

# Inverse Reinforcement Learning with Hybrid Weight Tuning and Trust Region Optimization for Autonomous Maneuvering

Yu Shen<sup>1</sup>, Weizi Li<sup>2</sup>, and Ming C. Lin<sup>1</sup>

<https://gamma.umd.edu/researchdirections/autonomousdriving/eirl/>

**Abstract**—Despite significant advancements, collision-free navigation in autonomous driving is still challenging. The navigation module needs to balance the use of machine learning and motion planning to achieve efficient and effective control of the vehicle. We propose a novel framework combining context-aware multi-sensor perception and inverse reinforcement learning with hybrid weight tuning (IRL-HWT) for autonomous maneuver. IRL-HWT incorporates several attributes including non-uniform prior for features, hybrid weight tuning based on trust-region optimization, parameter reuse for continuous training, and learning from accidents. Our method can incorporate both expert demonstration and domain knowledge (reward function) to train a good policy utilizing the non-uniform prior and the trust-region optimization. All these attributes help reduce the number of collisions of the vehicle up to 41%, increase the training efficiency by 2.5x, and obtain higher test scores up to two orders of magnitude. Overall, our method can enable the vehicle to drive 10x further than other methods, while achieving collision avoidance over both static and dynamic obstacles.

## I. INTRODUCTION

Autonomous vehicles (AVs) have the potential to contribute to a more efficient and safer transportation system by alleviating traffic congestion and reducing the number of accidents. In general, autonomous driving can be realized via either an end-to-end approach or a mediated-perception approach [1]. An end-to-end approach meaning that raw sensor data are directly mapped to control commands (e.g., steering angles), which usually results in a succinct training pipeline at the cost of model interpretability. A mediated-perception approach, on the other hand, decouples perception and navigation, thus offering better model interpretability and enhanced driving safety. However, a mediated-perception approach commonly adopts a planning algorithm for navigating an AV—a process that can be computationally expensive, given the requirement of holistic environment information for achieving global optimality and planning in high-dimensional state space.

We propose a novel mediated-perception approach for autonomous driving. Our approach consists of context-aware multi-sensor perception and inverse reinforcement learning with hybrid weight tuning mechanism (IRL-HWT). The perception module interprets unstructured information (i.e., images and point clouds) of an environment using multiple sensors with different viewpoints, extract context-aware

information, and produces structured information, such as the size and position of an object in the environment. IRL-HWT then takes the structured information along with expert trajectories as input to learn a control policy for autonomous driving. In particular, IRL-HWT addresses a fundamental limitation of the original IRL, which imposes a *uniform prior* on all features. This limitation can lead to undesired behaviors (e.g., frequent collisions) even when the feature expectation of the learned policy closely matches the feature expectation of the expert policy (see Sec. IV for a concrete example). IRL-HWT solves this problem by *incorporating a non-uniform prior for features and update the weights of features using an automatic hybrid weight tuning mechanism that is based on trust-region optimization*.

A key element for a successful reinforcement learning (RL) task is to design the reward function, which is learned through expert data in inverse reinforcement learning (IRL). Our method can incorporate both expert demonstrations and domain knowledge to train an effective policy. In many domains, we can find “experts” easily (e.g., a driver) and some domain requirements are clear (e.g., “reaching the goal” and “collision avoidance” for the driving task). We provide a *formulation that can take any reward function as an initial guess; then IRL-HWT will update the reward function within a trust region*.

Our approach offers several advantages: (1) As a mediated-perception approach, instead of using a planning algorithm for navigating an AV at all time steps, we only use a planning algorithm to generate expert demonstrations for IRL-HWT to imitate. Once learned, the control policy can operate in real time with a small number of features. This design choice greatly reduces the computational overhead of using a planning algorithm. (2) IRL-HWT complements the original IRL [2] with the flexibility to impose a non-uniform prior on important features for a specific problem, for example, collision for autonomous driving. (3) The hybrid weight tuning scheme allows the learned policy to adapt to different environments while maintaining desired behaviors and avoiding aggressive exploratory behaviors (which is essential for safety-critical applications such as driving). (4) Since our approach is based on reinforcement learning, compared to using supervised learning for imitation, IRL-HWT is more robust in rare situations and can generalize better in new environments [2]. (5) Similar to ADAPS [3], alternative safe trajectories are generated during the analysis of an accident. This approach enables the learning algorithm to *learn from accidents*, which is crucial as collecting accident

<sup>1</sup>Yu Shen and Ming C. Lin are with the Department of Computer Science, University of Maryland at College Park {yushen, lin}@cs.umd.edu

<sup>2</sup>Weizi Li is with the Department of Computer Science, University of Memphis wli@memphis.edu

data from the real world is impractical.

The effectiveness and efficiency of our approach are demonstrated in a variety of experiments. To show IRL-HWT can work beyond perfect perception, i.e., use ground-truth data from the environment, we run IRL-HWT in a complete autonomous system as described in Sec. III-A. Overall, our method can enable the AV to drive safely *10x further* than existing SOTA methods. The attribute non-uniform prior can assist in reducing the number of collisions of the AV up to 41%; the use of learned model parameters for continuous training can result in 2.5x faster training; and learning from accidents can help achieve higher test scores up to two orders of magnitude. In addition to quantitative results, we qualitatively show that our method can steer the AV to avoid both static and dynamic obstacles, even in the presence of a narrow passage (see supplementary video).

## II. RELATED WORK

In this section, we briefly discuss some related work.

### A. Autonomous Driving

Various methods have been proposed to solve the perception, planning, and control problems for autonomous driving [4], [5]. Examples of the end-to-end approach include end-to-end reinforcement learning [6] and end-to-end imitation learning [3], [7], [8], [9]. These approaches usually require a large amount of training data in order to be robust in rare cases, such as pre-accident scenarios. In addition, the use of deep neural networks in these approaches to directly map raw sensor data to control commands can lead to low model interpretability.

Examples of the mediated-perception approach include perception plus motion planning [10] and perception plus learning-based planning [7]. Because of the decomposition of perception and navigation, these approaches enjoy better model interpretability, hence improved driving safety. Recently, Li et al. propose ADAPS [3], an end-to-end imitation learning framework that enables an AV to learn from accidents. Compared to ADAPS, our work proposes a new architecture that can not only fuse data from different modalities but also generalize better, as it relies on reinforcement learning rather than using supervised learning to imitate the expert's behaviors.

### B. Inverse Reinforcement Learning (IRL)

As an effective technique for imitation learning, IRL involves two steps: 1) learning a reward function from experts' demonstrations and 2) using the acquired reward function for reinforcement learning to learn a control policy [2]. To provide some examples, Sharifzadeh et al. [11] apply Deep Q-Networks to extract a reward function in large state space. You et al. [12] use deep neural networks to approximate the latent reward function of the expert and then apply deep Q-learning to obtain the control policy.

Maximum entropy IRL [13] provides a probabilistic approach based on the entropy concept. Hierarchical guidance [14] leverages the hierarchical structure of the underlying problem to integrate different modes of expert

interactions. Guided Cost Learning [15] uses a nonlinear cost function and guided sampling strategy. GAIL [16] draws the analogy between imitation learning and generative adversarial networks and develops an effective model-free imitation learning for complex, high-dimensional tasks. Maximum entropy deep IRL [17] combines neural networks and traditional IRL. Lastly, a recent study [18] proposes a scalable IRL algorithm based on an adversarial reward learning process.

Compared to other IRL-based methods, our approach incorporates both expert demonstrations and domain knowledge to design an effective initial reward function for obtaining a good policy. In addition, our approach uses the parameters of a pre-trained policy as the start point for continuous training, which can improve training efficiency while maintaining the model performance.

## III. APPROACH

### A. Framework Overview

Our architecture combines context-aware multi-sensor perception and inverse reinforcement learning with hybrid weight tuning (IRL-HWT). The perception module takes multiple sensors' data as input and produces structured data, which are then used by IRL-HWT to learn a control policy.

We assume that the AV can obtain a global map from an external service, and compute its position and rotation with on-board GPS and Inertial Measurement Unit (IMU). We further assume that the AV knows beforehand a few sparse waypoints on its path to the goal, and the task of the AV is converted to reach the waypoints consecutively.

Our perception module can produce semantic-rich structured data to learn a reward function by utilizing the context-aware semantic information. Features with clear semantic interpretations can help construct non-linear features such as "whether there is a car in front within 3 meters", resulting in more flexible decision-making. This is particularly useful considering that IRL restricts the reward function to be a linear combination of the features.

The IRL-HWT training process, shown in Fig. 1, consists of an offline step and an online step. In the offline step, we use a planning algorithm as the expert to generate driving trajectories and compute the expert's feature expectation. In the online step, we learn the feature weights under a *non-uniform prior* given the expert's policy and the learned policy at the current iteration. Then, we construct a reward function using the learned feature weights. By further adopting the notion of *learning from accidents* [3], we use the resulting (additional) training data along with the newly constructed reward function for RL to update the learned policy. Inspired by transfer learning [19], the model parameters from the previous iteration are used as a starting point in the updating process.

### B. Context-aware Multi-sensor 3D Perception

We simulate three RGB cameras for the front-view, left-view, and right-view, respectively, as well as a 360-degree Lidar of the AV. We combine one RGB image and the point

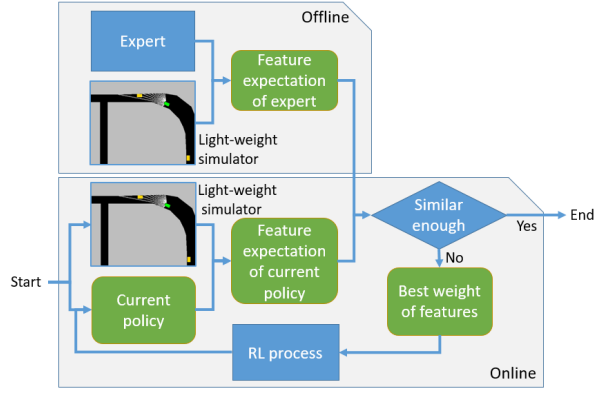


Fig. 1. IRL-HWT training process. We compute the expert’s trajectories and feature expectation offline. We then use the learned reward function to compute the feature expectation of a learned policy online. If the feature expectations of the expert and the learned policy differ too much, RL is used to update the learned policy.

cloud on each side to detect nearby vehicles. Then, we merge the results from all sides to obtain an overall view. The module works as follows.

First, it generates the front-view data, which contain segmentation and depth maps, and bird’s-eye-view image from the point cloud. By having the calibration data between the Lidar and camera, we can then align the front-view data with the RGB images. Next, we combine and feed the aligned front-view data with the RGB images into the feature extractor to obtain the front-view feature map. We apply the same procedure for the bird’s-eye-view image to obtain the bird’s-eye-view feature map. Then, we use the region proposal network (RPN) [20] and detection network from AVOD [21] to obtain the perception results for the front-view. Similarly, we obtain the perception results for the left-view and right-view. Finally, we merge all perception results together using the extrinsic calibration data of the three cameras.

The output of the perception module contains 3D bounding boxes and types of nearby dynamic obstacles. Our approach can also be extended to detect static obstacles when needed. The context-aware, multi-sensor 3D perception detects different types of information to produce features used by IRL.

### C. Inverse RL with Hybrid Weight Tuning (IRL-HWT)

The original IRL achieves imitation learning by first computing the expert’s feature expectation  $\hat{\mu}(\pi_E) = \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{\infty} \gamma^t \phi(s_t^{(i)})$ , given  $m$  trajectories  $\{s_0^{(i)}, s_1^{(i)}, \dots\}_{i=1}^m$  from the expert’s policy  $\pi_E$ , the discount factor  $\gamma$ , and the feature vector  $\phi(\cdot)$ . Then, IRL learns  $w$  ( $w \in \mathbb{R}^k$  and  $\|w\|_1 \leq 1$ ) given  $\pi$  at the current iteration and  $\pi_E$ , and synthesizes a reward function  $R(s) = w \cdot \phi(s)$ , where  $s$  is the state of the environment. Next, the policy  $\pi$  is re-learned using RL. This iterative process continues until  $\|\mu(\pi) - \hat{\mu}(\pi_E)\|_2 \leq \epsilon$ . The final policy is selected among all learned policies from all iterations.

The features used by IRL-HWT are based on structured data from the perception module, which include bounding boxes of nearby vehicles and static obstacles in the scene. See detailed feature list in Sec. IV.

One fundamental limitation of IRL is that it imposes a *uniform prior* on all features, causing small weights to be possibly assigned to some crucial features during the learning process. For example, in the context of driving, we find that the feature *collision* can receive a small weight as a result of any collision behavior of the AV will terminate a training episode. This limitation of IRL can lead to subpar performance of the task *learning to drive* (see Sec. IV-B). To give an example, the expert can drive the car safely (without any collision), while a randomly initialized policy can hardly drive the car far without collision. In this case, the feature expectation of the expert is calculated using extended, “global” trajectories, while the feature expectation of the learned policy is calculated using short, “local” trajectories. This discrepancy is likely to cause some important features to receive small weights during the minimization process of the feature expectations in IRL. A concrete example of this phenomenon from our experiments is shown in Table I.

TABLE I

LIMITATION OF IRL. IRL LEARNS FEATURE WEIGHTS  $w$  BY MINIMIZING THE DIFFERENCE OF FEATURE EXPECTATION BETWEEN THE EXPERT’S POLICY  $\mu(\pi_E)$  AND A RANDOM POLICY  $\mu(\pi_0)$ . WHILE BOTH POLICIES HAVE SMALL EXPECTATIONS FOR THE FEATURE *collision*, THE ACTUAL TRAJECTORY FROM THE EXPERT’S POLICY CAN BE MUCH LONGER THAN THE TRAJECTORY FROM A RANDOM POLICY. AS A RESULT, IRL ASSIGNS A NEGLIGIBLE WEIGHT TO *collision* (I.E.,  $-5.34e-06$ ).

feature	left dist.	front dist.	right dist.	...	<b>collision</b>
weight	0.113	0.184	-0.124	...	-0.00000534
$\mu(\pi_E)$	139.14	369.50	25.84	...	0
$\mu(\pi_0)$	40.31	216.06	142.87	...	0.005

To alleviate the aforementioned limitation of IRL, we propose IRL-HWT—an approach that not only incorporates a *non-uniform prior* on the features by allowing users to specify the weights of certain features for ensuring essential properties of a task, e.g., collision for driving, but also uses a hybrid weight tuning to update the policy so that the AV can adapt to different environments while maintaining desired behaviors. Formally, the overall weights  $w = [w_m, w_l]$  consist of empirically-initialized weights  $w_m = [w_1, \dots, w_k]$  and the weights to be learned from scratch  $w_l = [w_{k+1}, w_{k+2}, \dots, w_n]$ , where  $n$  is the total number of features. To achieve an optimal policy in diverse environments, we use trust-region optimization [22] to automatically tune  $w_m$  by avoiding violent exploratory behaviors. Specifically, in the  $i$ th iteration, we first try to update  $\epsilon$  and  $w$  by solving the following quasi-convex optimization program:

$$\epsilon^{(i)} = \max_{\substack{w^{(i+1)}: \|w^{(i+1)}\|_2 \leq 1 \\ \|w_m^{(i+1)} - w_m^{(i)}\|_2 \leq \Delta}} \left\{ \min_{j \in \{0, \dots, i\}} (w^{(i+1)})^T (\mu(\pi_E) - \mu(\pi^{(j)})) \right\}, \quad (1)$$

where  $\Delta$  is the trust-region radius ( $\Delta = \Delta^{(i)}$ ). The ratio  $\epsilon^{(i)}/\epsilon^{(i-1)}$  is used to determine the acceptance of the newly updated  $w$ . Since the predicted upper bound of  $\epsilon^{(i)}/\epsilon^{(i-1)}$  is  $\frac{n}{\sqrt{n^2 + (1-\gamma)^2 \epsilon^2}}$  [2], we set the acceptance condition to

$$\frac{\epsilon^{(i)}}{\epsilon^{(i-1)}} \leq \alpha \frac{n}{\sqrt{n^2 + (1-\gamma)^2 \epsilon^2}}, \quad (2)$$

where  $\alpha < 1$  is the coefficient parameter. If the condition is failed, we drop the newly found  $w$ , and update  $\epsilon$  and  $w$  by solving Eq. 1 with  $\Delta = 0$ . As the last step of one iteration, we update  $\Delta$  as follows:

$$\Delta^{(i+1)} = \begin{cases} \min(c_u \Delta^{(i)}, b_u) & \text{Eq. 2 met} \\ \max(c_l \Delta^{(i)}, b_l) & \text{Eq. 2 failed p times} \\ \Delta^{(i)} & \text{otherwise} \end{cases} \quad (3)$$

where  $b_u$  and  $b_l$  are the upper- and lower-bound of the trust-region radius, and  $c_u > 1$  and  $c_l < 1$  are the coefficients.

Additionally, IRL retrains  $\pi$  at each iteration, such a process can be inefficient. In IRL-HWT, we improve the training efficiency by using the learned model parameters  $\theta^{(i)}$  of  $\pi^{(i)}$  at the  $i$ th iteration as the initial parameters for training  $\pi^{(i+1)}$  at the  $(i+1)$ th iteration. The whole learning process stops when  $\|\mu(\pi) - \hat{\mu}(\pi_E)\|_2 < \epsilon$ . The complete IRL-HWT algorithm is shown in Algorithm 1.

---

**Algorithm 1:** Inverse Reinforcement Learning with Hybrid Weight Tuning (IRL-HWT)

---

**Result:** policy  $\pi^{(i)}$

Initialization: Calculate  $\mu(\pi_E)$  with expert trajectories;

Set  $i = 0$ , set  $\epsilon, \gamma, \alpha, b_u, b_l, c_u, c_l, p$ ;

Randomly set the model parameters  $\theta^{(0)}$  for  $\pi^{(0)}$ ;

Compute  $\mu(\pi^{(0)})$ ;

Set  $w_m^{(0)}$  such that  $\|w_m^{(0)}\|_2 < 1$  (initial reward weights),  $w_l^{(0)} = \mathbf{0}$ ;

Compute  $\epsilon^{(0)} = (w^{(0)})^T (\mu(\pi_E) - \mu(\pi^{(0)}))$ , where  $w^{(0)} = [w_m^{(0)} w_l^{(0)}]$ ;

Set  $\Delta^{(0)}$ ;

**while**  $\epsilon^{(i)} > \epsilon$  **do**

    Set  $i = i + 1$ ;

    Compute the reward function  $R = ((w^{(i-1)})^T \phi)$ ;

    Using  $R$  and  $\theta^{(i-1)}$  in RL to compute an optimal policy  $\pi^{(i)}$ ;

    Compute  $\mu(\pi^{(i)})$ ;

    Solve Optimization 1 with  $\Delta = \Delta^{(i-1)}$ , and get solution  $\epsilon^{(i)}$  at  $w^{(i)}$ ;

**if** Eq. 2 is True **then**

        Accept;

**else**

        Reject, solve Eq. 1 with  $\Delta = 0$ , and update  $\epsilon^{(i)}$  and  $w^{(i)}$ ;

**end**

    Set  $\Delta^{(i)}$  with Eq. 3;

**end**

---

In many cases, we can design a reward function to incorporate domain knowledge. For example, in autonomous driving we want the vehicle to reach the goal without collision; an example reward function could be  $r = w_1 g - w_2 c$ , where  $g$  is the “goal-reaching” flag,  $c$  is the collision flag, and  $w_1$  and  $w_2$  are positive coefficients. In complex scenarios such as our city environment, however, such a naive reward function does not work well: sometimes learned policy will let the car stop there forever to avoid collision, or run in a circle in

a very small area (keep running without collision but cannot achieve the goal). In this case, we can learn from expert demonstrations to complete the task.

We can represent any reward function as a sum of its linear/nonlinear components. To set the initial reward function for IRL-HWT, we add each component of the reward function into the feature vector  $\phi(\cdot)$  and add the corresponding coefficient into empirically-initialized weights  $w_m$ . For these components, since they already contain prior information, we do not want to aggressively update them (in fact in the beginning of IRL, aggressively updating feature components can deteriorate the learning performance because the initial policy is unstable)—this is achieved via the trust-region optimization.

#### D. Learning from Accidents & Parameter Reuse

Inspired by ADAPS [3], we adopt *learning from accident* to improve the training efficiency. Specifically, when the car crashes during the online training process, we will backtrack for certain frames and let the expert drive for certain frames to avoid collision. Both the crash and collision-free data will be (re-)used to train the policy. Essentially, by imposing a non-uniform prior on extracted features, we can introduce certain expert experience to the learning process directly. In addition, although the synthesized reward function is different at each iteration, they may embed useful information to share with each other. Thus, inspired by transfer learning [19], we reuse learned parameters for continuous training for faster convergence. Lastly, the process of *learning from accident* can result in both negative examples and positive examples in learning, leading to better performance and faster convergence.

## IV. EXPERIMENTS AND RESULTS

In this section, we detail our experiments and results. All experiments are conducted using an Intel(R) Xeon(TM) W-2123 CPU, an Nvidia GTX 1080 GPU, and 32G RAM.

### A. Overall Performance

We compare IRL-HWT to the original IRL, Generative Adversarial Imitation Learning (GAIL) [16], adversarial inverse reinforcement learning (AIRL) [18], and end-to-end imitation learning (IM) [23]. We use deep Q-learning with [164, 150] hidden units in each of the 2 dense layers and 20% dropout rate in all RL-based methods except for GAIL, where we use TRPO [24] as proposed in the original paper. We also experimented with Q-learning using more complex network architectures up to 6 hidden layers, but found [164, 150] with 2 hidden layers works best. The features used in IRL-based methods are listed in Fig. 2.

The action space contains 25 discrete (rotational speed, acceleration) action pairs: 5 levels of rotational speed (steering angle) and 5 levels of acceleration (throttle value and break value). We set policy similarity threshold  $\epsilon = 0.1$  and discount factor  $\gamma = 0.99$  for both IRL and IRL-HWT. We collect 20,000 steps of the expert trajectory and train all methods for 24 hours (we also observe that training

ID	Feature name	Description
1	Waypoint reaching flag	1 if a waypoint is reached 0 otherwise
2	Distance reduction to waypoint	Distance (meter) reduction to waypoint compared to the previous step
3	Direction away from waypoint	Relative angle (radian) between the forward direction of AV and direction to waypoint
4	Collision flag	1 if collide 0 otherwise
5~25	Depth list	Depth (meter) in 21 directions uniformly distributed in 120 degree $[-60, 60]$ of the forward direction in bird's-eye-view
26~46	Object type list	Object type of the first object that stops the ray from AV, in 21 directions uniformly distributed in 120 degree $[-60, 60]$ of the forward direction in bird's-eye-view 0 if nothing is detected or unknown object 1 if static object 2 if dynamic object

Fig. 2. Features used in IRL and IRL-HWT (46 in total).

12 more hours will not improve the performance). For IRL-HWT, specifically, we set the trust-region acceptance coefficient  $\alpha = 0.9$ , trust-region increasing and decreasing coefficients  $c_u = 1.1$ ,  $c_l = 0.9$ , upper- and lower-bound  $b_u = 0.05$ ,  $b_l = 0.001$ , and consecutive rejection number  $p = 5$ . We empirically initialize weights for the first 4 features in Fig. 2 with values  $[0.4, 0.01, -0.1, -0.8]$  (represents the initial reward function), and let IRL-HWT to learn the weights of depth- and object type-related features. The perception network is trained using 10,000 frames of simulated data.

Our evaluation criterion is *how far can the AV travel under a fixed number of steps*. The AV will stop if it finishes 1,000 steps or is in collision. Since the AV is assumed to know beforehand a series of waypoints on its path to the goal, we compute the score based on the number of waypoints reached by the AV:

$$s_{final} = (n_{reached} + (1 - \frac{dist_{next}}{dist_{last,next}})) \times s_{unit}, \quad (4)$$

where  $s_{final}$  is the final score,  $n_{reached}$  is the number of waypoints reached,  $dist_{next}$  is the distance between the last position on the car's trajectory and the next waypoint on the car's route,  $dist_{last,next}$  is the distance between the last reached waypoint to the next waypoint, and  $s_{unit}$  is the unit score by reaching a waypoint, which is set to 100. Note that a negative score may appear if  $dist_{next} > dist_{last,next}$ , which happens when the car drives away from the next waypoint. We use three scenes shown in Fig. 3 for evaluation:

- Scene 1: open space with only moving vehicles;
- Scene 2: city street with only static obstacles;
- Scene 3: city street with moving vehicles and static obstacles.



Fig. 3. Screenshots of the three scenes used in our evaluation. From left to right: Scene 1, Scene 2, Scene 3.

We record the final scores  $s_{final,2}$  and  $s_{final,3}$  from Scene 2 and 3; and average trajectory length  $l_{final,1}$ ,  $l_{final,2}$ ,  $l_{final,3}$  from Scene 1, 2, and 3. Because  $s_{final}$  is computed

based on the waypoint position and there is no explicit waypoint in Scene 1, we do not compute  $s_{final,1}$ . All data are obtained by running each learning method for 100 times with randomly initialized scene configurations such as the start position and direction of the AV, and the start position, direction and speed of obstacle vehicles. Our method achieves the highest scores and can enable the AV to drive safely 10x further (and longer) than the other methods. The full results are shown in Table II.

TABLE II

PERFORMANCE OF IRL-HWT VS. OTHER METHODS, USING THE SCORE DEFINED IN EQN. 4 AND SAFE-TRAJECTORY LENGTH. OUR METHOD NOT ONLY ACHIEVES THE *highest scores* BUT ALSO ENABLES THE AV TO DRIVE SAFELY 10x FURTHER THAN THE OTHER METHODS. WE ALSO SHOW THE ABLATION STUDY OF EACH COMPONENT (NUP, RoP, LFA STAND FOR 'NON-UNIFORM PRIOR', 'REUSE OF PARAMETERS', AND 'LEARN FROM ACCIDENTS', RESPECTIVELY), AND THE COMBINATION OF ALL TECHNIQUES LEADS TO THE BEST PERFORMANCE.

Method	$s_{final,2}$	$s_{final,3}$	$l_{final,1}$	$l_{final,2}$	$l_{final,3}$
IM [23]	77.4	60.1	105.6 m	53.7 m	44.7 m
IRL [2]	110.7	59.7	228.8 m	69.4 m	33.2 m
GAIL [16]	103.0	52.8	49.1 m	69.9 m	35.1 m
AIRL [18]	119.9	83.6	74.1 m	73.6 m	50.7 m
IRL + NUP	130.7	103.4	235.6 m	140.4 m	133.9 m
IRL + RoP	108.3	64.8	220.7 m	68.4 m	35.7 m
IRL + LFA	124.8	84.6	241.4 m	103.6 m	61.5 m
<b>Ours (IRL+ALL)</b>	<b>205.8</b>	<b>177.3</b>	<b>276.3 m</b>	<b>748.4 m</b>	<b>324.2 m</b>

We also compare performance between methods using reward functions vs. expert demonstration in Table III. The reward function for RL is the same as the initial reward function for our method. RL with reward functions cannot learn well in complex tasks, as it's difficult to design good rewards. IRL with only expert data cannot learn well either, while our method with both expert data and reward functions performs the best.

TABLE III

OVERALL PERFORMANCE COMPARISON BETWEEN METHODS USING REWARD FUNCTIONS VS. EXPERT DATA. RL WITH REWARD FUNCTIONS (R) OR IRL WITH ONLY EXPERT DATA (ED) CANNOT LEARN WELL ON COMPLEX TASKS, WHILE OURS WITH BOTH EXPERT DATA AND REWARD FUNCTIONS (ED+R) OUTPERFORMS BOTH.

Method	$s_{final,2}$	$s_{final,3}$	$l_{final,1}$	$l_{final,2}$	$l_{final,3}$
RL (R)	99.8	59.1	72.9 m	59.7 m	39.7 m
IRL (ED)	110.7	59.7	228.8 m	69.4 m	33.2 m
<b>ours (ED+R)</b>	<b>205.8</b>	<b>177.3</b>	<b>276.3 m</b>	<b>748.4 m</b>	<b>324.2 m</b>

**Analysis:** RL relies heavily on domain knowledge to design a good reward function in complex tasks. IM does not generalize well to new environments since it only mimics the expert trajectory. IRL, GAIL, and AIRL can only take the expert trajectory as information source. Our method addresses the limitations of these methods by combining expert trajectory and domain knowledge. The hybrid weight tuning and trust region mechanism further enable our method to adapt to different environments and achieve desired behaviors. The update of trust-region-based reward functions enables reuse of model parameters to further improve training efficiency. Learning from accidents helps RL process to obtain more effective samples in the high-reward areas than random sampling in otherwise sparse regions. *By combining all of them*



leads to a much higher overall performance improvement than the sum of improvements from each technique alone, thereby resulting in considerably better performance than other existing methods, as shown in Table II.

### B. Ablation Study

The first attribute of our approach is the non-uniform prior for features. To test this attribute, we empirically set the weight for *collision* to be  $-0.8$ . We count the number of collisions from the original IRL and IRL+non-uniform prior (IRL-HWT) by running both approaches for 10,000 steps. As shown in Table IV, IRL-HWT can reduce the number of collisions up to 41%.

TABLE IV

THE NUMBER OF COLLISIONS IN DIFFERENT SCENES WITHIN 10,000 STEPS: ORIGINAL IRL VERSUS IRL+NON-UNIFORM PRIOR.

IRL+NON-UNIFORM PRIOR CAN REDUCE THE NUMBER OF COLLISIONS UP TO 41%.

Model	Scene 1	Scene 2	Scene 3
IRL	35	58	111
IRL+non-uniform prior	33	41	93

The second attribute of our approach is *learned model parameters for continuous training*, which aims to improve training efficiency. From the results shown in Fig. 4, we can see that by having this attribute, we can achieve comparable model performance at 2.5x faster.

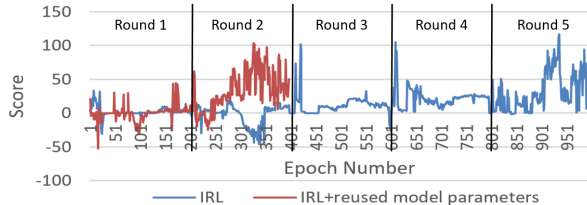


Fig. 4. By using previous model parameters for continuous training, we can achieve comparable score (model performance) 2.5x faster.

The last attribute of our approach is *learning from accidents*. We use ORCA [25] as the expert algorithm to generate alternative safe trajectories during the analysis of an accident. These trajectories are then used to generate additional training data for our algorithm. In Fig. 5, we show that the learning algorithm with this attribute added can achieve much higher scores up to *two orders of magnitude* in near collision scenarios under the same number of epochs.

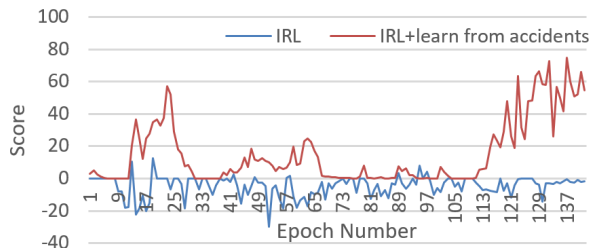


Fig. 5. Having the additional training data by *learning from accidents*, the learning algorithm achieves higher scores up to *two orders of magnitude* in near collision scenarios under the same number of epochs.

### C. Driving Cases

Our method can enable safe autonomous driving (without collision) in scenes with both static and dynamic obstacles.

We show the results of certain driving cases achieved by our method in Fig. 6. The full demo video can be found in the supplementary material.

In Fig. 6(a), we show that our method can steer the AV to make a left turn around the static obstacle while maintaining safe driving; In Fig. 6(b), we show that our method can lead the AV to avoid both static and dynamic obstacles. In particular, the AV (in green) steers to the right to avoid another vehicle coming from the opposite direction while moving away from the static obstacle; Lastly, we show that our approach can avoid multiple dynamic obstacles in Fig. 6(c). In this case, the AV (in green) is able to turn left to pass narrow space between two other vehicles.

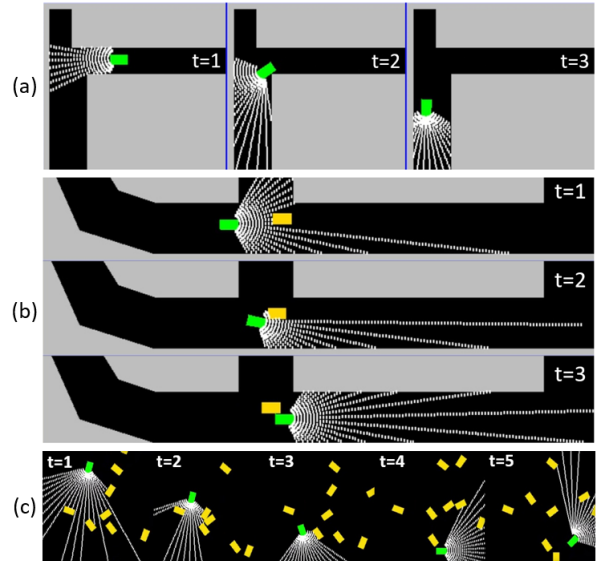


Fig. 6. (a) Static obstacle avoidance. Our method can lead the AV to make a left turn to avoid a static obstacle while maintaining safe driving. (b) Static and dynamic obstacle avoidance. The AV (in green) can avoid another vehicle (in yellow) coming from the opposite direction while steer away from the static obstacle. (c) Collision avoidance with multiple dynamic obstacles. Our method can direct the AV (in green) to avoid all nearby vehicles even when a narrow passage is presented.

## V. CONCLUSION AND FUTURE WORK

We propose a framework integrating context-aware multi-sensor perception and inverse reinforcement learning with hybrid weight tuning (IRL-HWT) for autonomous driving, which can incorporate both expert demonstrations and domain knowledge. We evaluate our approach using a variety of experiments, over the entire algorithm and each individual component. As shown in all comparisons, our method outperforms the state-of-the-art methods on all experiments.

**Limitations & Future Work:** We observe that the training results are relatively sensitive to the initial conditions. Our approach also inherits some limitations of IRL, e.g. information loss through encoding expert trajectories into a single feature expectation. We plan to develop a more efficient perception module by leveraging the sparsity embedded in the input data [26]. We also hope to alleviate the information loss during the expert's feature expectation. Lastly, we plan to test our approach in dense virtual traffic [27] reconstructed using real-world traffic data [28], [29], [30].

## REFERENCES

- [1] S. Ullman, "Against direct perception," *Behavioral and Brain Sciences*, vol. 3, no. 3, pp. 373–381, 1980.
- [2] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 1.
- [3] W. Li, D. Wolinski, and M. C. Lin, "ADAPS: Autonomous driving via principled simulations," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019, pp. 7625–7631.
- [4] W. Schwarting, J. Alonso-Mora, and D. Rus, "Planning and decision-making for autonomous vehicles," *Annual Review of Control, Robotics, and Autonomous Systems*, 2018.
- [5] Q. Chao, H. Bi, W. Li, T. Mao, Z. Wang, M. C. Lin, and Z. Deng, "A survey on visual traffic simulation: Models, evaluations, and applications in autonomous driving," *Computer Graphics Forum*, vol. 39, no. 1, pp. 287–308, 2019.
- [6] A. E. Sallab, M. Abdou, E. Perot, and S. Yogamani, "Deep reinforcement learning framework for autonomous driving," *Electronic Imaging*, vol. 2017, no. 19, pp. 70–76, 2017.
- [7] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2722–2730.
- [8] M. Pfeiffer, M. Schaeuble, J. Nieto, R. Siegwart, and C. Cadena, "From perception to decision: A data-driven approach to end-to-end motion planning for autonomous ground robots," in *2017 IEEE international conference on robotics and automation (icra)*. IEEE, 2017, pp. 1527–1533.
- [9] F. Codevilla, M. Müller, 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)*. IEEE, 2018, pp. 1–9.
- [10] M. Werling, S. Kammel, J. Ziegler, and L. Gröll, "Optimal trajectories for time-critical street scenarios using discretized terminal manifolds," *The International Journal of Robotics Research*, vol. 31, no. 3, pp. 346–359, 2012.
- [11] S. Sharifzadeh, I. Chiotellis, R. Triebel, and D. Cremers, "Learning to drive using inverse reinforcement learning and deep q-networks," *arXiv preprint arXiv:1612.03653*, 2016.
- [12] C. You, J. Lu, D. Filev, and P. Tsiotras, "Advanced planning for autonomous vehicles using reinforcement learning and deep inverse reinforcement learning," *Robotics and Autonomous Systems*, vol. 114, pp. 1–18, 2019.
- [13] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, *et al.*, "Maximum entropy inverse reinforcement learning," in *Aaai*, vol. 8. Chicago, IL, USA, 2008, pp. 1433–1438.
- [14] H. Le, N. Jiang, A. Agarwal, M. Dudík, Y. Yue, and H. Daumé, "Hierarchical imitation and reinforcement learning," in *International conference on machine learning*. PMLR, 2018, pp. 2917–2926.
- [15] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization," in *International conference on machine learning*. PMLR, 2016, pp. 49–58.
- [16] J. Ho and S. Ermon, "Generative adversarial imitation learning," *Advances in neural information processing systems*, vol. 29, pp. 4565–4573, 2016.
- [17] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum entropy deep inverse reinforcement learning," *arXiv preprint arXiv:1507.04888*, 2015.
- [18] J. Fu, K. Luo, and S. Levine, "Learning robust rewards with adversarial inverse reinforcement learning," *The International Conference on Learning Representations (ICLR)*, 2018.
- [19] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A comprehensive survey on transfer learning," *Proceedings of the IEEE*, 2020.
- [20] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [21] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3d proposal generation and object detection from view aggregation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1–8.
- [22] A. R. Conn, K. Scheinberg, and L. N. Vicente, *Introduction to derivative-free optimization*. SIAM, 2009.
- [23] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [24] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.
- [25] J. Van Den Berg, S. J. Guy, M. Lin, and D. Manocha, "Reciprocal n-body collision avoidance," pp. 3–19, 2011.
- [26] L. Lin, W. Li, and S. Peeta, "Efficient data collection and accurate travel time estimation in a connected vehicle environment via real-time compressive sensing," *Journal of Big Data Analytics in Transportation*, vol. 1, no. 2, pp. 95–107, 2019.
- [27] D. Wilkie, J. Sewall, W. Li, and M. C. Lin, "Virtualized traffic at metropolitan scales," *Frontiers in Robotics and AI*, vol. 2, p. 11, 2015.
- [28] W. Li, D. Nie, D. Wilkie, and M. C. Lin, "Citywide estimation of traffic dynamics via sparse GPS traces," *IEEE Intelligent Transportation Systems Magazine*, vol. 9, no. 3, pp. 100–113, 2017.
- [29] W. Li, D. Wolinski, and M. C. Lin, "City-scale traffic animation using statistical learning and metamodel-based optimization," *ACM Trans. Graph.*, vol. 36, no. 6, pp. 200:1–200:12, Nov. 2017.
- [30] W. Li, M. Jiang, Y. Chen, and M. C. Lin, "Estimating urban traffic states using iterative refinement and wardrop equilibria," *IET Intelligent Transport Systems*, vol. 12, no. 8, pp. 875–883, 2018.