

Driving with Style: Inverse Reinforcement Learning in General-Purpose Planning for Automated Driving

Sascha Rosbach^{1,2}, Vinit James¹, Simon Großjohann¹, Silviu Homoceanu¹ and Stefan Roth²

Abstract—Behavior and motion planning play an important role in automated driving. Traditionally, behavior planners instruct local motion planners with predefined behaviors. Due to the high scene complexity in urban environments, unpredictable situations may occur in which behavior planners fail to match predefined behavior templates. Recently, general-purpose planners have been introduced, combining behavior and local motion planning. These general-purpose planners allow behavior-aware motion planning given a single reward function. However, two challenges arise: First, this function has to map a complex feature space into rewards. Second, the reward function has to be manually tuned by an expert. Manually tuning this reward function becomes a tedious task. In this paper, we propose an approach that relies on human driving demonstrations to automatically tune reward functions. This study offers important insights into the driving style optimization of general-purpose planners with maximum entropy inverse reinforcement learning. We evaluate our approach based on the expected value difference between learned and demonstrated policies. Furthermore, we compare the similarity of human driven trajectories with optimal policies of our planner under learned and expert-tuned reward functions. Our experiments show that we are able to learn reward functions exceeding the level of manual expert tuning without prior domain knowledge.

I. INTRODUCTION

The trajectory planner in highly automated vehicles must be able to generate comfortable and safe trajectories in all traffic situations. As a consequence, the planner must avoid collisions, monitor traffic rules, and minimize the risk of unexpected events. General-purpose planners fulfill these functional requirements by optimization of a complex reward function. However, the specification of such a reward function involves tedious manual tuning by motion planning experts. Tuning is especially tedious if the reward function has to encode a humanlike driving style for all possible scenarios. In this paper, we are concerned with the automation of the reward function tuning process.

Unlike a strict hierarchical planning system, our planner integrates behavior and local motion planning. The integration is achieved by a high-resolution sampling with continuous actions [1]. Our planner, shown in Fig. 1, derives its actions from a vehicle transition model. This model is used to integrate features of the environment, which are then used to formulate a linear reward function. During every

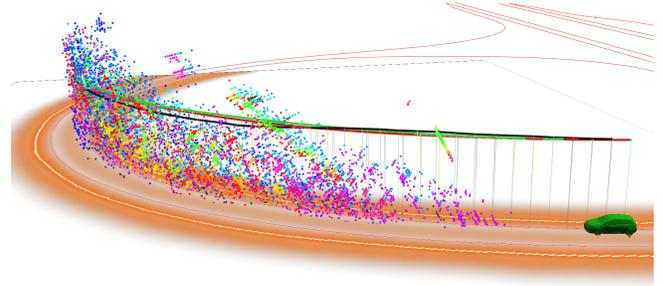


Fig. 1: This figure illustrates our general-purpose planner for automated driving. The color coding of the visualized state space indicates the state-action values. The z-axis corresponds to the velocity, while the groundplane depicts a subset of spatial features such as distance transformed lane centers and road boundaries. There are three color coded policies, black denotes the optimal policy of the planner, red the odometry of a human demonstration, and green the projection of the demonstration into the state space.

planning cycle of a model predictive control (MPC), the planning algorithm generates a graph representation of the high-dimensional state space. At the end of every planning cycle, the algorithm yields a large set of driving policies with multiple implicit behaviors, e.g., lane following, lane changes, swerving and emergency stops. The final driving policy has the highest reward value while satisfying model-based constraints. The reward function, therefore, influences the driving style of all policies without compromising safety.

Human driving demonstrations enable the application of inverse reinforcement learning (IRL) for finding the underlying reward functions, i.e., a linear combination of the reward weights. In this work, we utilize this methodology to automate the reward function tuning of our planner. Due to the planner’s exploration of a large set of actions, we are able to project demonstrated actions into our graph representation. Thereby, the demonstrations and associated features are efficiently captured. As a result, the learning algorithm enables the imitation of the demonstrated driving style. Most related work in IRL utilizes the state visitation frequency to calculate the gradient in *maximum entropy* IRL. However, the calculation of the state visitation is generally intractable in this high-dimensional state space. We utilize our graph representation to approximate the required empirical feature expectations to allow *maximum entropy* IRL.

¹The authors are with the Volkswagen AG, 38440 Wolfsburg, Germany {sascha.rosbach, vinit.james, simon.grossjohann, silviu.homoceanu}@volkswagen.de

²The authors are with the Visual Inference Lab, Department of Computer Science, Technische Universität Darmstadt, 64289 Darmstadt stefan.roth@visinf.tu-darmstadt.de

The main contributions of this paper are threefold: First, we formulate an IRL approach which integrates maximum entropy IRL with a model-predictive general-purpose planner. This formulation allows us to encode a humanlike driving style in a linear reward function. Second, we demonstrate the superiority of our automated reward learning approach over manual reward tuning by motion planning experts. We draw this conclusion on the basis of comparisons over various performance metrics as well as real world tests. Third, our automated tuning process allows us to generate multiple reward functions that are optimized for different driving environments and thereby extends the generalization capability of a linear reward function.

II. RELATED WORK

The majority of active planning systems for automated driving are built on the mediated perception paradigm. This paradigm divides the automated driving architecture into sub-systems to create abstractions from raw sensory data. The general architecture includes a perception module, a system to predict the intention of other traffic participants, and a trajectory planning system. The planning system is usually decomposed in a hierarchical structure to reduce the complexity of the decision making task [2]–[4]. On a strategic level, a route planning module provides navigational information. On a tactic and behavioral level, a behavioral planner derives the maneuver, e.g., lane-change, lane following, and emergency-breaking [5]. On an operational level, a local motion planner provides a reference trajectory for feedback control [6]. However, these hierarchical planning architectures suffer from uncertain behavior planning due to insufficient knowledge about motion constraints. As a result, a maneuver may either be infeasible due to over-estimation or discarded due to under-estimation of the vehicle capabilities. Furthermore, behavior planning becomes difficult in complex and unforeseen driving situations in which the behavior fails to match predefined admissibility templates. Starting with the work of McNaughton, attention has been drawn to parallel real-time planning [1]. This approach enables sampling of a large set of actions that respect kinematic constraints. Thereby a sequence of sampled actions can represent complex maneuvers. This kind of general-purpose planner uses a single reward function, which can be adapted online by a behavioral planner without the drawbacks of a hierarchical approach. However, it is tedious to manually specify and maintain a set of tuned reward functions. The process required to compose and tune reward functions is outside the scope of McNaughton’s work. We adopt the general-purpose planning paradigm in our approach and focus on the required tuning process.

Reward functions play an essential role in general-purpose planning. The rewards encode the driving style and influence the policy selection. Recently, literature has been published on the feature space of these reward functions. Heinrich et al. [7] propose a model-based strategy to include sensor coverage of the relevant environment to optimize the vehicle’s future pose. Gu et al. [8] derive tactical features from the

large set of sampled policies. So far, however, there has been little discussion about the automated reward function tuning process of a general-purpose planner.

Previous work has investigated the utilization of machine learning in hierarchical planning approaches to predict tactical behaviors [5]. Aside of behavior prediction, a large and growing body of literature focuses on finding rewards for behavior planning in hierarchical architectures [9], rewards associated with spatial traversability [10], and rewards for single-task behavior optimization of local trajectory planners [11]. The IRL approach plays an import role in finding the underlying reward function of human demonstrations for trajectory planning [12]. Similar to this work, several studies have investigated IRL in high-dimensional planning problems with long planning horizons. Shiarlis et al. [13] demonstrate maximum margin IRL within a randomly-exploring random tree (RRT*). Byravan et al. [14] focus on a graph-based planning representation for robot manipulation, similar to our planning problem formulation. Compared to previous work in IRL, our approach integrates IRL directly into the graph construction and allows the application of *maximum entropy* IRL for long planning horizons *without* increasing the planning cycle time.

Compared to supervised learning approaches such as direct imitation and reward learning, reinforcement learning solves the planning problem through learning by experience and interaction with the environment. Benefits of reinforcement learning are especially notable in the presence of many traffic participants. Intention prediction of other traffic participants can be directly learned by multi-agent interactions. Learned behavior may include complex negotiation of multiple driving participants [15]. Much of the current literature focuses on simulated driving experience and faces challenges moving from simulation to real-world driving, especially in urban scenarios. Another challenge includes the formulation of functional safety within this approach. Shalev-Shwartz et al. [16] describe a safe reinforcement learning approach that uses a hierarchical options graph for decision making where each node within the graph implements a policy function. In this approach, driving policies are learned whereas trajectory planning is not learned and bound by hard constraints.

Most of the current work in IRL utilizes the *maximum entropy* principle by Ziebart et al. [17] that allows training of a probabilistic model by gradient descent. The gradient calculation depends on the state visitation frequency, which is often calculated by an algorithm similar to backward value iteration in reinforcement learning. Due to the curse of dimensionality, this algorithm is intractable for driving style optimization in high-dimensional continuous spaces. Our work extends previous work by embedding IRL into a general-purpose planner with an efficient graph representation of the state space. The design of the planner effectively enables the driving style imitation without a learning task decomposition. As a result, we utilize the benefits of a model-based general-purpose planner and reward learning to achieve nuanced driving style adaptations.

III. PRELIMINARIES

The interaction of the agent with the environment is often formulated as a Markov Decision Process (MDP) consisting of a 5-tuple $\{\mathcal{S}, \mathcal{A}, T, R, \gamma\}$, where \mathcal{S} denotes the set of states, and \mathcal{A} describes the set of actions. A continuous action a is integrated over time t using the transition function $T(s, a, s')$ for $s, s' \in \mathcal{S}, a \in \mathcal{A}$. The reward function R assigns a reward to every action \mathcal{A} in state \mathcal{S} . The reward is discounted by γ over time t .

In this work, a model of the environment M returns a feature vector \mathbf{f} and the resultant state s' after the execution of action a in state s . The reward function R is given by a linear combination of K feature values f_i with weights θ_i such that $\forall (s, a) \in \mathcal{S} \times \mathcal{A} : R(s, a) = \sum_{i \in K} -\theta_i f_i(s, a)$. A policy π is a sequence of time-continuous transitions T . The feature path integral f_i^π for a policy π is defined by $f_i^\pi = \int_t \gamma_t f_i(s_t, a_t) dt$. The path integral is approximated by the iterative execution of sampled state-action sets \mathcal{A}_s in the environment model M . The value V^π of a policy π is the integral of discounted rewards during continuous transitions $V^\pi = \int_t \gamma_t R(s_t, a_t) dt$. An optimal policy π^* has maximum cumulative value, where $\pi^* = \arg \max_\pi V^\pi$. A human demonstration ζ is given by a vehicle odometry record. A projection of the odometry record ζ into the state-action space allows us to formulate a demonstration as policy π^D . For every planning cycle, we consider a set of demonstrations Π^D , which are geometrically close to the odometry record ζ . The planning algorithm returns a finite set of policies Π with different driving characteristics. The final driving policy π^S is selected and satisfies model-based constraints.

IV. METHODOLOGY

The planning system in Fig. 2 uses MPC to address continuous updates of the environment model. A periodic trigger initiates the perception system for which the planner returns a selected policy. In the following, we give an overview of the general-purpose planner. We use the nomenclature of reinforcement learning to underline the influence of reward learning in the context of search-based planning. Furthermore, we propose a *path integral maximum entropy* IRL formulation for high-dimensional reward optimization.

A. General-Purpose Planner for Automated Driving

Our planning algorithm for automated driving in all driving situations is based on [7], [18]. The planner is initialized at state s_0 , either by the environment model or in subsequent plans by the previous policy, and designed to perform an exhaustive forward search of actions to yield a set of policies Π . The set Π implicitly includes multiple behaviors, e.g., lane following, lane changes, swerving, and emergency stops [1]. Fig. 2 visualizes the functional flow of the planning architecture during inference and training.

Algo. 1 formally describes our search-based planning approach. The planner generates trajectories for a specified planning horizon H . Trajectories for the time horizon H are iteratively constructed by planning for discrete transition

Algorithm 1: General-Purpose Planner

Input: planning horizon H , model M ,
reward function R , reward discount factor γ
Output: policies Π , planning solution π^S

```

1 function SearchAlgorithm( $H, M, R, \gamma$ )
2   for  $t$  in  $H$  do
3      $\mathcal{S}_t \leftarrow$  get set of states
4     forall  $s \in \mathcal{S}_t$  do
5        $\mathcal{A}_s \leftarrow$  sample set of actions
6       forall  $a \in \mathcal{A}_s$  do
7         execute action  $a$  in  $M(s, a)$ 
8         observe resultant state  $s'$ 
9         observe resultant transition  $T$ 
10        observe resultant features  $\mathbf{f}(s, a)$ 
11        construct labels  $c(T)$ 
12         $R(s, a) \leftarrow \sum_{i \in K} -\theta_i f_i(s, a)$ 
13         $V(s') \leftarrow V(s) + \gamma_t R(s, a)$ 
14       $\mathcal{S}_{t+1} \leftarrow$  prune  $\mathcal{S}_t$ 
15   $\Pi \leftarrow$  get policies in  $\mathcal{S}, \mathcal{A}$ 
16   $\pi^S \leftarrow$  select model-based from  $\Pi$ 

```

lengths. The planner uses the parallelism of a graphics processing unit (GPU) to sample for all states $s \in \mathcal{S}_t$ a discrete number of continuous actions \mathcal{A}_s , composed of accelerations and wheel angles. The sampling distribution for each state is based on feasible vehicle dynamics. The actions itself are represented by time-continuous polynomial functions, where order and coefficients are derived from actor-friendly continuity constraints. This results in longitudinal actions described by velocity profiles up to fifth order, and lateral actions described by wheel angle profiles up to third order.

The search algorithm calls the model of the environment M for all states $s \in \mathcal{S}_t$ to observe the resultant state s' , transition T , and features \mathbf{f} for each state-action tuple. The feature vector \mathbf{f} is generated by integrating the time-continuous actions in the environment model. A labelling function assigns categorical labels to transitions, e.g., a label associated with collision. A pruning operation limits the set of states \mathcal{S}_{t+1} for the next transition step $t + 1 \in H$. Pruning is performed based on the value $V(s)$, label c , and properties of the reachable set \mathcal{S}_t to terminate redundant states with low value $V(s)$. This operation is required, first to limit the exponential growth of the state space, and second to yield a policy set Π with maximum behavior diversity. The algorithm is similar to parallel breadth first search and forward value iteration. The final driving policy π^S is selected based on the policy value $V(\pi)$ and model-based constraints.

B. Inverse Reinforcement Learning

The driving style of a general-purpose motion planner is directly influenced by the reward function weights θ . The goal of IRL is to find these reward function weights θ that enable the optimal policy π^* to be at least as good

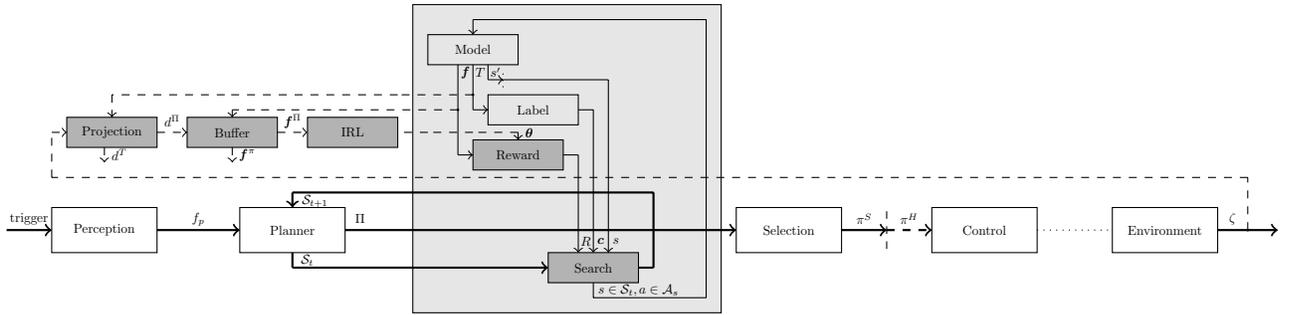


Fig. 2: Functional flow block diagram: The left input of a block corresponds to the output of the previous function. The inputs on top of the blocks denote intermediate outputs of previous functions. A thick solid line indicates the main flow from the environment perception f_p to the driven trajectory ζ . The vehicle control architecture is outside the scope of this work. In this work, we focus on the dark grey blocks of the architecture that influence the reward learning. Dashed connections between the blocks indicate the information flow during the training procedure. During data collection, we record the environment as well as the odometry ζ of the hidden driving policy of a human π^H .

as the demonstrated policy π^D , i.e., $V(\pi^*) \geq V(\pi^D)$. Thereby, the planner indirectly imitates the behavior of a demonstration [19]. However, learning a reward function given an optimal policy is ambiguous since many reward functions may lead to the same optimal policy [20]. Early work in reward learning for A* planning and dynamic programming approaches utilized structured maximum-margin classification [21], yet this approach suffers from drawbacks in the case of imperfect demonstrations [17]. Over the past decade, most research in IRL has focused on maximizing the entropy of the distribution on state-actions under the learned policy, which is known as *maximum entropy* IRL. This problem formulation solves the ambiguity of imperfect demonstrations by recovering a distribution over potential reward functions while avoiding any bias [17]. Ziebart et al. [17] propose a state visitation calculation, similar to backward value iteration in reinforcement learning, to compute the gradient of the entropy. The gradient calculation is adopted by most of the recent work in IRL for low-dimensional, discrete action spaces, which is inadequate for driving style optimizations. Our desired driving style requires high-resolution sampling of time-continuous actions, which produces a high-dimensional state space representation. In the following, we describe our intuitive approach, which combines search-based planning with *maximum entropy* IRL.

C. Path Integral Maximum Entropy IRL

In our IRL formulation, we maximize the log-likelihood L of expert behavior in the policy set Π by finding the reward function weights θ that best describe human demonstrations $\pi^D \in \Pi^D$ within a planning cycle, which is given by

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \sum_{\pi^D \in \Pi^D} \ln p(\pi^D | \theta) \quad (1)$$

$$= \arg \max_{\theta} \sum_{\pi^D \in \Pi^D} \ln \frac{1}{Z} \exp(-\theta f^{\pi^D}), \quad (2)$$

where the partition function is defined by $Z = \sum_{\pi \in \Pi} \exp(-\theta f^{\pi})$.

Similar to Aghasadeghi et al. [22], we optimize under the constraint of matching the feature path integrals f^{π} of the demonstration and feature expectations of the explored policies,

$$\forall i \in 1, \dots, k : \sum_{\pi \in \Pi} p(\pi | \theta) f_i^{\pi} = \frac{1}{m} \sum_{\pi^D \in \Pi^D} f_i^{\pi^D} = \hat{f}_i^{\Pi^D}, \quad (3)$$

where $\hat{f}_i^{\Pi^D}$ references the empirical mean of feature i calculated over m demonstrations in Π^D . The constraint in Eq. 3 is used to solve the non-linear optimization in Eq. 2.

The gradient of the log-likelihood can be derived as,

$$\nabla L(\theta) = \sum_{\pi \in \Pi} p(\pi | \theta) f^{\pi} - \hat{f}^{\Pi^D}, \quad (4)$$

and allows for gradient descent optimization.

The calculation of the partition function Z in Eq. 2 is often intractable due to the exponential growth of the state-action space over the planning horizon. The parallelism of the action sampling of the search-based planner allows us to explore a high-resolution state representation S_t for each discrete planning horizon increment t . A pruning operation terminates redundant states having sub-optimal behavior in the reachable set S_t , which is denoted by a lower value $V(s)$. Therefore, the pruning operation ensures multi-behavior exploration of the reachable set S_t that is evaluated with a single reward function. Thereby our sample-based planning methodology allows us to approximate the partition function similar to Markov chain Monte Carlo methods.

Once we obtain the new reward function, the configuration of the planner is updated. Hence, policies that have similar features as the human demonstration acquire a higher value assignment. This implies that they are more likely to be chosen as driving policy.

V. EXPERIMENTS

We assess the performance of *path integral maximum entropy* IRL in urban automated driving. We focus on a base feature set for static environments, similar to the manual

tuning process of a motion planning expert. After this process more abstract reward features are tuned relative to the base features.

A. Data Collection and Simulation

Our experiments are conducted on a prototype vehicle, which uses a mediated perception architecture to produce feature maps as illustrated in Fig. 1. We recorded data in static environments and disabled object recognition and intention prediction. The data recordings include features of the perception system as well as odometry recordings of the human driver’s actions. The training of our algorithm is performed during playbacks of the recorded data. After every planning cycle of the MPC, the position of the vehicle is reset to the odometry recording of the human demonstration.

B. Projection of Demonstration in State Space

The system overview in Fig. 2 includes a projection function that transfers the actions of a manual drive into the state-action space of the planning algorithm. The projection metric d is calculated during the graph construction between odometry ζ and continuous transitions $T(s, a, s')$ of all policies π in the set Π :

$$d(\zeta, \pi) = \int_t \alpha_t \|\zeta_t - \pi_t\| dt. \quad (5)$$

The norm is based on geometrical properties of the state space, e.g., the Euclidean distance in longitudinal and lateral direction as well as the squared difference in the yaw angle. Further, the metric includes a discount factor α_t over the planning horizon. The policy π^D has the least discounted distance towards the odometry record. There are multiple benefits of using the projection metric: First, the projected trajectory includes all constraints of the planner. If the metric surpasses a threshold limit, the human demonstrator does not operate in the actor’s limits of the vehicle and therefore can not be used as a valid demonstration. Second, the projection metric allows for an intuitive evaluation of the driving style based on the geometrical proximity to the odometry. Third, we may augment the number of demonstrations by loosening the constraint of the policy π^D to have least discounted distance towards the odometry. Thereby, multiple planner policies qualify as demonstration $\pi^D \subseteq \Pi^D$.

C. Reward Feature Representation

In this work, the reward function $R(s, a)$ is given by a linear combination of K reward features. The features describe motion and infrastructural rewards of driving. The discount factor γ is manually defined by a motion planning expert and is not optimized at this stage of the work. Our perception system in Fig. 2 provides normalized feature maps with spatial information of the environment. The feature path integral f^π of a policy π is created by transitioning through the feature map representation of the environment. We concentrate on a base reward set consisting of $K = 12$ features, which are listed in the legend of Fig. 4a. Heinrich et al. formally described a sub-set of our features [18]. Seven

of our feature values describe the motion characteristics of the policies, which are given by derivatives of the lateral and longitudinal actions. They include the difference between the target and policy velocity, and the acceleration and jerk values of actions. The target in the velocity may change depending on the situation, e.g., in a situation with a traffic light the target velocity may reduce to zero. Furthermore, the end direction feature is an important attribute for lateral behavior that specifies the angle towards the driving direction at the end of the policy. The creeping feature suppresses very slow longitudinal movement in situations, where a full stop is more desired. Infrastructural features include proximity measures to the lane center and curbs, cost potentials for lanes, and direction. Furthermore, we specify a feature for conflict areas, e.g., stopping at a zebra crossing.

D. Implementation Details

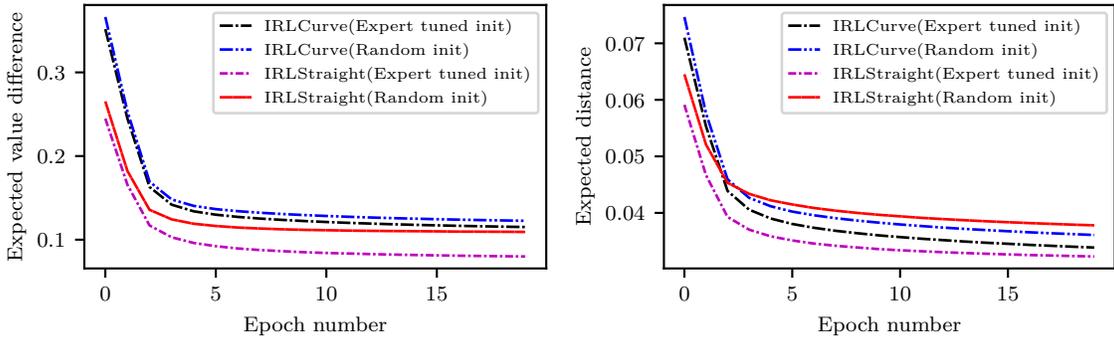
During the playback of a human demonstration, the path integral feature vectors f^Π of the policy set Π are approximated for every planning cycle and stored within a replay buffer. By including our projection metric in the action sampling procedure, we associate each policy π with the distance to the odometry of the human demonstration. During training, we query demonstrations Π^D , which are policies with a low projection metric value, from our replay buffer where $\pi^D \subseteq \Pi^D \subseteq \Pi$. Hence, the replay buffer contains features of demonstrations for each planning cycle denoted as $f^{\Pi^D} \subseteq f^\Pi$. Fig. 2 describes the information flow from the odometry record of the demonstration to the feature query from the replay buffer. Due to actor constraints of the automated vehicle’s actions, the planning cycles without demonstrations are not considered for training. We utilize experience replay similar to reinforcement learning and update on samples or mini-batches of experience, by drawing randomly from the buffered policies. This process allows us to efficiently use previous experience, which can be trained on multiple times. Further, stability is provided by not altering the representation of the expert demonstrations within the graph representation.

VI. EVALUATION

We aim to evaluate the utility of our automated reward function optimization in comparison to manual expert tuning. First, we analyze our driving style imitation in terms of value convergence towards the human demonstration. Second, we compare the driving style of our policies under random, learned, and expert-tuned reward functions against human driving on dedicated test route segments.

A. Training Evaluation

We analyze the convergence for different training initializations and road segment types, namely straight and curvy. Due to the linear combination of reward weights, one expects a segment-specific preference of the reward function. As a reference, a motion planning expert generated a tuned reward function for general driving. We perform two drives per training segment, one with a random and one with an



(a) Difference between expected value of human driving demonstration and expected value of planner policies under learned reward functions.

(b) Expected distance of planner policies towards the human driving demonstration under learned reward functions.

Fig. 3: Illustration of training and validation metrics for multiple segments and training initializations. Convergence of maximum entropy IRL over training epochs. Validation of the training by indicating the reduction of the expected distance towards the human demonstration. The probability is calculated independently for every planning cycle of the MPC, whereas the policy set includes an average of approx. 4000 policies.

expert-tuned reward function. The policies to be considered as human demonstrations are chosen based on our projection metric and therefore depend on our chosen reward function initialization. The expert initialization yields demonstrations with a mean projection error 7% lower as compared to random initialization. During every planning cycle on the segments, we trace the policies of the planner in replay buffers. We generate four tuned reward functions which are referred to in Fig. 4a by training on our replay buffers.

The convergence of the training algorithm is measured by the expected value difference (EVD) over training epochs between learned and demonstrated policies. The EVD is calculated for every planning cycle and averaged over the segment. The EVD is given by

$$\begin{aligned} & \mathbb{E}[V(\Pi)] - \mathbb{E}[V(\Pi^D)] \quad (6) \\ &= \sum_{\pi \in \Pi} p(\pi|\theta)V(\pi) - \sum_{\pi^D \in \Pi^D} p(\pi^D|\theta)V(\pi^D). \quad (7) \end{aligned}$$

The performance of the random and expert-tuned reward functions is given by the EVD at epoch zero. The initial and final EVD differences between the straight and curvy segment is 30% and 19% respectively. A preference of the straight segments by both reward functions is visible in the initial EVD difference Fig. 3. The learned reward functions show a large EVD reduction of 67% for curvy and 63% for straight segments at the end of training.

We can interpret the training results in the following way:

- (a) The projection metric depends on the quality of the reward function.
- (b) Improved reward functions lead to improved action sampling and therefore produce better demonstrations.
- (c) Learning reward functions without prior knowledge is possible, e.g. generating a replay buffer with a randomly initialized reward function and training with a random initialization.
- (d) Unsuitable reward functions improve more significantly during training.

Hence, continuously updating the policies in the replay buffer generated from an updated reward function should lead to faster convergence.

The desired driving style is given by the actions of a human driving demonstration. Therefore, the projection error in Eq. 5, which we use to select driving demonstrations, extends itself as a direct validation metric of the actions. Due to our goal of optimizing the likelihood of human driving behavior within the policy set, we calculate the expected distance (ED) in the policy set, given by

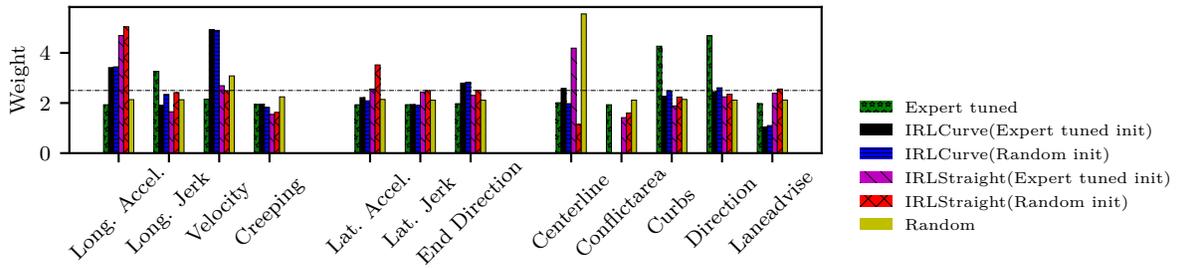
$$\mathbb{E}[d(\zeta, \Pi)] = \sum_{\pi \in \Pi} p(\pi|\theta)d(\zeta, \pi). \quad (8)$$

The learned reward functions in Fig. 3b show a large ED reduction of 54% for curvy and 44% for straight segments at the end of training. The ED reduction trends have high similarity to the above mentioned EVD trend and therefore this validates the premise of a high correlation between value and distance to the demonstration. An improved expected distance ensures a high likelihood of selecting policies which are similar to humanlike driving demonstrations.

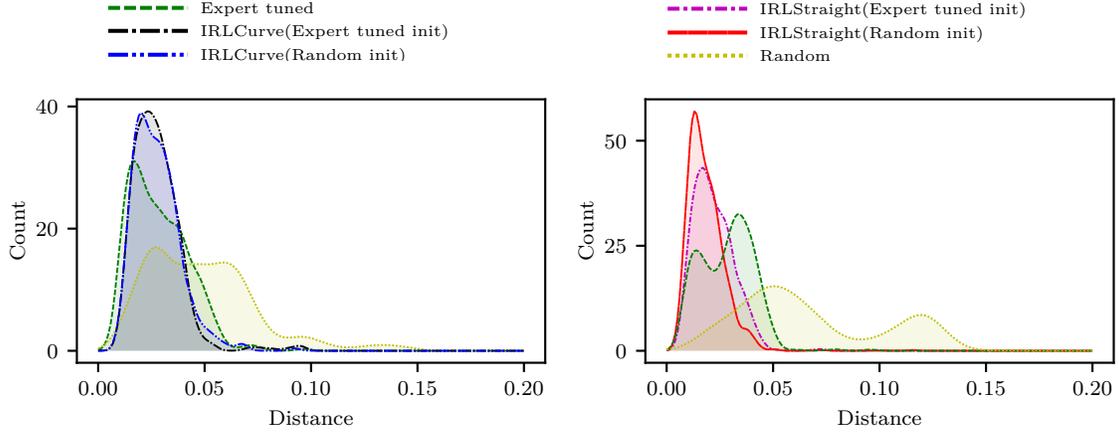
B. Driving Style Evaluation

In this part of the evaluation, we compare the driving style of the random, learned, and expert-tuned reward functions shown in Fig. 4a to manual human driving. The parameters of the reward functions allow for introspection and reasoning about the segment-specific preference. The reward weight is inversely proportional to the preference of that feature value in the policy. Learned reward functions are of two types:

- (a) IRL with random initialization, hereby referred as IRL(random). Both the training trajectory set and the learning task are randomly initialized.
- (b) IRL with expert initialization, hereby referred as IRL(expert). Both the training trajectory set and the learning task are initialized by expert tuning.

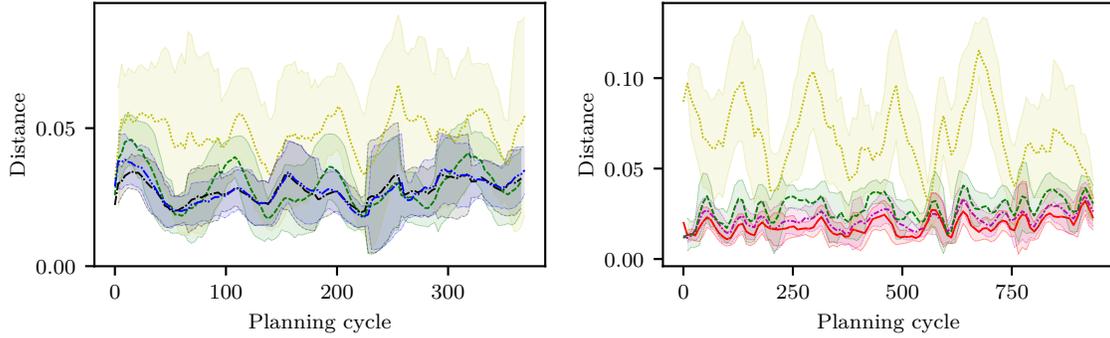


(a) Feature weights of tested reward functions.



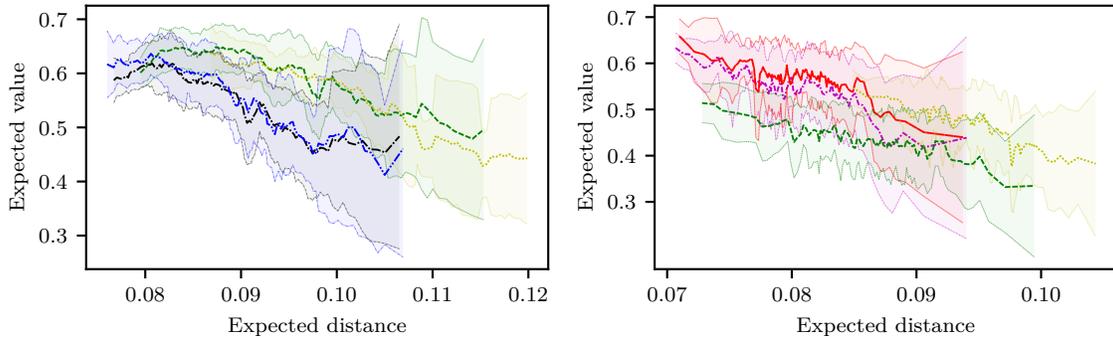
(b) Distances of the optimal policies on a curvy test segment.

(c) Distances of the optimal policies on a straight test segment.



(d) Distances of the optimal policy over planning cycles on a curvy test segment.

(e) Distances of the optimal policy over planning cycles on a straight test segment.



(f) Expected value and distance under the expert tuned reward function of every planning cycle on a curvy test segment.

(g) Expected value and distance under the expert tuned reward function of every planning cycle on a straight test segment.

Fig. 4: The tests contrast the driving style of random, learned, and expert-tuned reward functions. The graphs present the results of independent playbacks on a dedicated test track. The probability is calculated independently for every planning cycle of the MPC, whereas the policy set includes on average 4000 policies.

Using these reward functions, we run our planning algorithm on dedicated test route segments to verify the generalized performance of the optimal policies. We carry out multiple drives over the test segments to generate representative statistics. Fig. 4b and Fig. 4c present the projection metric distribution, which is the distance of the optimal policy to the odometry of a manual human drive for every planning cycle. We fit a Gaussian distribution over the histogram with 200 bins of size 0.001 with 944 planning cycles for the straight and 369 planning cycles for the curvy segment. The learned reward functions improve the driving style on all segments even in the case of random initialization. Our evaluation metric, which is the mean distance of the optimal policy to the odometry, decreases for IRLStraight(random) by 73% and for IRLCurve(random) by 43%. In case of expert-tuned initialization, IRLStraight(expert) decreased by 22% and IRLCurve(expert) by 4%. The strong learning outcome in the straight segment can be attributed to the easier learning task as compared to the curvy segment. Even though the expert-tuned reward functions do not improve substantially in terms of mean distance, they show a lower variance in distance of the optimal policy to the odometry over planning cycles after training as is shown in Fig. 4d and Fig. 4e. Here we indicate variance in the distance of the optimal policy over planning cycles by one standard deviation. The variance reduction of learned reward function depicts higher stability over planning cycles. Hence, we are able to encode the human driving style through IRL without applying prior domain knowledge as done by motion planning experts.

Fig. 4f and Fig. 4g present the expected value of our evaluated reward functions r^A under the expert-tuned reward function r^E , given by $\mathbb{E}[V(\Pi)] = \sum_{\pi \in \Pi} p(V^\pi(r^A))V^\pi(r^E)$. The overall trend indicates an inverse relationship between expected value and expected distance. The learned reward functions have lower expected distance as compared to expert tuned and random reward functions, while having a higher rate of value reduction with increasing expected distance. This ensures that the learned reward functions induce a high degree of bias in the policy evaluation such that the humanlike demonstrated behavior is preferred.

VII. CONCLUSION AND FUTURE WORK

We utilize *path integral maximum entropy* IRL to learn reward functions of a general-purpose planning algorithm. Our method integrates well with model-based planning algorithms and allows for automated tuning of the reward function encoding the humanlike driving style. This integration makes maximum entropy IRL tractable for high dimensional state spaces. The automated tuning process allows us to learn reward functions for specific driving situations. Our experiments show that learned reward functions improve the driving style exceeding the level of manual expert-tuned reward functions. Furthermore, our approach does not require prior knowledge except the defined features of the linear reward function. In the future, we plan to extend our IRL approach to update the reward function dynamically.

REFERENCES

- [1] M. McNaughton, "Parallel Algorithms for Real-time Motion Planning," Ph.D. dissertation, Carnegie Mellon University, 2011.
- [2] C. Katrakazas, M. Qaddus, W.-H. Chen, and L. Deka, "Real-time motion planning methods for autonomous on-road driving: State-of-the-art and future research directions," *Transportation Res. Part C: Emerging Technologies*, vol. 60, pp. 416–442, 2015.
- [3] 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 Trans. Intelligent Vehicles*, vol. 1, no. 1, pp. 33–55, 2016.
- [4] W. Schwarting, J. Alonso-Mora, and D. Rus, "Planning and Decision-Making for Autonomous Vehicles," *Annu. Rev. Control Robot. Auton. Syst.*, vol. 1, no. 1, pp. 187–210, 2018.
- [5] S. Ulbrich and M. Maurer, "Towards Tactical Lane Change Behavior Planning for Automated Vehicles," in *Proc. IEEE Int. Conf. Intell. Transp. Syst. (ITSC)*, 2015.
- [6] M. Werling, J. Ziegler, S. Kammel, and S. Thrun, "Optimal trajectory generation for dynamic street scenarios in a frenet frame," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2010.
- [7] S. Heinrich, J. Stubbemann, and R. Rojas, "Optimizing a driving strategy by its sensor coverage of relevant environment information," in *IEEE Intell. Vehicles Symp.*, 2016, pp. 441–446.
- [8] T. Gu, J. M. Dolan, and J.-W. Lee, "Automated tactical maneuver discovery, reasoning and trajectory planning for autonomous driving," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Syst. (IROS)*, Daejeon, South Korea, 2016.
- [9] P. Abbeel, D. Dolgov, A. Y. Ng, and S. Thrun, "Apprenticeship learning for motion planning with application to parking lot navigation," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Syst. (IROS)*, 2008.
- [10] M. Wulfmeier, D. Rao, D. Z. Wang, P. Ondruska, and I. Posner, "Large-scale cost function learning for path planning using deep inverse reinforcement learning," *Int. J. Robotics Research*, vol. 36, no. 10, pp. 1073–1087, 2017.
- [11] M. Kuderer, S. Gulati, and W. Burgard, "Learning driving styles for autonomous vehicles from demonstration," in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2015.
- [12] S. Arora and P. Doshi, "A survey of inverse reinforcement learning: Challenges, methods and progress," in *arXiv Preprint arXiv:1806.06877*, 2018.
- [13] K. Shiarlis, J. Messias, and S. Whiteson, "Inverse Reinforcement Learning from Failure," in *Proc. Int. Conf. Autonomous Agents Multi-Agent Syst.*, 2016.
- [14] A. Byravan, M. Monfort, B. Ziebart, B. Boots, and D. Fox, "Graph-Based Inverse Optimal Control for Robot Manipulation," in *Proc. Int. Joint Conf. Artificial Intell. (IJCAI)*, vol. 15, 2015.
- [15] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "Safe, multi-agent, reinforcement learning for autonomous driving," in *Learning, Inference and Control of Multi-Agent Syst. Workshop (NIPS)*, 2016.
- [16] —, "On a Formal Model of Safe and Scalable Self-driving Cars," in *arXiv Preprint arXiv:1708.06374*, 2017.
- [17] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum Entropy Inverse Reinforcement Learning," in *Proc. Nat. Conf. Artificial Intell. (AAAI)*, vol. 8, 2008.
- [18] S. Heinrich, "Planning Universal On-Road Driving Strategies for Automated Vehicles," Ph.D. dissertation, Freie Universität Berlin, 2018.
- [19] A. Y. Ng and S. J. Russell, "Algorithms for Inverse Reinforcement Learning," in *Proc. Int. Conf. Machine Learning (ICML)*, 2000.
- [20] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proc. Int. Conf. Machine Learning (ICML)*. ACM, 2004.
- [21] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich, "Maximum margin planning," in *Proc. Int. Conf. Machine Learning (ICML)*, 2006.
- [22] N. Aghasadeghi and T. Bretl, "Maximum entropy inverse reinforcement learning in continuous state spaces with path integrals," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Syst. (IROS)*. IEEE, 2011.