lead to a significant program refactorization. Aiming at a more flexible alternative, this work explores a learning-based decision-making module whose downstream modules, i.e. planning and control, remain encapsulated. Essentially, we draw inspiration from a prior work [3], which, in stark contrast to end-to-end control policies, reinforces to learn a decision module to switch between controllers designed with Lyapunov domain knowledge.

The primary challenge we have to address for a learning-based decision module is then the source of supervision. Here *driving decision* is defined as the high-level abstraction about what lane the driver wants the vehicle to be with what velocity in  $T$  seconds. This decision determines the execution of downstream planning and control modules. However, it is both difficult and ambiguous for human drivers to describe such a driving decision. For example, when different drivers try to merge into traffic at a roundabout, they may have different habitual behaviors to execute. Even the same individual may make different decisions in almost identical scenarios. To this end, there is hardly a set of criteria that are institutional or standardized for evaluating the quality of human-annotated driving

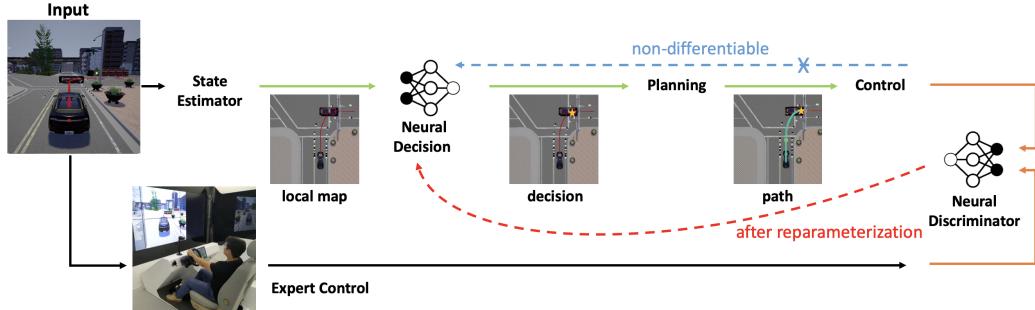


Figure 1: **Framework.** The green line shows the modular pipeline of the autonomous driving system. Solid orange lines indicate the offline training of the neural discriminator, while the dashed blue line shows that the planning and control modules are not differentiable. But the decision policy can be trained with reward assigned on control action by the discriminator after reparameterization, as indicated by the dashed red line.

decisions. In contrast, drivers’ physical behaviors, such as stepping on throttle/brake and steering the wheel, implicitly reflect their high-level decisions. Even though it is a weak supervision in the sense that the decision is not directly supervised, it is an easily accessible and technically reliable one. As an analogy, we can consider the behavioral data as in pixel space of images and the decision data as in the semantic space. Humans may have a different understanding on the same image. But it is easy to tell how different two images are, or how realistic a synthesized image is.

The setup of weak, indirect supervision invites a mechanism to infer human intentions from their behaviors and imitate these cognitive decisions. Intention understanding and behavior imitation have been extensively studied in the field of artificial intelligence and cognitive science, under the banner of Inverse Planning [4] and Inverse Reinforcement Learning [5]. But the former one mainly focuses on inferring the unobservable high-level intention by observing human behaviors from a third-person view [6]. And the latter one always adopts end-to-end neural policies and does not distinguish intentions from actions [7, 8]. Thus the indirect supervision in this work was barely discussed. Even though intentions are implicitly represented in the learned neural control policies, for physical systems like autonomous driving vehicles, our treatment of hierarchical modularization is appealing since end-to-end neural policies do not guarantee the fulfillment of low-level constraints in states or actions [9], nor do they have the robustness and stability guarantee as classical controllers [10]. Different from all these prior arts, given human driving behaviors, we learn a high-level decision module in the hierarchical pipeline where downstream modules follow the classical paradigm. This modular design requires a tailored learning method to handle the issue that the downstream programs are not differentiable for automated gradient calculators [11, 12].

In this work, we propose an imitation learning framework for decision making in the modular driving pipeline that *learns neural decisions from human behavioral demonstrations*. Although perception, planning, and control modules are included in the modeling, this learning framework is agnostic to their mechanisms. Policy learning is conducted in a generative adversarial manner [13]. The generator (green line in Fig. 1) simulates the modularized driving pipeline. With information from a local map, decisions are generated by a neural policy and processed by programmed planning and control modules to interact with the environment sequentially. The neural discriminator takes in the generated trajectories and compares them with drivers’ behavior data (orange lines in Fig. 1). Due to the absence of gradients from planning and control modules (blue line in Fig. 1), we derive a learning objective to propagate credits on control action back to the decision directly (red line in Fig. 1). We evaluate this framework on various simulated urban driving scenarios in CARLA, including following, merging at the crossing, merging at the roundabout, and overtaking, showing its superior performance over rule-based modular systems and end-to-end control policies.

We summarize our contributions as:

- The proposed framework combines the programming-based system with a generative learning method, where the high-level decision policy is data-driven while the low-level planning and control modules remain configurable and physically constrained. These low-level constraints also improve sample efficiency in interactive learning.

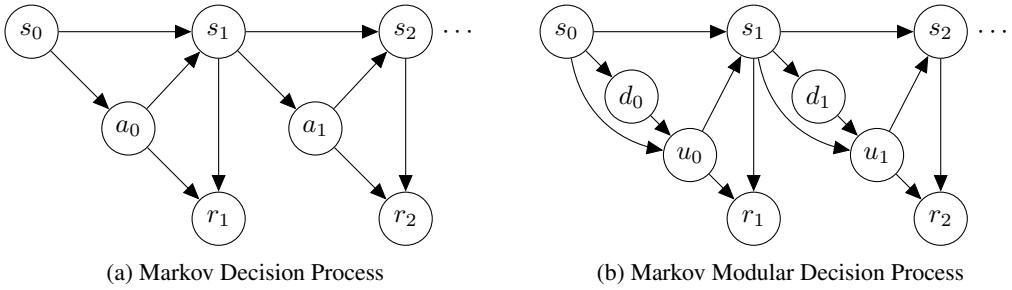


Figure 2: Comparison between *Markov Decision Process* and *Markov Modular Decision Process*

- The tailored generative adversarial learning method is weakly supervised and trained end-to-end. It learns high-level driving decisions from low-level control data, circumventing ambiguous annotation of human intentions.
- This framework successfully distills behaviors in different scenarios into one decision module. Empirical results demonstrate that the learned decision module fits into the pipeline, endowing the vehicle with smoother trajectories. It can also generalize to novel complex scenarios.

## 2 Related Work

**Classical autonomous driving system.** The design of the autonomous driving pipeline could be traced back to DARPA Urban Challenge 2007 [14]. The survey from [1] provides an overview of the hierarchical structure of perception, decision, planning, and control. Specifically, finite state machine (FSM) was used in the decision level, or equivalently upper planning level, in [15]. Classical autonomous driving systems are organized in this way such that it is friendly to testing and diagnosis. Based on this modular pipeline, our work aims at developing the data-driven decision module to replace the FSM.

**Imitation learning-based autonomous driving system.** With the popularity of deep generative learning, methods have been proposed to learn end-to-end neural control policies from driver’s data. Previously, Ziebart *et al.* [8] and Ross *et al.* [7] proposed general methods in Inverse Reinforcement Learning and Interactive Learning from Demonstration, with an empirical study on a driving game. More recently, Kuefler *et al.* [16] and Behbahani *et al.* [17] learn an end-to-end policy in a GAIL[18]-like manner. Codevilla *et al.* [19] and Liang *et al.* [20] share similar hierarchical perspective as us, but their control policies are completely neural. One thing worth mentioning is Codevilla *et al.* [19] proposed to imitatively learn a downstream policy given human’s high-level command, whose motivation is supplementary to ours. Rhinehart *et al.* [21] propose to learn a generative predictor for vehicles’ coordinate sequence. Li *et al.* [22] imitate the mapping from predicted waypoints to control commands via neural network policy. Our work differentiates from all of them as we imitate the driver’s low-level behaviors to learn high-level decisions, whose execution is conducted by a transparent and configurable program.

**Learning human intention.** Baker *et al.* [4] formalized human’s mental model of planning with Markov Decision Process (MDP) and proposed to learn human’s mental states with Bayesian inverse planning. A similar probabilistic model of decisions is adopted by us but the learning is done with a novel generative adversarial method. Wang *et al.* [6] proposed a Hidden Markov Model (HMM) for human behaviors and materialized it with Gaussian Process (GP), where the intent is the hidden variable. Jain *et al.* [23] model the interaction between human intentional maneuvers and their behaviors in a driving scenario, both in-car and outside, as an Auto-Regressive Input-Output HMM and learn it with EM algorithm. However, they didn’t explicitly distinguish humans from the environment. Besides, rather than inferring the human’s mental state, we focus more on how a robot, *i.e.* can make good decisions.

### 3 Modular Driving System

#### 3.1 System Overview

Fig 1 shows an overview of our proposed learning framework. We focus on the decision-making module and its connection to the downstream planning and control modules. The *state estimator* (usually a perception module) constructs a local map by processing the sensory data and combining them with prior knowledge from the map and rules. The *decision* module then receives *observation* and decides a legal local driving task that steers the car towards the destination. To complete this task, the *planning* module searches for an optimal trajectory under enforced physical constraints. The *control* module then reactively corrects any errors in the execution of the planned motion.

In our proposed framework, the *decision* module is a conditional probability distribution parameterized by neural networks. The downstream *planning* and *control* module are encapsulated. Because they are not the main contribution of this work, we adopt minimal settings to showcase our learning frameworks. These encapsulated modules could be replaced by arbitrary alternatives.

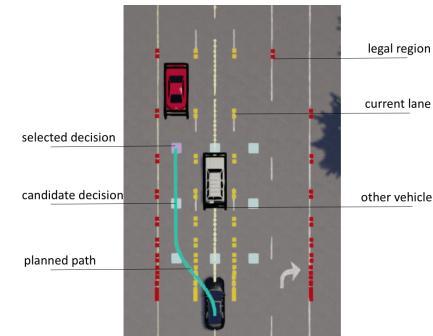
#### 3.2 The Decision Making Module

We assume that a local map centered at the ego vehicle could be constructed from either the state estimator or an offline global HD map. This local map contains critical driving information such as the routing guidance to the destination, legal lanes, lane lines, other vehicle coordinates and velocity in the ego vehicle's surroundings, as well as ego vehicle's speed and acceleration.

The interaction between the decision module and the environment through its downstream is modeled as Markov Decision Process (MDP). MDP is normally represented as a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{D}, \mathcal{P}, r, \rho_0, \gamma)$ , with a state set  $\mathcal{S}$ , a decision set  $\mathcal{D}$ , a transitional probability distribution  $\mathcal{P} : \mathcal{S} \times \mathcal{D} \times \mathcal{S} \rightarrow \mathbf{R}$ , a bounded reward function  $r : \mathcal{S} \times \mathcal{D} \rightarrow \mathbf{R}$ , an initial state distribution  $\rho_0 : \mathcal{S} \rightarrow \mathbf{R}$ , a discount factor  $\gamma \in [0, 1]$  for infinite horizon. Decision policy  $\pi_\theta : \mathcal{S} \times \mathcal{D} \rightarrow \mathbf{R}$  takes current state  $s_t \in \mathcal{S}$  from local map and generates a high-level decision  $d_t$ . Note that the notation of  $\mathcal{D}$  and  $d$  are unconventional, we explicitly denote *decision* with  $d_t$  to differentiate it from *control action*  $u_t$ . This decision directs the execution of planning and control module, and makes the autonomous driving system proceed in the environment to acquire next state  $s_{t+1}$ . Therefore there is a modification of the decision process, as shown in Fig.2. The optimization objective of this policy is to maximize the expected discounted return  $\mathbf{E}_\tau[\sum_{t=0}^T \gamma^t \mathcal{R}(s_t, u_t)]$ , where  $\tau = (s_0, d_0, u_0, \dots)$  denotes the whole trajectory,  $s_0 \sim \rho_0(s_0)$ ,  $d_t \sim \pi_\theta(d_t|s_t)$ ,  $u_t \sim P(u_t|s_t, d_t)$  and  $s_{t+1} \sim \mathcal{P}(s_{t+1}|s_t, u_t)$ .

Below we introduce the *observation* space and *decision* space, which are the interfaces in our modular system. Note that we assume full access to necessary information on the local map.

To make the decision-making policy more generalizable to different driving scenarios, the output of the state estimator, also as the input to the neural decision module, is the local map and traffic state. As shown in Fig. 3, a part of the observation is the 3D coordinates of sample points of lane lines. We select two sets of lanes for the decision policy. One for the ego vehicle's current lane and the other for the edge of *legal* regions. Note that here *legal* means driving in that region abides by both the traffic rules and a global routing, which would have a crucial effect at the crossing. To reduce the dimension of input, these lane points are sampled with exponentially increasing slots towards the further end. In some sense, it mimics the effect of Lidar, with the observation more accurate in the near end. The observation also includes the coordinates and velocities of the 6 nearest vehi-



**Figure 3: Illustration of input and output of the decision module when the driving system conducts overtaking.** The trained policy decides to pick the left lane, in terms of the target point and the target speed. The planner generates a trajectory, which is tracked by the controller. Observations are shown on the right. Yellow lines are the ego vehicle's current lane lines and the red lines are the edge of *legal* regions.

cles in traffic within 70m range of the ego vehicle.

We define the *decision* as three independent categorical variables, one for lateral, one for longitudinal and one for velocity. Combining with the local map, each of them is assigned a specific semantic meaning. The lateral decision variable has three classes as *changing to the left lane*, *keeping current lane* and *changing to the right lane*. Note that at a crossing, global routing information has been implanted into the local map as introduced in the last paragraph. The lateral decision set is complete since these three are the only possible lateral decisions for a vehicle. The longitudinal decision variable has four different classes, with each indicating *the distance along with waypoints in time interval T*. The exact coordinate of the endpoint could be extracted from the local map. Combining with the predicted target speed at the endpoint from the velocity decision variable, which discretizes the allowed speed range equally into four baskets, this decision module provides a quantitative configuration of goal states for the downstream trajectory planner.

### 3.3 Planning and Control Modules

In this subsection we introduce a minimal design of the planning module and the control module. These are not major contributions of this work and are introduced for the sake of completeness. In classical autonomous driving system, the planning module calculates the trajectories under constraints such as collision avoidance. Its optimization objective is to reach the goal state specified by the decision module with minimal cost. In our minimal setting, the planning module processes path and velocity separately. Paths are planned with cubic Bezier curves, while velocity is planned with Quadratic Programming (QP) to minimize the acceleration, jerk, and discontinuity in curvature. The control module drives the vehicle to trace the planned trajectories. And we implement it minimally with PID. Details are provided in Appx. A.

It is flexible to further extend each module of the proposed framework, as long as the alternatives are deterministic or stationary. Note that while the interface between decision and planning is the location of a goal point and the speed ego vehicle is expected to maintain when reaching it, the planning module could be replaced by more sophisticated search-based or optimization-based trajectory planner, to fulfill more practical requirements in execution time and constraint enforcement. Similarly, the model-free controller could be replaced with alternatives like Model Predictive Control (MPC). Other practical constraints and concerns over the controller such as robustness and stability could also be taken into consideration if necessary.

## 4 Learning Method

### 4.1 From IRL to GAIL

Under the MDP modeling, we choose imitation learning over reinforcement learning. In principle, imitation learning could save practitioners from *ad hoc* reward engineering and focus their effort on reusable infrastructure development. Among all the imitation learning methods, behavior cloning (BC) [24] is the most straightforward one. The policy is trained as a regressor in a purely supervised manner, ignoring the temporal effect of each action in an MDP trajectory. Thus it suffers from covariate shift when presented with states which are not covered by the training data. By contrast, imitation learning such as Inverse Reinforcement Learning (IRL) explicitly models the interaction between policy and environment and approximates for exact value iteration [8] or Monte Carlo approximation [18]. Thus they can generalize better in states not covered in the demonstration.

Fig.2 illustrates the difference in the graphical model of a trajectory from an MDP when the modular system is adopted for the agent. In this modified probabilistic model, a trajectory is a sequence  $\tau = (s_0, d_0, u_0, s_1, d_1, u_1 \dots)$ . However, what is observable from demonstration is  $\hat{\tau} = (s_0, u_0, s_1, u_1 \dots)$ . According to the formulation of MaxEnt IRL [8], when maximum entropy principle is applied for the tie breaking between cost functions that can generate identical optimal trajectories, the probability of a trajectory is

$$p(\hat{\tau}|R) = \frac{1}{Z} \exp(R(\hat{\tau})), \quad (1)$$

where  $Z$  is the partition function,  $R$  is the reward function, which is normally Markovian  $R(\hat{\tau}) = \sum r(s_i, u_i)$ . Given a set of demonstration  $\mathcal{D}_E$ , we can infer the reward function by Maximum

Likelihood (MLE):

$$\mathcal{L}_{IRL} = \mathbf{E}_{\mathcal{D}_E}[\log p(\hat{\tau}|R)] = \mathbf{E}_{\mathcal{D}_E}[\sum_{\hat{\tau}} r(s_i, u_i) - \log Z]. \quad (2)$$

In normal IRL, with the learnt reward function, we can solve for the policy with any reinforcement learning methods. However, as the partition function is always intractable if the state space is continuous, methods have been proposed to approximate it [8, 25]. Basically, there is a background distribution  $q(\hat{\tau})$  for the estimation of  $Z$ . In [26], it is further illustrated how a delicately designed importance sampling scheme can connect this IRL loss to a specific type of discriminator over trajectories. Ho *et al.* [18] propose Generative Adversarial Imitation Learning (GAIL) to further approximate this discriminator with a step-wise one, remedy the trajectory-wise information with advantage accumulation and learn a policy with off-the-shelf reinforcement learning methods. Intuitively, it learns a policy whose actions at given states are indistinguishable from demonstration data. More formally, with expert demonstration  $\mathcal{D}_E$ , a neural discriminator  $D_\phi : \mathcal{S} \times \mathcal{U} \rightarrow (0, 1)$  is introduced and the reward function is defined as  $r(s, u) = -\log(D_\phi(s, u))$ . The training objective is:

$$\mathcal{L}_{GAIL} = \min_{\theta} \max_{\phi} \{ \mathbf{E}_{p_\theta(u|s)}[\log(D_\phi(s, u))] + \mathbf{E}_{\mathcal{D}_E}[\log(1 - D_\phi(s^E, u^E))] \}. \quad (3)$$

$D_\phi$  is optimized with cross entropy loss, while  $\pi_\theta$  is optimized with policy gradient [27].

## 4.2 Handling Non-differentiable Downstream Modules

Different from GAIL, in our framework the generation process is modularized, following the structure of classical autonomous driving systems. Because the downstream planning and controls modules are not necessarily differentiable, a separation occurs between decisions from neural policy and the control data from drivers, as illustrated by the blue dotted line in Fig 1: policy  $\pi_\theta$  generates *decisions*  $d_t$  while discriminator only distinguishes *actions*  $u_t$ . And as shown in Fig 2,  $p_\theta(u_t|s_t) = \sum_{d_t} p(u_t|d_t)\pi_\theta(d_t|s_t)$ . Since planning and control module are both deterministic, at every state, once the decision is chosen, the transformation from decision to control is a deterministic and correspondence mapping. Hence, one important insight of our work is that the transformation at each state  $s$

$$g_s : \{d_t\} \xrightarrow[1\text{-to-1}]{\text{onto}} \{u_t\}, \quad (4)$$

which is the abstraction of the planning and control module, is a deterministic and correspondence mapping. Therefore, it is a candidate for *reparametrization* [28] or *push-forward* [29]. Reparametrize  $u$  in the first expectation in Eq.3 with  $d$ :

$$\begin{aligned} \mathbf{E}_{p_\theta(u|s)}[\log(D_\phi(s, u))] &= \sum p_\theta(u|s) \log(D_\phi(s, u)) \\ &= \sum \pi_\theta(d|s) \log(D_\phi(s, g_s(d))) = \mathbf{E}_{\pi_\theta(d|s)}[\log(D_\phi(s, g_s(d)))]. \end{aligned} \quad (5)$$

Intuitively, we only use this black-box function  $g_s$  for Monte Carlo sampling, thus there is no need to know its analytical form. In the training stage, we do not actually need to calculate  $u_t$  with the push-forward function as in Eq. 4. Instead, we calculate the true  $g_s(d_t)$  according to Sec. 3 during sampling, and save both  $d_t$  and  $u_t$  for training. During training,  $d_t$  is fed to  $\pi_\theta$  and the corresponding  $u_t$  is fed to  $D_\phi$ . Since  $\mathcal{D}$  is discrete, it is differentiated with policy gradient. We then have this modified learning objective:

$$\mathcal{L} = \min_{\theta} \max_{\phi} \{ \mathbf{E}_{\pi_\theta(d|s)}[\log(D_\phi(s, u))] + \mathbf{E}_{\mathcal{D}_E}[\log(1 - D_\phi(s^E, u^E))] \}. \quad (6)$$

The whole algorithm is described in Algorithm 1.



Figure 4: Trajectories from our proposed framework (green lines) are smoother than the ones from the end-to-end policy (orange lines), because of the explicit enforcement of geometrical constraints in the downstream planning module. ( $\times$  indicates the car violates routing rules.)

---

**Algorithm 1** Learning executable decisions from human behavioral data

---

**Require:** Expert demonstration  $\mathcal{D}_E$ , batch size  $B$   
 Initialize policy  $\pi_\theta$ , discriminator  $D_\phi$ , (optionally)  
 value function  $V_w$   
**while** policy iteration not converged **do**  
 Initialize data buffer  $\mathcal{B}$  with length  $B$   
**while**  $\neg\mathcal{B}.full()$  **do**  
 Sample expert trajectory:  $\tau_E \sim \mathcal{D}_E$   
 Initialize  $sim$  with  $s_0^E$   
 $traj = \text{GENERATE-TRAJ}$   
 Append  $traj$  to  $\mathcal{B}$   
**end while**  
 Compute reward for samples in  $\mathcal{B}$  with  $D_\phi$   
 Update  $D_\phi$  with data in  $\mathcal{B}$  and  $\mathcal{D}_E$   
 Policy iteration on  $\pi_\theta$  with tuples in  $\mathcal{B}$ , (optionally with value iteration on  $V_w$ )  
**end while**

---

**Algorithm 2** GENERATE-TRAJ

---

**Require:** Simulator (or ROS offline replayer) with local map  $sim$ , planning module  $planner$ , control module  $controller$ , maximum trajectory horizon  $\mathcal{H}$   
 Initialize placeholder  $traj$ ,  $t = 0$   
**while**  $t < \mathcal{H}$  or  $\neg sim.done$  **do**  
 Get  $s_t$  from  $sim$   
 Sample  $d_t$  from  $\pi_\theta(o_t)$   
 Get  $pln_t$  from  $planner(sim, d_t)$   
 Get  $u_t$  from  $controller(sim, pln_t)$   
 Append  $(s_t, d_t, u_t)$  to  $traj$   
 Send  $u_t$  to  $sim$   
 $t++$   
**end while**  
 Yield  $traj$

---

## 5 Experiments

### 5.1 Rudimentary Driving in Empty Town

Experiments in *Empty Town* show the difference between a modular driving system and an end-to-end control policy. The training time on an 8-GPU cluster for the end-to-end control policy is nearly 31 hours versus 15 hours for our modular driving system. Different training time in the same platform implies different numbers of interactions these two methods need to converge, showing the improvement in exploration from constraint enforcement.

As shown in Fig. 4, trajectories from our proposed framework are smoother than the ones from the end-to-end policy. This demonstrates the effect of explicit enforcement of geometrical constraints in the downstream planning module. Interestingly, there are some road structures where the learning-based modular system can pass while the end-to-end control policy cannot. For example, Fig. 4(d) shows a narrow sharp turn where end-to-end control policy drives across the legal lane line. We believe this is a difficult temporally-consistent exploration problem. An agent needs some successful trials of consecutive right-turning action samples to learn, which is challenging for generative interactive learning. The modular system, in contrast, only explores behaviors that are temporally plausible like human drivers. With the planning module that searches for a geometrically constrained path, the agent can effortlessly make this sharp turn.

Table 1 shows the statistics of behaviors from the learning-based modular system, end-to-end control policy, and rule-based method. The comparison shows that our method can drive more safely (0% collision rate), much faster (less time to finish), and with higher comfort (less acceleration). This exhibits the advantage of having both the learned decision policy and the classic planning and control module. Besides, the decision module trained in Town 3 works similarly well in Town 2, showing its generalization ability.

## 5.2 Socially Interactive Driving in Traffic Scenarios

In this section, we test our model’s performance in each traffic scenario. The statistics of Traffic Scenarios in Table 1 are the average of 100 evaluations. Traffic scenarios include basic scenarios such as *Car Following* and *Crossroad Merge* where vehicle’s interactions with zombies are relatively simple and monotonic; and complex ones such as *Crossroad Turn Left*, *Roundabout Merge* and *Overtake*. In complex traffic scenarios, the dynamics around the ego vehicle is of higher variation, due to either more traffic users or more complex relations between them.

As shown in Table 1, in basic traffic scenarios, rule-based and learning-based modular systems are somehow on par in terms of time taken to finish, while rule-based agent seems to drive a little bit more comfortably. End-to-end agent drives more rudely in most of the scenarios, except in *Crossroad Merge*, where it takes much longer time to finish. Since the modular pipeline enforces smoother planning and control, the learning-based modular system offers more driving comfort than the end-to-end system. In basic traffic scenarios, none of them drives as well as experts.

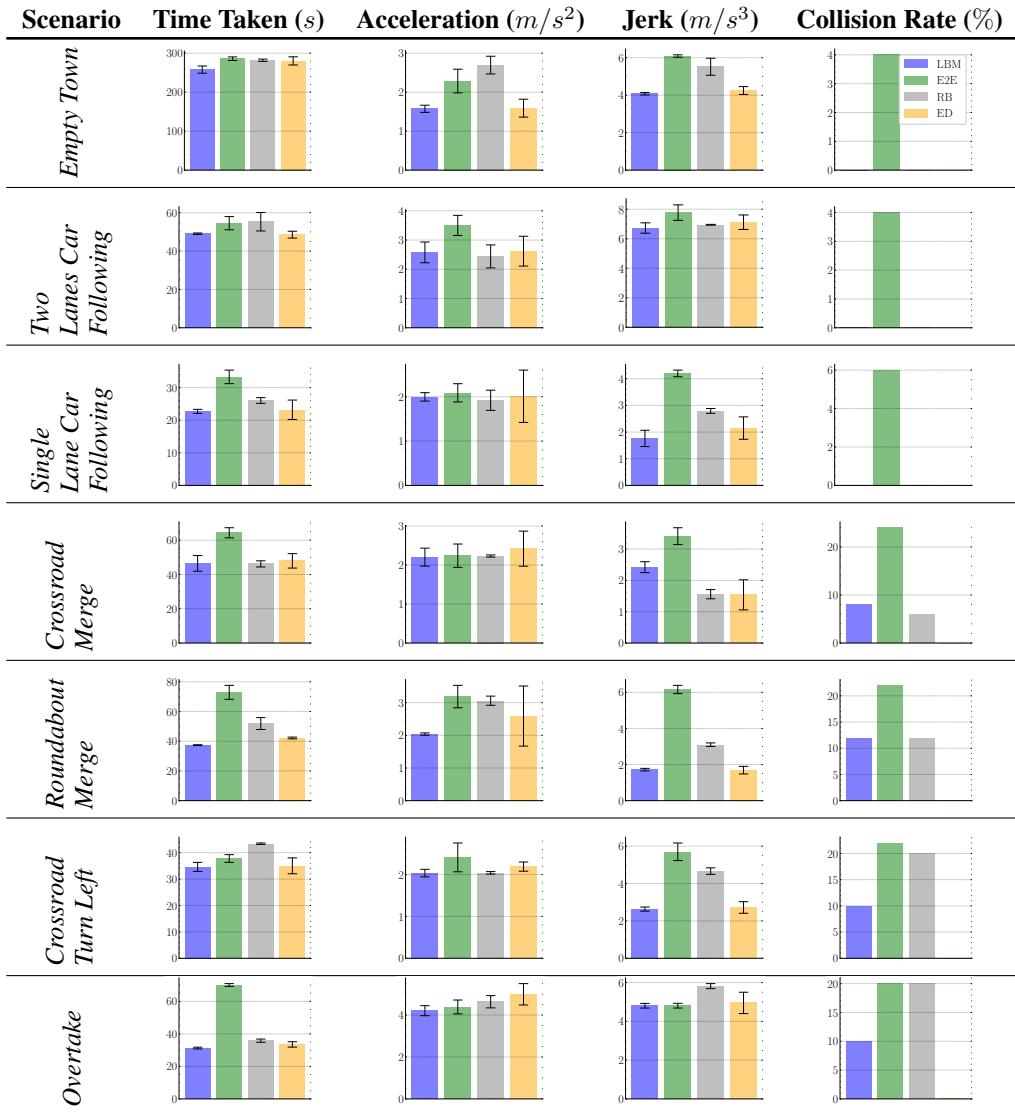
In complex traffic scenarios, rule-based agent configured in some scenarios fails to drive safely in others while learning-based modular agent is safer and smarter. This is because rule-based decision module is based on FSM, which is composed of a huge number of rules. In order to drive safely in all environments, the rules of safety have priority over other factors such as driving comfort. In complex environments, the rules in FSM would be more conservative to handle corner cases and uncertainty(*e.g.*, more hard brakes or accelerations). Unlike the rule-based system, the learning-based system can minimize the acceleration and jerk by imitating human experiences. Furthermore, since the rule-based system is program-based, its configurations might vary in different scenarios, which makes it difficult to generalize across scenarios. When compared with end-to-end control policy, a similar conclusion could be drawn as in the previous subsection that learning-based modular agent offers more comfort because the modular pipeline enforces smoother planning and control. Interestingly, in complex social scenarios, the learning-based modular system achieves higher comfort (lower acceleration) than experts. This may be attributed to human errors, and a learning-based modular agent somehow fixes it. Experiments of generalization capability can be found in Appendix.

## 6 Conclusion and Future Work

This work introduces a flexible framework for learning decision making for autonomous driving. To learn the driver’s high-level driving decision, the proposed framework imitates driver’s control behavior with a modularized generation pipeline. Being agnostic to the design of the non-differentiable downstream planning and control modules, this framework can train the neural decision module through reparametrized generative adversarial learning. We evaluate its effectiveness in simulation environments with human driving demonstrations.

This work can be extended to more complex environments where other vehicles are also learning agents, *i.e.* multi-agent learning system. Alternatively, if we can collect data in the real-world and replay them in the simulator, this framework could be tested in real driving environments.

Table 1: Performance Comparision (LBM: Learning-based Module, E2E: End-to-end neural policy, RB: Rule-based method, ED: Expert Data)



## References

- [1] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli. A survey of motion planning and control techniques for self-driving urban vehicles. *IEEE Transactions on intelligent vehicles*, 1(1):33–55, 2016.
- [2] S. Thrun, W. Burgard, and D. Fox. *Probabilistic robotics*. MIT press, 2005.
- [3] T. J. Perkins and A. G. Barto. Lyapunov design for safe reinforcement learning. *Journal of Machine Learning Research*, 3(Dec):803–832, 2002.
- [4] C. L. Baker, R. Saxe, and J. B. Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.
- [5] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- [6] Z. Wang, K. Mülling, M. P. Deisenroth, H. Ben Amor, D. Vogt, B. Schölkopf, and J. Peters. Probabilistic movement modeling for intention inference in human–robot interaction. *The International Journal of Robotics Research*, 32(7):841–858, 2013.
- [7] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635, 2011.
- [8] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. 2008.
- [9] J. F. Fisac, A. K. Akametalu, M. N. Zeilinger, S. Kaynama, J. Gillula, and C. J. Tomlin. A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions on Automatic Control*, 2018.
- [10] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel. Adversarial attacks on neural network policies. *arXiv preprint arXiv:1702.02284*, 2017.
- [11] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>. Software available from tensorflow.org.
- [12] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [14] M. Buehler, K. Iagnemma, and S. Singh. *The DARPA urban challenge: autonomous vehicles in city traffic*, volume 56. Springer, 2009.
- [15] M. Montemerlo, J. Becker, S. Bhat, H. Dahlkamp, D. Dolgov, S. Ettinger, D. Haehnel, T. Hilden, G. Hoffmann, B. Huhnke, et al. Junior: The stanford entry in the urban challenge. *Journal of field Robotics*, 25(9):569–597, 2008.
- [16] A. Kuefler, J. Morton, T. Wheeler, and M. Kochenderfer. Imitating driver behavior with generative adversarial networks. In *2017 IEEE Intelligent Vehicles Symposium (IV)*, pages 204–211. IEEE, 2017.
- [17] F. Behbahani, K. Shiarlis, X. Chen, V. Kurin, S. Kasewa, C. Stirbu, J. Gomes, S. Paul, F. A. Oliehoek, J. Messias, et al. Learning from demonstration in the wild. *arXiv preprint arXiv:1811.03516*, 2018.

- [18] J. Ho and S. Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, pages 4565–4573, 2016.
- [19] F. Codevilla, M. Miiller, A. López, V. Koltun, and A. Dosovitskiy. End-to-end driving via conditional imitation learning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–9. IEEE, 2018.
- [20] X. Liang, T. Wang, L. Yang, and E. Xing. Cirl: Controllable imitative reinforcement learning for vision-based self-driving. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 584–599, 2018.
- [21] N. Rhinehart, R. McAllister, and S. Levine. Deep imitative models for flexible inference, planning, and control. *arXiv preprint arXiv:1810.06544*, 2018.
- [22] G. Li, M. Mueller, V. Casser, N. Smith, D. L. Michels, and B. Ghanem. Oil: Observational imitation learning. *arXiv preprint arXiv:1803.01129*, 2018.
- [23] A. Jain, H. S. Koppula, B. Raghavan, S. Soh, and A. Saxena. Car that knows before you do: Anticipating maneuvers via learning temporal driving models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3182–3190, 2015.
- [24] D. A. Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural Computation*, 3(1):88–97, 1991.
- [25] C. Finn, S. Levine, and P. Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pages 49–58, 2016.
- [26] C. Finn, P. Christiano, P. Abbeel, and S. Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint arXiv:1611.03852*, 2016.
- [27] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- [28] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [29] N. Rhinehart, K. M. Kitani, and P. Vernaza. R2p2: A reparameterized pushforward policy for diverse, precise generative path forecasting. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 772–788, 2018.
- [30] J. Choi, R. Curry, and G. Elkaim. Path planning based on b  zier curve for autonomous ground vehicles. In *Advances in Electrical and Electronics Engineering - IAENG Special Edition of the World Congress on Engineering and Computer Science 2008*, pages 158–166, 2008.
- [31] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017.
- [32] S. Shalev-Shwartz, S. Shammah, and A. Shashua. On a formal model of safe and scalable self-driving cars. *arXiv preprint arXiv:1708.06374*, 2017.
- [33] S. Bansal, V. Tolani, S. Gupta, J. Malik, and C. Tomlin. Combining optimal control and learning for visual navigation in novel environments. *arXiv preprint arXiv:1903.02531*, 2019.

## Appendices

### A Planning and Control Modules

The planning module processes path and velocity separately. Paths are planned with cubic Bezier curve [30]:

$$B(t) = (1-t)^3 P_s + 3(1-t)^2 t P_1 + 3(1-t)t^2 P_2 + t^3 P_g, \quad t \in [0, 1] \quad (7)$$

where  $P_s$ ,  $P_g$  are the starting and goal point,  $P_1$  and  $P_2$  are intermediate control points, imposing the geometrical constraints such as end pose aligning with lane direction.

The velocity planning module reparameterizes Eq.(7) with arc-length, introducing temporal unit. This reparameterization is not done directly, rather,  $n$  spline segments with identical time unit  $\Delta T$  are used to fit sample points from (7). Each of the segments is defined as:

$$s_i = a_i + b_i t + c_i t^2 + d_i t^3 + e_i t^4, \quad t \in [0, \Delta T]. \quad (8)$$

Here we consider the constraint on velocity planning, by minimizing the acceleration and jerk while fitting to the sample points from (7) and maintaining the continuity of curvature. Formally, the fitting is done by Quadratic Programming (QP):

$$\begin{aligned} & \min \sum_{i=1}^n \left[ \int_0^{\Delta T} \ddot{s}_i^2 dt + \int_0^{\Delta T} \dot{\ddot{s}}_i^2 dt \right], \\ & \text{s.t. } \dot{s}_0(0) = v_{now}, \ddot{s}_0(0) = a_{now}, \dot{s}_n(\Delta T) = v_{goal} \\ & \quad \dot{s}_k(\Delta T) = \dot{s}_{k+1}(0), \ddot{s}_k(\Delta T) = \ddot{s}_{k+1}(0), \text{ for } k = 0, 1 \dots n-1 \end{aligned} \quad (9)$$

The control module is a basic PID controller. Here to reconcile with the MDP assumption in decision, we adopt the discrete time form:

$$u(t) = K_p e(t_k) + K_i \sum_{i=0}^k e(t_k) \Delta t + K_d \frac{e(t_k) - e(t_{k-1})}{\Delta t}, \quad (10)$$

where  $K_p$ ,  $K_i$  and  $K_d$  denote the coefficients for the proportional, integral, and derivation terms respectively.  $e(t)$  is the error function. In this work we have two independent error functions  $e_{lat}(t)$  and  $e_{lon}(t)$  for lateral controller steering wheel  $u_{steer}(t)$  and longitudinal controller throttle/brake  $u_v(t)$ , which are both capped by mechanical constraints.

### B Training Details

#### B.1 Experimental Setup

**Training System** We evaluate the proposed framework in CARLA simulation platform [31]. As the proposed framework focuses on the decision and its downstream, the state estimator is neglected and the observations are extracted directly from the simulator. Since the most recent version of CARLA is not efficient enough for large-scale simulation and learning, we build a parallel training system. The parallel training system shown in Fig. 5 consists of two parts. The first part (denoted as a blue dotted line box named Rollout) collects the transition data and the second part updates the model's parameters with the transition data. Every iteration, in Rollout, an actor synchronizes with the Policy node and obtains the model's parameters. Then the actor retrieves the transition data via the CARLA frontend. After that, the transition data would be stored in the Trajectory node. Here the transition data is generated by the interaction between the CARLA frontend and the CARLA backend in a Docker container. After the collection of the transition data finished, we start to update our model's parameters. The Discriminator and Policy nodes would update their parameters by following Algorithm 1 and using the data in the Trajectory and Demo nodes. To improve the robustness of the system, a Daemon reboot node is injected to reboot the Docker container when the system hangs.

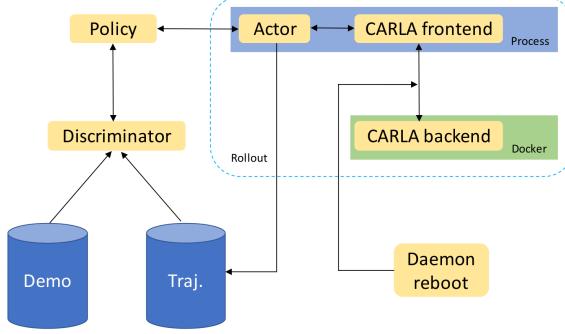


Figure 5: **Parallel Training System.** The first part (denoted as a blue dotted line box named Rollout) collects the transition data and the second part updates the model’s parameters with the transition data.



Figure 6: **Autonomous Driving Simulation Platform**, which is used to collect the low-level drivers’ control behavior. We learn high-level driving decisions by imitating the low-level drivers’ control commands.

**Scenarios Setup** Multiple common driving scenarios of varying complexity are constructed for evaluation. All driving scenarios are deployed in the Map Town 3 except Overtake. Since Town 3 provides only two-lane roads and Overtake should be demonstrated on a three-lane road, we choose Town 4 to design Overtake. To increase the complexity of the environment, the speed of the zombie car is set with randomness, and we give its average speed and design standard deviation.

- *Empty Town* is a special traffic scenario where there is no traffic. It covers urban road structures like a straight road, curve road, crossing, and roundabout.
- *Car Following* (Fig. 7a). There are two sub-scenarios in this scenario: *Two Lanes Car Following* (Fig. 7a) and *Single Lane Car Following*. *Two Lanes Car Following* is designed based on a four-lane dual carriageway. Two zombie cars are running side by side in two lanes, and the ego-car follows them in one of the lanes. The speed of the zombie cars is  $18 \pm 1$  km/h. In *Single Lane Car Following*, only one zombie car is running ahead in a single lane, and the ego-car follows the zombie car. The speed of the zombie car is  $15 \pm 1$  km/h.
- *Crossroad Merge* (Fig. 7b). In the lateral lane, no less than three zombie cars as a traffic stream drive through the intersection in series, and ego-car merges into their traffic flow from a lane perpendicular to the lateral lane. The speed of zombie cars is randomly chosen from  $21 \pm 1$  km/h,  $25 \pm 1$  km/h. There are two sub-scenarios.
- *Crossroad Turn Left* (Fig. 7c). Multiple zombie cars drive in the lateral lane and the opposite lane. The ego-car needs to cross the intersection and then head to the left lane, during which it needs to avoid collision with the zombie cars and obey the traffic rules. The speed of zombie cars is randomly chosen from  $26 \pm 1$  km/h,  $42 \pm 1$  km/h. There are two sub-scenarios.
- *Roundabout Merge* (Fig. 7d). Four zombie cars are entering the roundabout one by one, bypassing 270 degrees counterclockwise and then driving out of the roundabout. The ego-car enters the roundabout and joins the traffic flow at about 90 degrees. After following the zombie cars around the roundabout for about 180 degrees, the ego-car drives out with them together. The speed of zombie cars is randomly chosen from  $21 \pm 1$  km/h,  $25 \pm 1$  km/h,  $35 \pm 1$  km/h. There are three sub-scenarios.
- *Overtake* (Fig. 3). Ego Car drives in the middle lane (medium speed lane) of the three lanes with Zombie Car I in front of it. Zombie Car II (different from Zombie Car I), which is on the left front of Zombie Car I, drives in the high-speed lane. Under the premise of obeying traffic rules, Ego Car has to go beyond Zombie Car I in front of it by changing to the high-speed lane. Due to the blocking of Zombie Car II, Ego Car will change back to the medium speed lane after exceeding Zombie Car I, then completing the entire overtaking process.

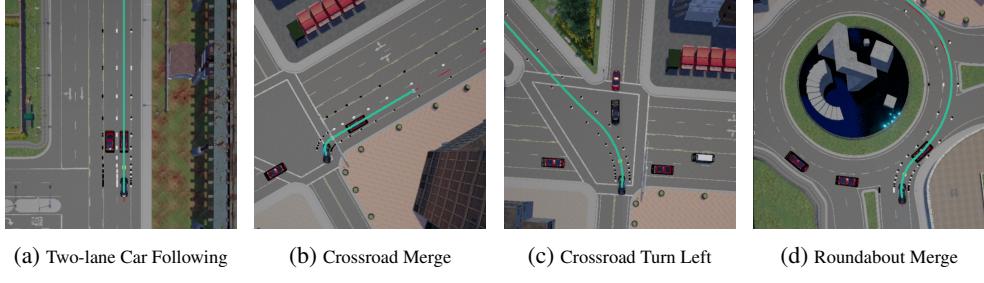


Figure 7: Screenshots of four scenarios.

The maneuver of zombies in all scenarios are defined by two items: configurations of zombies, and a CARLA built-in waypoint-following controller. Configurations of zombies for experiments in Table 1 is the footnote part in Table 2.

Besides the car velocity variance mentioned above, the evaluation environments for each scenario also vary in seeds and vehicle models. For example, Audi TT, Chevrolet, Dodge (the police car), Etron, Lincoln, Mustang, Tesla 3S provided by CARLA.

Table 2: Configurations of Scenarios

Scenarios	Number of Cars	Starting Velocity of Ego Car (km/h)	Starting Velocity of Zombie Cars (km/h)	Target Velocity of Zombie Cars (km/h)	Starting Direction (°)	Expected Direction (°)
<i>Single Lane Car Following</i>	2	25.20	10.00	$15.00 \pm 1.00^{\ddagger}$	0.00	180.00
<i>Two Lanes Car Following</i>	3	32.40	11.00	$18.00 \pm 1.00^{\ddagger}$	0.00	0.00
<i>Crossroad Merge</i>	> 4	25.20	10.00	$21.00 \pm 1.00^{\ddagger}$ $25.00 \pm 1.00$	0.00	90.00
<i>Roundabout Merge</i>	> 4	25.20	10.00	21.00 ± 1.00, 25.00 ± 1.00, 35.00 ± 1.00 <sup>‡</sup>	0.00	0.00
<i>Crossroad Turn Left</i>	> 7	25.20	10.00	26.00 ± 1.00, $42.00 \pm 1.00^{\ddagger}$	0.00	-42.00
<i>Overtake</i>	3	33.48	11.00	25.00 ± 1.00 <sup>‡</sup> (Zombie Car I), $40.00 \pm 1.00^{\ddagger}$ (Zombie Car II)	0.00	0.00

<sup>‡</sup>The footnote part is the speed configuration for experimental results in Table 1 and Table 4.

**Data Collection** The demonstration data is collected by human driving ego car in an autonomous driving simulation platform. The simulation platform shown in Fig. 6 includes CARLA simulator with Logitech G29 Driving Force Steering Wheels & Pedals. The PC configuration is Ubuntu 16.04 x86, Intel Core i7-8700 CPU, GeForce GTX 1060 and 16 GB memory. The following criteria are used to set to filter out low-quality demonstrations:

- There is no collision;
- No dangerous moves and obeying traffic rules.

For each scenario, experts repeatedly drive in the simulator to finish episodes with 200 steps. Trajectories violating the aforementioned criteria are rejected. After collecting 1000 trajectories in each scenario, we each choose 100 trajectories that meet the above criteria as a demonstration. To

improve the generalization ability of the learned policy, stochasticity was added into the expert’s steering and throttle during the collecting procedure. Namely, a random configuration is sampled at the beginning of each episode, the expert would choose different strategies to finish the current scenario. Aggressive strategies may include rash steering and throttle.

The statistic results of expert data are listed as Table 3. Note that the ‘speed’ column in the Table 3 is not the expert data statistics but only used to distinguish sub-environments. For example  $26.00 \pm 1.00$  indicates the expert data collected under the *Crossroad Turn Left* ( $26.00 \pm 1.00$ ) sub-environment. Statistical indicators include Collision Rate, Time Taken, Accel. and Jerk.

Table 3: Statistics of Expert Data

Scenarios	Speed (km/h)	Collision Rate (%)	Time Taken (s)	Accel. ( $m/s^2$ )	Jerk ( $m/s^3$ )
<i>Empty Town</i>	/	0.00	$280.21 \pm 10.48$	$1.59 \pm 0.23$	$4.25 \pm 0.21$
<i>Single Lane Car Following</i>	$15.00 \pm 1.00$	0.00	$23.21 \pm 2.98$	$2.01 \pm 0.59$	$2.15 \pm 0.42$
<i>Two Lanes Car Following</i>	$18.00 \pm 1.00$	0.00	$48.56 \pm 1.79$	$2.62 \pm 0.51$	$7.12 \pm 0.49$
<i>Crossroad Merge</i>	$21.00 \pm 1.00$	0.00	$47.95 \pm 4.21$	$2.42 \pm 0.45$	$1.54 \pm 0.48$
	$25.00 \pm 1.00$	0.00	$48.87 \pm 1.21$	$2.63 \pm 0.43$	$1.69 \pm 0.19$
	$21.00 \pm 1.00$	0.00	$36.23 \pm 3.69$	$2.11 \pm 0.19$	$1.69 \pm 0.60$
<i>Roundabout Merge</i>	$25.00 \pm 1.00$	0.00	$34.91 \pm 0.34$	$2.01 \pm 0.32$	$1.65 \pm 0.19$
	$35.00 \pm 1.00$	0.00	$42.17 \pm 0.59$	$2.59 \pm 0.92$	$1.69 \pm 0.21$
<i>Crossroad Turn Left</i>	$26.00 \pm 1.00$	0.00	$34.91 \pm 3.33$	$1.93 \pm 0.21$	$2.45 \pm 0.21$
	$42.00 \pm 1.00$	0.00	$35.01 \pm 3.01$	$2.19 \pm 0.11$	$2.72 \pm 0.31$
<i>Overtake</i>	$25.00 \pm 1.00,$ $40.00 \pm 1.00$	0.00	$33.52 \pm 1.65$	$4.99 \pm 0.51$	$4.95 \pm 0.55$

**Evaluation Metrics** We compare the learned decision module with both the rule-based module and end-to-end neural control policy. Rule-based module is composed following [32]. The end-to-end policy is trained with vanilla GAIL. Following [33], we conduct quantitative comparison including metrics like collision rate, time to accomplish tasks, average acceleration, and jerk after 100 trials in evaluation environments. These metrics reflect how safe and smooth the agent drives. The trajectory from different approaches is also visualized and compared qualitatively. We provide demo videos in the supplementary.

## B.2 Baseline Methods

**Rule-based Method** The rule-based system in our paper is based on the formal model provided in the paper of MobileEye [32]. This paper proposed the definitions of the safe longitudinal and lateral distance and designed a car-following model. For example, the vehicle should keep the safe longitudinal distance to the front car and avoid collision with the rear car. Since the paper didn’t mention the overtake scenario, we divided the overtake scenario into three stages and design rules with the instructions of the formal model. In the first stage, the car-following stage, the vehicle would make decisions instructed by the car-following model. In the second stage, the turning-left stage, the agent would decide when to turn left according to the safe longitudinal and lateral distance, then turning-left with high speed. In the last stage, the turning-right stage, similar to the turning-left decision, the agent would turn right back at a reasonable speed.

**End-to-end Neural Policy** The end-to-end method used in our paper is GAIL [18]. Nothing has changed except that the state space and action space have been adapted to the CARLA environment.

## B.3 Other Training Details

**Reparameterization** The planning and control module are both deterministic. Hence at every state, once the decision is chosen, the transformation from decision to control is a deterministic and

correspondence mapping. The reparameterization trick can be seen in Equation 5. We reparametrize  $u$  in the first expectation with  $d$ . In other words, the conditional probability of  $d$  is the same as  $u$ .

**Decision Space** Human decisions can be divided into three parts: which lane to go, how far to go, at which speed to go. We model the three parts as three categorical variables to simplify its representations. To make a tractable inversion of the collected data, a 10 meters interval is used in the longitudinal goal while a 10km/h interval in the target speed. The interval is based on the feedback from our human drivers.

**Policy Architecture and Hyperparameters** The policy architecture of the decision module in our method is exactly the same as the architecture in the paper of GAIL [18]. For a fair comparison, all hyperparameters are also the same. Learning-based Module method and End-to-end Neural Policy baseline are trained for 500 iterations before evaluation while Rule-based method baseline does not require training and can be evaluated directly. The maximum time step for each iteration is 1000. We use a two-layer ReLU network with 32 units for the discriminator of GAIL and our method. For the policy, we use a two-layer (32 units) ReLU neural network. Entropy regularizer weight is set to be 0.1 across our method and GAIL. We use a batch size of 512 steps per update. The learning rate for generator and discriminator is 0.0003.

## C More Experiment Results

### C.1 Policy Distillation (Generalization Capability)

The disadvantage of end-to-end learning is not just that there are no guarantees, but also that generalization is tricky. Generalization capability is crucial for our method. There are three experiments that provide evidence for the generalization capability of our method.

- Different towns: The Empty town model trained in Town3 can be applied in Town2 and Town1 (cf. attached video).
- Different scenarios: The model trained in *Crossroad Merge* also performs well in *Roundabout Merge* and *Single Lane Car Following* (cf. *One Policy for Various Scenarios*).
- Different velocities: The *Crossroad Merge* model trained with the velocity of 25km/h also performs well at 15km/h (cf. *One Policy for Various Speeds*).

The aforementioned experiments train one specific model for each scenario separately. To test the practical potential of the proposed framework, we train one policy with demonstration data for multiple scenarios and multiple speeds. This policy distillation problem is divided into: I) *One Policy for Various Scenarios* and II) *One Policy for Various Speeds*.

For *One Policy for Various Scenarios*, the policy distillation process runs as follows: We leverage the same demonstration data as the previous experiments in four scenarios (*Single Lane Car Following*, *Crossroad Merge*, *Roundabout Merge* and *Crossroad Turn Left*) together to train one policy for these scenarios jointly. Four CARLA simulators run in parallel, each corresponding to one scenario, while only one learner is deployed to learn a policy from all data sampled simultaneously. Then we can obtain one policy model for all scenarios.

For *One Policy for Various Speeds*, we leverage the same demonstration data for each scenario. However, the demonstration data for each scenario contains various speeds (as Table 2 shows). We train a policy for each scenario using corresponding demonstration data and environments with various speeds, while only one learner is deployed to learn one policy from all data sampled from different speeds simultaneously. Finally, one policy is learned for various speeds.

Table 4 shows the statistics of behaviors of the learning-based modular system for various scenarios. After the training is finished, the learned policy is evaluated on the evaluation environment of different scenarios respectively (100 trials each). We also evaluated the learned policy on the other two unseen environments and reported the average results of 100 trials in Table 4. Table 5 6 7 shows the statistics of behaviors of the learning-based modular system for various speeds. The learned policy is evaluated on the evaluation environment of various speeds respectively (100 trials each). The distilled policy shows a slight performance loss compared to the policies previously trained separately in the case of one specific scenario. But the performance degradation is within the acceptable

Table 4: Generalized to Various Scenarios

Scenarios	Learning-based Module			
	Collision Rate (%)	Time Taken (s)	Accel. ( $m/s^2$ )	Jerk ( $m/s^3$ )
<i>Single Lane Car Following</i>	2.00	$24.35 \pm 1.01$	$2.41 \pm 0.56$	$1.92 \pm 0.22$
<i>Crossroad Merge</i>	8.00	$52.83 \pm 0.61$	$2.92 \pm 0.21$	$5.31 \pm 0.33$
<i>Roundabout Merge</i>	14.00	$40.99 \pm 1.26$	$2.21 \pm 0.67$	$1.82 \pm 0.32$
<i>Crossroad Turn Left</i>	14.00	$35.21 \pm 2.29$	$4.81 \pm 0.95$	$5.22 \pm 0.68$
<i>Two Lane Car Following (Unseen)</i>	6.00	$50.35 \pm 4.81$	$2.81 \pm 0.61$	$7.92 \pm 0.53$
<i>Overtake (Unseen)</i>	18.00	$35.62 \pm 1.16$	$5.63 \pm 0.55$	$5.43 \pm 0.72$

Table 5: Generalized to Various Speeds (Crossroad Merge)

Scenarios	Speed ( $km/h$ )	Learning-based Module			
		Collision Rate (%)	Time Taken (s)	Accel. ( $m/s^2$ )	Jerk ( $m/s^3$ )
<i>Crossroad Merge</i>	$21.00 \pm 1.00$	10.00	$51.82 \pm 4.32$	$2.99 \pm 0.44$	$5.52 \pm 0.79$
	$25.00 \pm 1.00$	12.00	$47.83 \pm 5.21$	$2.83 \pm 0.35$	$5.69 \pm 0.58$

Table 6: Generalized to Various Speeds (Roundabout Merge)

Scenarios	Speed ( $km/h$ )	Learning-based Module			
		Collision Rate (%)	Time Taken (s)	Accel. ( $m/s^2$ )	Jerk ( $m/s^3$ )
<i>Roundabout Merge</i>	$21.00 \pm 1.00$	14.00	$41.52 \pm 5.22$	$2.63 \pm 0.31$	$1.72 \pm 0.35$
	$25.00 \pm 1.00$	12.00	$39.56 \pm 6.21$	$2.09 \pm 0.39$	$1.70 \pm 0.41$
	$35.00 \pm 1.00$ (Unseen)	14.00	$37.15 \pm 5.59$	$2.01 \pm 0.49$	$1.72 \pm 0.42$

Table 7: Generalized to Various Speeds (Crossroad Turn Left)

Scenarios	Speed ( $km/h$ )	Learning-based Module			
		Collision Rate (%)	Time Taken (s)	Accel. ( $m/s^2$ )	Jerk ( $m/s^3$ )
<i>Crossroad Turn Left</i>	$26.00 \pm 1.00$	14.00	$36.49 \pm 3.21$	$2.15 \pm 0.29$	$2.99 \pm 0.32$
	$42.00 \pm 1.00$	16.00	$38.53 \pm 3.18$	$2.09 \pm 0.81$	$2.72 \pm 0.69$

range. The above policy distillation results (*One Policy for Various Scenarios* and *One Policy for Various Speeds*) exhibits the potential that our framework can handle scenarios in different levels of complexity concurrently without compromising performance.

## C.2 Model Predictive Control as Controller

Our focus is on assembling learning-based module with programming-based modules in an autonomous driving system, and handle non-differentiable downstream modules via the reparameterization trick. We adopt PID controller and Bezier curve because they are easy to implement and widely used in autonomous driving systems. They can be replaced with any other controllers (*e.g.*, MPC) or planning algorithms. In Table 8, we demonstrate the quantitative results of our method with MPC as a controller. It can be seen that our method with MPC as a controller still outperforms the baselines (Rule-based method and End-to-end neural policy). Compared with Table 1 which uses PID as a controller, the performance of our method is improved because MPC controller is a more advanced method of process control.

Table 8: Performance Comparision [MPC Controller] (LBM: Learning-based Module, E2E: End-to-end neural policy, RB: Rule-based method, ED: Expert Data)

