# Teaching Autonomous Vehicles to Express Interaction Intent during Unprotected Left Turns: A Human-Driving-Prior-Based Trajectory Planning Approach

Jiaqi Liu, Xiao Qi, Ying Ni, Jian Sun and Peng Hang*

*Department of Traffic Engineering & Key Laboratory of Road and Traffic Engineering, Ministry of Education, Tongji University, China*

## ABSTRACT

With the integration of Autonomous Vehicles (AVs) into our transportation systems, their harmonious coexistence with Human-driven Vehicles (HVs) in mixed traffic settings becomes a crucial focus of research. A vital component of this coexistence is the capability of AVs to mimic human-like interaction intentions within the traffic environment. To address this, we propose a novel framework for Unprotected left-turn trajectory planning for AVs, aiming to replicate human driving patterns and facilitate effective communication of social intent. Our framework comprises three stages: trajectory generation, evaluation, and selection. In the generation stage, we use real human-driving trajectory data to define constraints for an anticipated trajectory space, generating candidate motion trajectories that embody intent expression. The evaluation stage employs maximum entropy inverse reinforcement learning (ME-IRL) to assess human trajectory preferences, considering factors such as traffic efficiency, driving comfort, and interactive safety. In the selection stage, we apply a Boltzmann distribution-based method to assign rewards and probabilities to candidate trajectories, thereby facilitating human-like decision-making. We conduct validation of our proposed framework using a real trajectory dataset and perform a comparative analysis against several baseline methods. The results demonstrate the superior performance of our framework in terms of human-likeness, intent expression capability, and computational efficiency. Limited by the length of the text, more details of this research can be found at https://shorturl.at/jqu35

## 1. Introduction

The advent of Autonomous Vehicles (AVs) is expected to usher in transformative changes to the entire transportation system. In the prevailing mixed-traffic scenario, the peaceful coexistence of AVs and Human-driven Vehicles (HVs) is a topic of paramount importance. Contrary to the tacit intent communication and interaction among human drivers through driving trajectories (such as trajectory deviation, speed adjustment, etc.) (De Ceunynck, Polders, Daniels, Hermans, Brijs and Wets (2013); Lee, Madigan, Giles, Garach-Morcillo, Markkula, Fox, Camara, Rothmueller, Vendelbo-Larsen, Rasmussen et al. (2021)), existing AVs struggle with expressing social intent at the trajectory behavior level (Wang, Wang, Zhang, Liu, Sun et al. (2022)). Their driving behaviors often seem ambiguous and challenging for humans to interpret, thus raising substantial concerns for interactive safety in mixed-driving environments (Wang et al. (2022)).

While discrete optimization-based trajectory planning methods can satisfy requirements such as ease of construction, trajectory smoothness, and continuous differentiability, they often fall short in accurately reflecting real human driving behaviors. Consequently, the trajectories generated by these methods significantly differ from actual human driving trajectories. In an attempt to address this discrepancy, prior research (Huang, Wu and Lv (2021); Dongjian, Bing, Jian, Jiayi and Yanchen (2022)) has put forth trajectory planning techniques that emulate human-like behaviors by learning from and imitating expert human trajectories.

Although these methods can employ real human-driven vehicle (HV) trajectories to instruct autonomous vehicle (AV) trajectory planning, they leave several issues unresolved. One key issue is that high-level decision-making

intentions play a pivotal role in constraining the trajectory planning space. However, existing learning methods often deduce AV decisions from the outcomes of trajectory selection, contradicting the genuine interactive operational principle of 'decision precedes execution'. Furthermore, these methods tend to disregard the encapsulation of intent expression information concealed within trajectory movements under the influence of high-level decisions. This oversight is particularly glaring in complex intersection interaction scenarios with diverse trajectories (Wang et al. (2022)). The intent information embedded within trajectory movements is critical during the interaction between both parties, with noticeable differences in trajectory motion strategies and anticipated spaces under varying decisions (Dongjian et al. (2022)). These challenges amplify in complexity and importance when an AV navigates an intersection, one of the most demanding interaction scenarios (Zhao, Knoop, Sun, Ma and Wang (2023)). This complexity often hinders an AV from executing unprotected left turns in the same way a human driver would.

To address the aforementioned challenges, we select the most representative unprotected left-turn scenario and propose a trajectory planning framework that learns from human trajectory prior knowledge and expresses social intent. This framework includes three stages: trajectory generation, evaluation, and selection. We initially extract implicit interaction intent rules from human trajectory data, and during the trajectory generation stage, we establish constraints on the expected trajectory space under various decisions based on human implicit intent experiences, generating candidate motion trajectories that consider intent expression. In the evaluation stage, we construct a reward function based on efficiency, comfort, and safety, employing maximum entropy inverse reinforcement learning (ME-IRL) to learn and evaluate human trajectory preferences. In the selection stage, we set up a probability distribution of reward gains and candidate trajectories based on the Boltzmann distribution, enabling human-like action selection.

We validate our framework using a real trajectory dataset and compare it against multiple baseline methods. The results demonstrate that our algorithm excels in human-likeness, intent expression capability, and computational efficiency.

In summary, our contributions are shown as follows:

- We propose an unprotected left-turn trajectory planning framework capable of expressing social intent, which includes three stages: trajectory generation, trajectory evaluation, and trajectory selection.

- During the trajectory generation phase, we extract implicit intent prior information from real human trajectory data, proposing a method for generating an expected trajectory space constrained by human interaction priors.

- We propose a method for learning and evaluating human trajectory preferences based on maximum entropy inverse reinforcement learning, and adopt a human-like action selection method based on the Boltzmann distribution.

## 2. Related Works

### 2.1. Trajectory Planning Methods

Once a decision intent is formulated during the motion process, Autonomous Vehicles (AVs) necessitate a suitable trajectory planning algorithm to create a collision-free, executable path. This planning ensures the vehicle can traverse from the start point to the destination efficiently and safely.

Discrete optimization-based trajectory planning methods, such as the Frenet trajectory planning method (Werling, Ziegler, Kammel and Thrun (2010)), are widely employed. This method operates on the Frenet coordinate system, converting the problem into two-dimensional S-T and L-T problems. Building on this, several studies (Zhang, Sun, Zhou, Pan, Hu and Miao (2020); Hu, Deng, Cao, Zhang, Khajepour, Zeng and Wu (2022); Hu, Chen, Tang, Cao and He (2018)) have enhanced and optimized the Frenet trajectory planning method. Apollo researchers (Zhang et al. (2020)), drawing from the Frenet trajectory planning method, proposed a quadratic optimization method to cater to the non-holonomic constraints of the vehicle. Hu et al. (Hu et al. (2018)) suggested a real-time dynamic path planning method for autonomous driving, computing collision risk via the Gaussian convolution algorithm to ensure safe interaction with both static and dynamic obstacles.

Other methods, like the Archimedean spiral (Ma, Sun and Wang (2017)) and the Bezier curve (Zhou, Ma, Zhang and Sun (2022)), are also used to characterize motion trajectories during intersection turns or high-speed lane changes. However, these methods necessitate more parameters such as trajectory control points and control parameters, and unlike polynomial methods that already incorporate the planning results of linear speed and acceleration, these methods need to further resolve motion planning parameters.

## 2.2. Social Behavior in Planning Algorithms

While discrete optimization-based trajectory planning methods can satisfy basic motion planning needs, significant disparities persist between its planned trajectories and human trajectories (Huang et al. (2021)). Moreover, its inability to express intent impedes effective interaction between AVs and HVs (Wang et al. (2022)).

In numerous existing studies (Hu et al. (2022); Werling et al. (2010)), feature weights in the reward function of trajectory planning are assigned manually or optimized through numerous experiments, resulting in a lack of alignment with driver cognitive characteristics. Inverse Reinforcement Learning (IRL) offers a solution to this issue by reconstructing the reward function through learning from expert example trajectories (Arora and Doshi (2021)). IRL has been extensively applied to problems such as route planning (Ziebart, Maas, Bagnell, Dey et al. (2008); Wulfmeier, Rao, Wang, Ondruska and Posner (2017)) and the learning of human behavior trajectory features (Abbeel and Ng (2004)). Notable IRL methods include Maximum Entropy IRL (Ziebart et al. (2008); Wulfmeier et al. (2017)), Apprentice IRL (Abbeel and Ng (2004)), and Bayesian IRL (Ramachandran and Amir (2007)). Among these, Maximum Entropy IRL (ME-IRL) holds substantial benefits in studying human expert trajectory behavior, including uncertainty modeling, learning behavior diversity, and robust generalization ability.

To address the traditional trajectory planning algorithms' deficiency in expressing social intent, we utilize ME-IRL to learn the human driving behavior trajectory selection mechanism and construct the expected trajectory space considering the intent expression mode under varying decisions.

## 3. Methodology

In this section, we introduce a social trajectory planning method framework with inherent intent expression capabilities. We begin with an overview of the entire left-turn trajectory planning framework, followed by an analysis of the interaction features of human drivers during left turns. Subsequently, we elaborate on our trajectory generation method, which incorporates intent expression.

### 3.1. Overview of Framework

The comprehensive left-turn trajectory planning framework proposed in this study is illustrated in Fig.1. The specific process and underlying philosophy of our trajectory planning method are as follows:

We scrutinize the characteristics of human trajectory intent expression in left-turn scenarios based on the human driving dataset, SIND (Xu, Shao, Li, Yang, Wang, Huang, Lv and Wang (2022)), providing prior knowledge of human drivers for trajectory generation. Given the decision-making intent of proceeding or yielding, coupled with prior human knowledge of trajectory intent expression, we generate multiple candidate trajectories to form an expected trajectory space for decision-making. We then formulate a trajectory reward function to evaluate the features of all generated candidate trajectories considering aspects such as traffic efficiency, driving comfort, and dynamic interactive safety. Subsequently, we establish a probability distribution of candidate trajectories based on the Boltzmann distribution, employing the Boltzmann noise rational model to emulate driver trajectory selection behavior. Lastly, we learn the criteria for trajectory evaluation and selection during human driving via the method of maximum entropy inverse reinforcement learning.

### 3.2. Analysis of Human Trajectory Intent Expression Characteristics

Unlike existing autonomous driving algorithms, human drivers often subtly express their intentions to interactive objects through speed adjustments and trajectory deviations in some interactive scenarios. As the movement trajectories in the intersection driving space are less confined and more flexible, left-turning HVs can communicate their intent decisions through distinct trajectory progression methods during interaction. We select unprotected left-turn interaction as a representative scenario and analyze the features of human trajectory intent expression using the SIND dataset. The distribution of unprotected left-turn interaction trajectories in the SIND dataset is shown in Fig. 2, where the orange line symbolizes the yielding left-turn trajectory and the blue line signifies the proceeding left-turn trajectory. It is evident that human drivers employ diverse intent expression methods for left-turn trajectories under varying decisions:

- As depicted in the left part of Fig. 2, the preceding trajectory employs a pre-turn behavior to pivot the vehicle head in advance and laterally displace it in the upstream section, swiftly concluding the turn in a direct turn manner within the intersection.
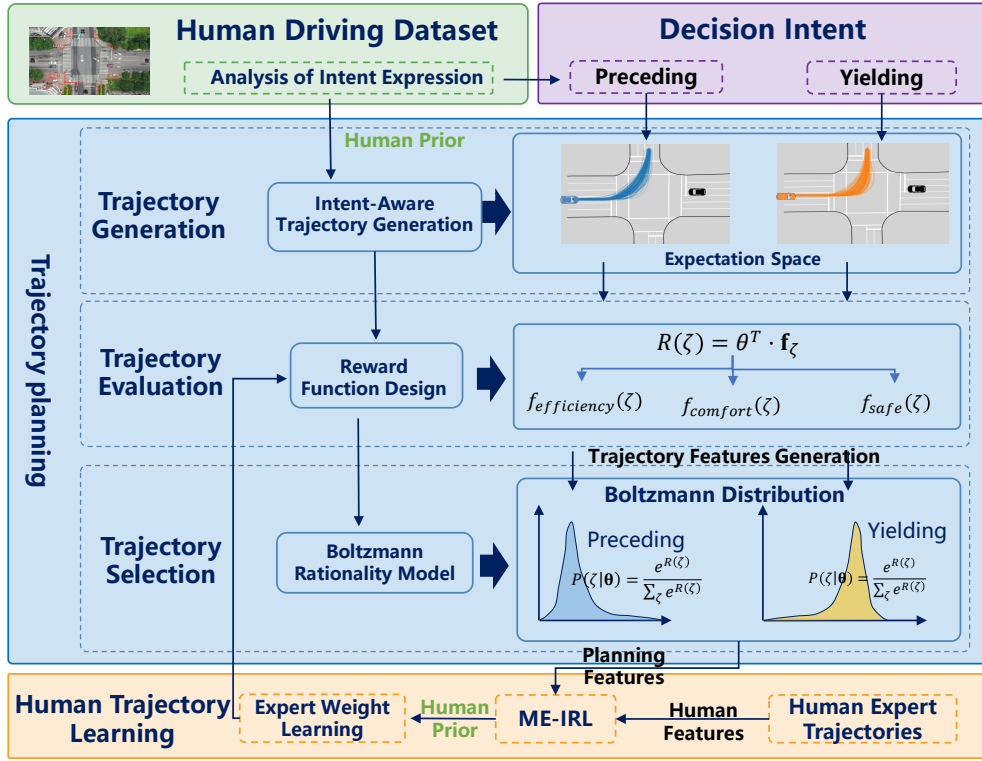
**Figure 1:** The overview framework of our methods.

- As illustrated in the right part of Fig. 2, the yielding trajectory maintains a straight course with no notable deviation in the upstream section, completing the yielding turn in a two-stage turn manner, i.e., "proceed straight first and then complete the turn," within the intersection.

By integrating the intent expression methods of these differing decisions, we will establish an expected trajectory space constraint capable of realizing implicit intent expression. This enables the trajectory to operate within the human expectation space, bestowing it with intent expression capabilities.
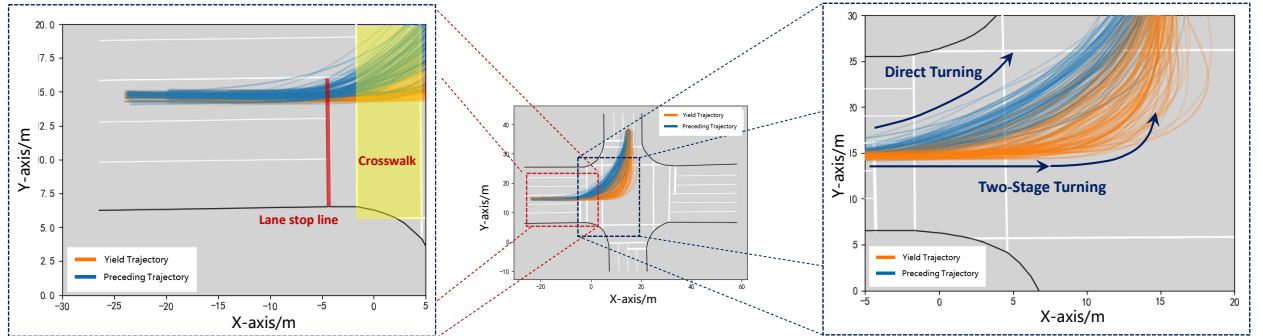


**Figure 2:** Trajectory characteristics under different HV decision-making

### 3.3. Generating Left-turn Candidate Trajectories at Intersections Based on the Frenet Coordinate System

To efficiently generate smooth and flexible trajectories at intersections, we assume that drivers typically make short-term plans based on the current interaction motion state, considering both lateral and longitudinal expected target points. Since the polynomial trajectory generation method can independently establish expressions for longitudinal and lateral displacement, and it has advantages like smooth trajectories, curvature continuity, and derivability of mathematical expressions, we employ the polynomial method to generate left-turn candidate trajectories based on the Frenet coordinate system.

The boundary conditions for trajectory generation using the polynomial method can be represented as follows:

$$\begin{cases} s\left(t_0\right) = s_0, \dot{s}\left(t_0\right) = v_{s0}, \ddot{s}\left(t_0\right) = a_{s0}, \dot{s}\left(T\right) = v_{sT}, \ddot{s}\left(T\right) = a_{sT} \\ l\left(t_0\right) = l_0, \dot{l}\left(t_0\right) = v_{l0}, \ddot{l}\left(t_0\right) = a_{l0}, l\left(T\right) = l_T, \dot{l}\left(T\right) = v_{lT}, \ddot{l}\left(T\right) = a_{lT} \end{cases} \tag{1}$$

where $s_0$, $v_{s_0}$, and $a_{s_0}$ represent the longitudinal displacement, velocity, and acceleration at the initial state, respectively. $l_0$, $v_{l0}$, and $a_{l_0}$ denote the lateral displacement, velocity, and acceleration at the initial state. $v_{s_T}$, $a_{s_T}$ represent the longitudinal velocity and acceleration at the terminal state. Similarly, $l_T$, $v_{l_T}$, and $a_{l_T}$ denote the lateral displacement, velocity, and acceleration at the terminal state.

Meanwhile, left-turning usually requires determining a reasonable sampling space to generate a set of trajectories, from which the optimal trajectory is selected. The sampling space for generating the trajectories is determined by the range of the terminal state parameters in the boundary conditions. To avoid the dimension explosion problem caused by too many parameters, we simplify the variables in the sampling space by setting the terminal moment longitudinal acceleration $a_{s_T} = 0$ and lateral acceleration $a_{l_T} = 0$. Then, we sample the trajectories over the time length $T$. The state variables that actually affect the trajectory sampling space are as follows:

$$\begin{cases} v_{sT} \in \left[v_{s0} - \Delta v_s, v_{s0} + \Delta v_s\right] \\ v_{lT} \in \left[v_{l0} - \Delta v_l, v_{l0} + \Delta v_l\right] \\ l_T = f_l\left(s_T\right) \\ T = f_T\left(s_T\right) \end{cases} \tag{2}$$

The change range of the longitudinal and lateral velocities, $v_{s_T}$ and $v_{l_T}$, under the trajectory terminal state is constrained based on the longitudinal and lateral velocities, $v_{s_0}$ and $v_{l_0}$, under the initial state of the trajectory. The lateral deviation $l_T$ under the terminal state is constrained based on the feature distribution of the expected trajectory space in subsequent research on human left-turn trajectory decisions. In determining the terminal position, we generate trajectories using an indefinite time length $T$. When the terminal position is uncertain, we generate trajectories using a fixed time length $T = 5s$. These two methods are applied to upstream road sections and internal intersection trajectory generation, respectively.

### 3.4. Constraint of Desired Trajectory Space Based on the Expression of Interaction Intent
#### 3.4.1. Desired Trajectory Space Constraint for Upstream Sections

It was analyzed in the previous sections that for left turn trajectories in the upstream sections, the preceding trajectories tend to make a pre-turn while the yielding trajectories continue to maintain straight-line motion along the lane. As shown in Fig. 3, the position, speed, and acceleration $(s_0, l_0, v_{s0}, v_{l0}, a_{s0}, a_{l0})$ are determined as the initial states of the scene (not trajectory). Based on the analysis of pre-turning behavior characteristics, it is determined that the trajectory will experience a lateral deviation and a change in heading angle when reaching the stop line. Hence, the pre-turning state of the trajectory at the stop line can be represented as $(s_{pre}, l_{pre}, \theta_{pre})$ Here, $s_{pre} = s_{stop}$, the longitudinal position of the stop line on the reference line, is known. The trajectory's lateral deviation $l_{pre}$ and heading

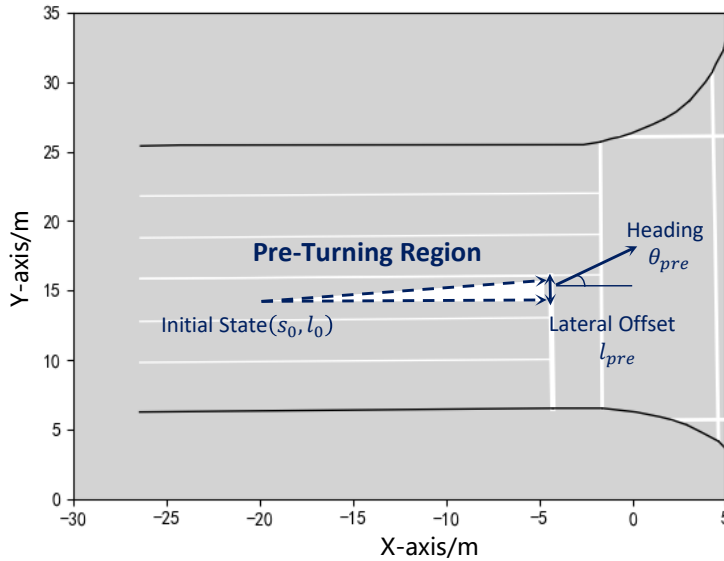angle $\theta_{pre}$ at the stop line are determined as follows:

$$
\begin{cases}
\theta_0 = \frac{v_{l0}}{v_{s0}} \\
\theta_{pre} = \frac{v_{lpre}}{v_{spre}} \\
l_{pre} \in \left[l_0, l_{stop}^{max}\right] \\
\theta_{pre} \in \left[\theta_0, \theta_{stop}^{max}\right]
\end{cases}
\tag{3}
$$

where $l_{stop}^{max}$ is the maximum trajectory deviation at that point. $\theta_{pre}$ determines the ratio of the lateral speed $v_{lpre}$ to the longitudinal speed $v_{spre}$ of the trajectory at the stop line. $\theta_{stop}^{max}$ is the maximum heading angle of the trajectory at that point. Lateral deviation and heading angle are determined through a uniformly distributed random sampling method.

Therefore, assuming the initial trajectory state at time $(s_t, l_t, v_{st}, v_{lt}, a_{st}, a_{lt})$, the trajectory sampling time length is set as $T_{pre}$. The pre-turning behavior characteristic state $(s_{pre}, l_{pre}, \theta_{pre})$ corresponds to the trajectory generation control terminal state $(s_{pre}, l_{pre}, v_{spre}, v_{lpre}, a_{spre}, a_{lpre})$ at this moment, where acceleration items are simplified to default zero. Then, when controlling the trajectory generation space variable duration $T_{pre}$, longitudinal speed $v_{spre}$, and lateral speed $v_{lpre}$, they can be determined as the following calculation results:

$$
\begin{cases}
T_{pre} \in \left[\frac{s_{pre}-s_t}{v_{st}} - \Delta T_{pre}, \frac{s_{pre}-s_t}{v_{st}} + \Delta T_{pre}\right] \\
\sqrt{v_{spre}^2 + v_{lpre}^2} \in \left[\sqrt{v_{st}^2 + v_{lt}^2} - \Delta v_{pre}, \sqrt{v_{st}^2 + v_{lt}^2} + \Delta v_{pre}\right] \\
\frac{v_{lpre}}{v_{spre}} = \theta_{pre}
\end{cases}
\tag{4}
$$

where $\Delta T_{pre}$ and $\Delta v_{pre}$ are control variables of trajectory duration and speed variation interval, respectively. The space variables of the yielding trajectory generation are consistent with the preceding trajectory, both controlling the trajectory sampling duration $T_{pre}$, longitudinal speed $v_{spre}$, and lateral speed $v_{lpre}$. The calculation method is consistent with the aforementioned one.



**Figure 3:** Space constraint of pre-steering trajectory in upstream section

### 3.4.2. Constraining the Desired Trajectory Space within Intersections

Analysis from previous sections determined that, within an intersection, left-turning vehicles following a trajectory that has the right of way will complete a rapid turn directly, while those yielding will adopt a two-phase turning method, maintaining a straight path before completing the turn.

As shown in Fig.2, different decision trajectories lead to distinctly different trajectory distribution spaces within the intersection due to the different turning methods employed. In this study, we model the relationship between the longitudinal displacement, s, and lateral displacement, l, of a trajectory using a continuously varying normal distribution. The mean and standard deviation of the preemptive trajectory's normal distribution are denoted as $\mu_{\text{preempt}}(s)$ and $\sigma_{\text{preempt}}(s)$ respectively. The mean and standard deviation of the yielding trajectory's normal distribution are denoted as $\mu_{\text{yield}}(s)$ and $\sigma_{\text{yield}}(s)$ respectively.

By determining the longitudinal displacement at time t, $s(t)$, we can deduce that $l_t \sim N(\mu(s(t)), \sigma(s(t)))$. Further, based on the 3-sigma rule of normal distribution, we can determine the range of the trajectory's lateral displacement at the end, $l_t \in [\mu(s(t)) - 2\sigma(s(t)), \mu(s(t)) + 2\sigma(s(t))]$.

Simultaneously, the longitudinal speed, $v_{sT}$, and the lateral speed, $v_{lT}$, at the end of the trajectory are determined within a changing interval based on the initial speeds $v_{s0}$ and $v_{l0}$. These speeds are uniformly sampled within the interval to determine $v_{sT}$ and $v_{lT}$. Assuming an initial trajectory state within the intersection at a particular time, $(s_0, l_0, v_{s0}, v_{l0}, a_{s0}, a_{l0})$, the longitudinal displacement at the end of the trajectory, $s_T$, can be computed given $v_{sT}$. With $s_T$ determined, the lateral displacement at the end of the trajectory, $l_T$, can be decided based on the established relationship between $s$ and $l$.

Therefore, the expected trajectory space within the intersection can be controlled and generated by the following variables:

$$
\begin{cases}
v_{sT} \in [v_{s0} - \Delta v_s, v_{s0} + \Delta v_s] \\
v_{lT} \in [v_{l0} - \Delta v_l, v_{l0} + \Delta v_l] \\
l_T \in [\mu(s_T) - 2\sigma(s_T), \mu(s_T) + 2\sigma(s_T)]
\end{cases}
\tag{5}
$$

where $\Delta v_s$ and $\Delta v_l$ are the control variables for the changing interval of longitudinal and lateral speeds respectively.

In order to ensure that all generated trajectories reflect normal kinematic characteristics exhibited by human driving and avoid collisions, we apply the following constraints to the dynamics of the trajectory:

$$
\begin{cases}
v_{min} \leq v_t \leq v_{max} \\
a_{min} \leq a_t \leq a_{max} \\
c_{min} \leq c_t \leq c_{max}
\end{cases}
\tag{6}
$$

where $v_t$, $a_t$, and $c_t$ represent the speed, acceleration, and curvature of a generated trajectory $\zeta_I^i \in \zeta_I$ at time $t \in [0, T]$. The $(v_{\max}, v_{\min})$, $(a_{\max}, a_{\min})$, and $(c_{\max}, c_{\min})$ are the minimum and maximum speeds, accelerations, and curvatures respectively, as determined from the statistical analysis of real human trajectories.

Furthermore, we conduct collision checks to ensure that the generated trajectories do not result in collisions with interactive objects at any time. This study includes the consideration of a safety margin, determining the safe bounding box $Boundingbox_{\text{safe}}$. The collision constraint can be expressed as:

$$
Boundingbox_{\text{safe}}^t (\text{left}) \cap Boundingbox_{\text{safe}}^t (\text{straight}) = \emptyset
\tag{7}
$$

where $Boundingbox_{\text{safe}}^t$ is the safe bounding box at time t, which is calculated by adding 0.5m to the front and rear spaces and 0.3m to the left and right spaces of the vehicle's boundary rectangle. If, at any point, a trajectory $\zeta_{I,K}^i \in \zeta_{I,K}$ fails to meet the collision constraint, that trajectory is removed from the set $\zeta_{I,K}$. The final set of desired trajectories that meet the collision constraints is $\zeta_{I,K,C}$.

## 3.5. Human Prior Learning Based on ME-IRL

In traditional trajectory planning research, after defining the trajectory feature function, the optimal trajectory selection is achieved by manually setting or experimentally determining the trajectory reward function feature weights.

However, in real-world interaction scenarios, it's challenging to accurately specify a reward function that captures all aspects of safe and efficient driving. Inverse Reinforcement Learning (IRL) can solve this problem. In this section, we consider the trajectory generation method expressing intentions and use ME-IRL to learn the human expert's trajectory behavior selection strategy under different decisions.

### 3.5.1. Maximum Entropy Inverse Reinforcement Learning

We assume that the total reward of a trajectory is a linear expression of the trajectory reward function, which is the weighted sum of selected features. Furthermore, we postulate that human drivers' preferences or behaviors under the same decision do not exhibit noticeable time variability and individual heterogeneity. Therefore, the total reward $R(\zeta)$ of a trajectory $\zeta$ can be expressed as:

$$R(\zeta) = \theta^T \cdot \mathbf{f}\zeta \tag{8}$$

where $\theta = [\theta_1, \theta_2, \cdots, \theta_K]$ is the weight vector, and $K$ depends on the number of trajectory reward functions $\mathbf{f}\zeta$. Given a human driving demonstration dataset $D = \{\zeta_1, \zeta_2, \cdots, \zeta_N\}$ composed of $N$ trajectories, the ME-IRL algorithm is used to infer the reward weight $\theta$, which can then be used to generate driving strategies that match the human expert demonstration trajectories.

Simultaneously, to simulate the randomness of human driver's trajectory selection, we employ the Boltzmann noise theory model to construct the candidate trajectory distribution. Under the Boltzmann distribution, all features expected to match the expert's demonstration have the maximum entropy principle, corresponding to the maximum entropy IRL. Therefore, the probability of a trajectory is proportional to the return of that trajectory,

$$P(\zeta | \theta) = \frac{e^{R(\zeta)}}{Z(\theta)} = \frac{e^{\theta^T \mathbf{f}\zeta}}{Z(\theta)} \tag{9}$$

where $\theta$ is the feature weight, $\mathbf{f}\zeta$ is the feature vector of trajectory $\zeta$, $P(\zeta | \theta)$ is the probability of trajectory $\zeta$, and $Z(\theta)$ is the partition function. As the partition function $Z(\theta)$ represents the integral sum of the rewards of all possible trajectories, it is challenging to directly calculate in continuous and high-dimensional spaces. We reduce the trajectory space to the previously researched and mined human prior expected trajectory space $\zeta_{I,K,C}$ and use a finite number of discretely generated feasible trajectories to approximate the partition function. Therefore, the probability expression of a trajectory yields that

$$P(\zeta | \theta) \approx \frac{e^{\theta^T \mathbf{f}\zeta}}{\sum_{i=1}^{N} e^{\theta^T \mathbf{f}\widetilde{\zeta}^i}} \tag{10}$$

where $\widetilde{\zeta}^i \in \zeta I, K, C$ is a generated trajectory with the same initial state as trajectory $\zeta$, $\mathbf{f}_{\widetilde{\zeta}^i}$ is the feature vector of trajectory $\widetilde{\zeta}^i$, and $N$ is the total number of generated trajectories.

The aim of maximum entropy IRL is to maximize the likelihood of expert demonstration trajectories by adjusting the feature weight $\theta$. The optimization objective function can be expressed as:

$$\max_{\theta} \mathcal{J}(\theta) = \max_{\theta} \sum_{\zeta \in D} \log P(\zeta|\theta) \tag{11}$$

where $D = \{\zeta_i\}i = 1^N$ represents the set of human expert demonstration trajectories. By substituting $P(\zeta | \theta)$ into the above equation, we obtain the optimized objective function $\mathcal{J}(\theta)$ as follows:

$$\mathcal{J}(\theta) = \sum_{\zeta \in D} \left[ \theta^T \mathbf{f}_\zeta - log \sum_{i=1}^{M} e^{\theta^T \mathbf{f}_{\widetilde{\zeta}^i}} \right] \tag{12}$$

The above equation can be optimized using gradient-based methods. The gradient of the optimization objective function $\mathcal{J}(\theta)$, $\nabla_\theta \mathcal{J}(\theta)$, can be expressed as follows:

$$\nabla_\theta \mathcal{J}(\theta) = \sum_{\zeta \in D} \left[ \mathbf{f}_\zeta - \sum_{i=1}^{M} \frac{e^{\theta^T \mathbf{f}_{\widetilde{\zeta}^i}}}{\sum_{i=1}^{M} e^{\theta^T \mathbf{f}_{\widetilde{\zeta}^i}}} \mathbf{f}_{\widetilde{\zeta}^i} \right] \tag{13}$$

The gradient can be viewed as the difference in feature expectations between human demonstration trajectories and generated trajectories:

$$\nabla_\theta \mathcal{J}(\theta) = \sum_{\zeta \in D} \left[ \mathbf{f}_\zeta - \sum_{i=1}^{M} P\left(\widetilde{\zeta}^i \mid \theta\right) \mathbf{f}_{\widetilde{\zeta}^i} \right] \tag{14}$$

Following the process outlined by (Huang et al. (2021)), we use a gradient ascent method to iteratively update the trajectory feature weights and compute the optimization objective until the loss converges. To prevent overfitting, we incorporate L2 regularization into the objective function $\mathcal{J}(\theta)$. Consequently, the gradient $\nabla_\theta \mathcal{J}(\theta)$ includes the difference in feature expectations plus a regularization term, as shown below:

$$\nabla_\theta \mathcal{J}(\theta) = \sum_{\zeta \in D} \left[ \mathbf{f}_\zeta - \sum_{i=1}^{M} P\left(\widetilde{\zeta}^i \mid \theta\right) \mathbf{f}_{\widetilde{\zeta}^i} \right] - 2\lambda\theta \tag{15}$$

where $\lambda > 0$ is the regularization parameter.

### 3.5.2. Reward Function Design

When designing the reward function, we considered three aspects: traffic efficiency, driving comfort, and interactive safety.

**(1) Traffic Efficiency**

We set the reward function for traffic efficiency as the loss in speed, which is the difference between the trajectory speed and the expected speed. We determine the traffic efficiency feature $f_{\text{efficiency}}(\zeta)$ of the trajectory $\zeta$ to be the mean speed loss at all times, represented as follows:

$$f_{\text{efficiency}}(\zeta) = -\frac{\sqrt{\sum_{t=0}^{T} \left(v_t - v_{\text{target}}\right)^2}}{T} \tag{16}$$

where $v_t$ is the scalar speed of the trajectory $\zeta$ at time $t$, $T$ is the total duration of trajectory sampling, and $v_{\text{target}}$ is the expected speed of the vehicle turning left.

**(2) Driving Comfort**

We select the jerk (rate of change of acceleration) to establish the comfort feature function $f_{\text{comfort}}(\zeta)$ for assessing whether the trajectory $\zeta$ is comfortable. This feature is determined as the mean of the jerk vector sum longitudinally and laterally at all times for the trajectory $\zeta$, calculated as follows:

$$\begin{cases} f_{\text{comfort}}(\zeta) = -\dfrac{\sum_{t=0}^{T} \sqrt{((Jerk_s^t)^2 + (Jerk_l^t)^2)}}{T} \\ Jerk_s^t = s'''(t) \\ Jerk_l^t = l'''(t) \end{cases} \tag{17}$$

where $Jerk_s^t$ and $Jerk_l^t$ are the longitudinal and lateral jerks of the trajectory $\zeta$ at time $t$ respectively. The reward function $f_{\text{comfort}}(\zeta)$ takes into account the smoothness in both the longitudinal and lateral directions.
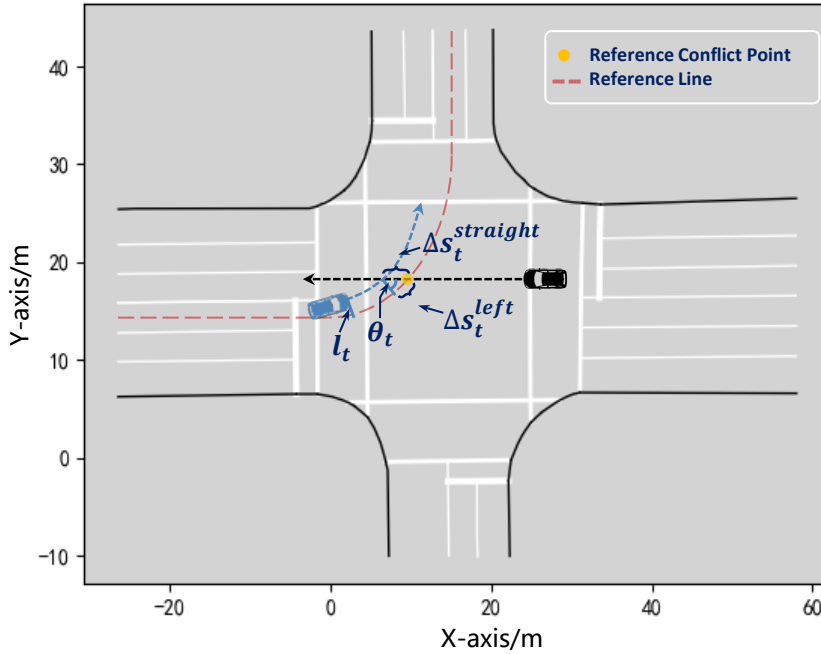
**(3) Interaction Safety**

We quantify interaction safety by calculating the time difference $\Delta TTCP$ between the interacting parties reaching the conflict point. $\Delta TTCP$ takes into account the relative relationships of the positions and speeds of the interacting parties, characterizing the time interval for both parties to leave the conflict point under the current state.

We deconstruct interaction safety into longitudinal progress and lateral deviation. The impact on interaction safety from the longitudinal progress at time $t$ is determined as the time difference for the interacting parties to reach the conflict point without considering the lateral deviation of the left-turning vehicle, denoted as $\Delta TTCP_{st}$. Its calculation is as follows:

$$\Delta TTCP_{st} = TTCP_{st}^{left} - TTCP_{st}^{straight} \tag{18}$$

$$\begin{cases} TTCP_{st}^{left} = \frac{s_t^{left} - s_{cp}^{left}}{v_{st}^{left}} \\ TTCP_{st}^{straight} = \frac{s_t^{straight} - s_{cp}^{straight}}{v_{st}^{straight}} \end{cases} \tag{19}$$

where $TTCP_{st}^{left}$ and $TTCP_{st}^{straight}$ are the times for the left-turning and straight-driving vehicles to pass the conflict point at the longitudinal level at time $t$ respectively. $s_t^{left}$ and $s_t^{straight}$ are the longitudinal positions of the left-turning and straight-driving vehicles at time $t$, $s_{cp}^{left}$ and $s_{cp}^{straight}$ are the longitudinal positions of the reference conflict points on the reference lines for the left-turning and straight-driving vehicles respectively, $v_{st}^{left}$ and $v_{st}^{straight}$ are the longitudinal velocities of the left-turning and straight-driving vehicles at time $t$.



**Figure 4**: Influence of lateral displacement of left-turning vehicle on interactive safety

After calculating $\Delta TTCP_{st}$, we further determine the feature function $f_{safe,s}(\zeta)$ on the longitudinal level of interaction safety as the mean of $\Delta TTCP_{st}$ at all times on trajectory $\zeta$, represented as follows:

$$f_{safe,s}(\zeta) = \frac{\sum_{t=0}^{T} |\Delta TTCP_{st}|}{T} \tag{20}$$

On the other hand, the impact of lateral deviation at time $t$ on interaction safety, as analyzed above, is determined by the impact of the left-turning vehicle's lateral deviation on the time difference for both interacting parties to reach the conflict point, denoted as $\Delta TTCP_{lt}$. Its calculation is as follows:

$$\Delta TTCP_{lt} = \Delta TTCP_{l}^{left} + \Delta TTCP_{l}^{straight} \tag{21}$$

where $\Delta TTCP_{l}^{left}$ and $\Delta TTCP_{l}^{straight}$ are the impacts caused by the left-turn vehicle's lateral deviation on the times for itself and the oncoming straight-driving vehicle to pass through the conflict point, respectively. Their calculations consider the variables as shown in Fig. 4 and are as follows:

$$\begin{cases} \Delta TTCP_l^{left} = \frac{l_t \cdot tan\theta_t}{v_{st}^{left}} \\ \Delta TTCP_l^{straight} = \frac{l_t \cdot cos\theta_t}{v_{st}^{straight}} \end{cases} \qquad (22)$$

where $l_t$ is the lateral deviation of the left-turn vehicle at time $t$, and $\theta_t$ is the angle produced by the expected trajectory of the oncoming straight-driving vehicle and the projection point of the expected conflict point on the reference line along the $l$ axis at time $t$.

After calculating $\Delta TTCP_{lt}$, we further determine the feature function $f_{safe,l}(\zeta)$ on the longitudinal level of interaction safety as the mean of $\Delta TTCP_{lt}$ at all times on trajectory $\zeta$, represented as follows:

$$f_{safe,l}(\zeta) = \frac{\sum_{t=0}^{T} |\Delta TTCP_{lt}|}{T} \qquad (23)$$

The interaction safety feature function of trajectory $\zeta$ is characterized from both longitudinal and lateral perspectives, represented as $f_{safe,s}(\zeta)$ and $f_{safe,l}(\zeta)$.

## 4. Experiment and Analysis

This section introduces the dataset used, describes the implementation details, and compares the experimental results while analyzing a few cases.

### 4.1. Dataset

We employ the SIND (Xu et al. (2022)), a drone dataset, for our analysis of human driver interaction behavior. The SIND dataset, curated by the SOTIF research team at Tsinghua University, represents an intersection dataset collected from a two-phase traffic signal-controlled intersection in Tianjin, China. This two-phase signal control scheme enables simultaneous movement of left-turning and straight-going vehicles, thereby resulting in pronounced interactions and frequent conflicts. After meticulous selection, we extract 268 sets of one-on-one interaction events between left-turning and straight-going vehicles. Out of these events, 132 instances feature a yielding left-turning vehicle, and 136 instances involve a left-turning vehicle proceeding first.
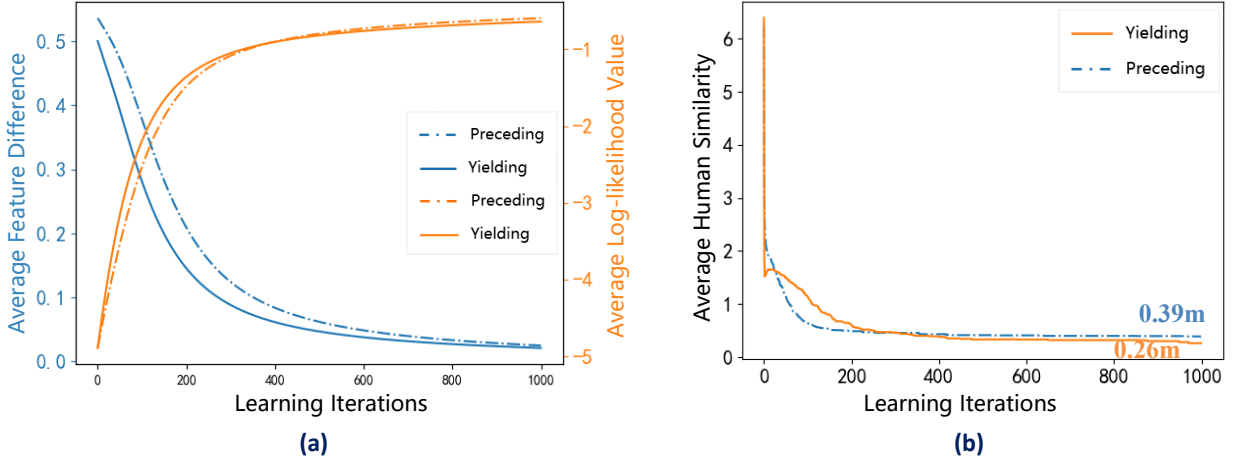
### 4.2. Implementation Details

In the inverse reinforcement learning process, we initially extract the left-turning trajectory of human experts from the SIND dataset. The duration of these trajectories varies between approximately 8 to 15 seconds. Beginning from the initial moment of the left-turning and straight-moving interaction event, we choose a trajectory segment of length T=5s every 0.5 seconds as the demonstration trajectory, denoted as $\zeta$, of the human expert. Based on the initial state of this trajectory segment, we generate a trajectory, denoted as $\zeta_{I,K,C}$, which considers intent expression. We randomly select 80% of the generated trajectories and their corresponding human expert demonstration trajectories as training data for learning the reward function, while the remaining 20% serve as test data.

During the training process, the trajectory generation sample space is determined by three control variables, namely, the longitudinal speed $v_{sT}$ at the end time, the lateral speed $v_{lT}$, and the lateral displacement $l_T$. Here, $v_{sT} \in [v_{s0} - \Delta v_s, v_{s0} + \Delta v_s]$, where $\Delta v_s = 3\,m/s$, and within this interval, we uniformly select 6 values for $v_{sT}$. $v_{lT} \in [v_{l0} - \Delta v_l, v_{l0} + \Delta v_l]$, where $\Delta v_l = 1\,m/s$, and within this interval, we uniformly select 5 values for $v_{lT}$. $l_T \in [\mu(s_T) - 2\sigma(s_T), \mu(s_T) + 2\sigma(s_T)]$, within this interval, we uniformly select 10 values for $l_T$. The time length of the trajectory planning is T=5s, and the time interval between each trajectory point is 0.1s. For the learning process of ME-IRL, we set the number of iterations E to 1000, with $\alpha = 0.05$ and $\lambda = 0.01$.

For comparative experiments with our proposed method, we select three kinds of trajectory planning methods, namely, **Frenet trajectory planning method** (Werling et al. (2010); Zhang et al. (2020)), **decision-based ME-IRL (DB-ME-IRL) trajectory learning method** (Dongjian et al. (2022)) that does not consider intent expression space, and **motion planning method based on potential field** (Xu, Ma and Sun (2019)). To compare the differences among these trajectory planning methods, we select the decision results of real interactive events as high-level decision information. More details about the method and results can be found at the site.[1]

---

[1]See https://shorturl.at/jqu35

**Figure 5:** The learning results of IRL, (a) Convergence Verification of IRL Results, (b) Accuracy Verification of IRL

## 4.3. Analysis of Inverse Reinforcement Learning Results

The training progression of the Maximum Entropy Inverse Reinforcement Learning (ME-IRL) model is depicted in Fig. 5. Fig. 5 (a) illustrates the iterative process that shows the average feature discrepancy between the reward function policy learned under yield and priority decisions and human drivers, as well as the change in the average log likelihood value of the human expert demonstration trajectory. These correspond to the gradient $\nabla_\theta \mathcal{J}(\theta)$ and the optimization target function $\mathcal{J}(\theta)$. The changing curves in the figure indicate that the average log likelihood value of the human expert demonstration trajectory steadily increases with the number of iterations under different decision trajectories, eventually reaching convergence.

We employ the Average Human Trajectory Similarity to gauge the closeness of the trajectory planning results to the human driving trajectory. We define the Average Human Trajectory Similarity (denoted as AHL) as the average final displacement error of the n trajectories with the highest selection probability in the generated trajectory distribution:

$$\text{AHL} = \frac{1}{n} \min_{i=1}^{n} |\hat{\zeta}_i(T) - \zeta(T)|_2 \tag{24}$$

where $\hat{\zeta}_i(T)$ $(i = 1, 2, ..., n)$ are the trajectories with the highest selection probability in the generated trajectory distribution, and $\zeta(T)$ is the real trajectory of human drivers. The curves in Fig.5 (b) respectively represent the accuracy of ME-IRL under the preceding and yielding decisions. The trajectory planning results of the priority and yield decisions have average errors of only 0.39m and 0.26m respectively when compared to the actual human trajectories, thus substantiating the effectiveness of our method.

## 4.4. Evaluation of Planning Results

We perform a comprehensive evaluation of our method across four dimensions: overall trajectory distribution, intent expression capability based on lateral offset, safety and efficiency of trajectory motion interaction features, and computational efficiency and learning effect.

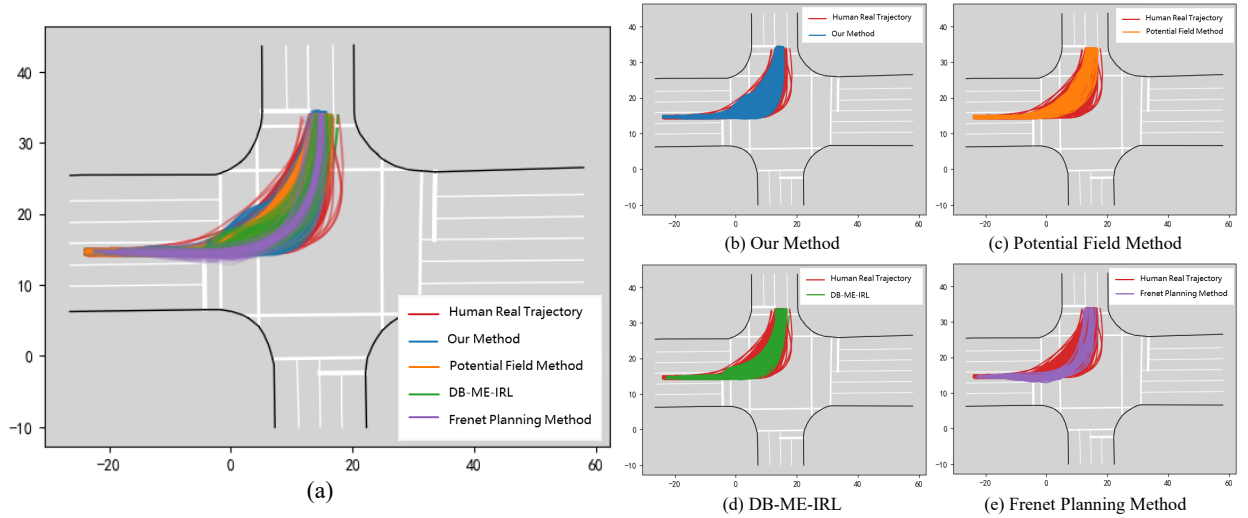### 4.4.1. Overall Trajectory Distribution

The social capability of trajectory planning can be analyzed from the overall distribution and spatial expansion of the trajectory planning. Trajectory planning methods with weaker sociality lack consideration of intent expression, resulting in a more single-choice trajectory selection and smaller trajectory space. The difference in trajectories under priority and yield decisions is not pronounced, making it difficult for interactive objects to determine their intent through trajectory behavior.

To intuitively compare the differences in the trajectory planning space of different methods, we drew comparisons between the trajectory planning distribution of four methods and the actual trajectory distribution, as shown in Fig.6.

**Table 1**
Comparison of trajectory coverage of different methods

| Distribution characteristic index | Trajectory planning method | | | | |
|---|---|---|---|---|---|
| | Real trajectory | Our method | Frenet Planning Method | DB-ME-IRL | Potential field method |
| Meta area coverage | 1066 | 824 | 576 | 710 | 737 |
| Ratio to true trajectory | - | 77% | 54% | 67% | 69% |



**Figure 6:** The planning trajectories from different methods. (a): Comparison of four trajectory methods with real trajectory distribution; (b) Our method; (c) Potential field method; (d) DB-ME-IRL; (e)Frenet planning method.

Through qualitative comparison, it can be found that the trajectory space distribution of our method is closest to the actual human trajectory distribution. The potential field method and DB-ME-IRL method's trajectory planning space distribution is slightly inferior to our method, with the Frenet planning method showing the most significant discrepancy from actual human trajectories.

We quantitatively analyze the differences in trajectory distribution using trajectory coverage. The specific calculation of trajectory coverage involves dividing the upstream entrance lane and the area inside the intersection into multiple subregions at 0.5m intervals in all directions. If there is a trajectory point in a subregion, it is determined that the subregion has been covered. Based on this concept, we calculated the coverage of actual human trajectories and the trajectories of the four planning methods, as shown in Table 1. The calculation results show that our method can achieve 77% coverage of the actual trajectory distribution. The potential field method and DB-ME-IRL method show a decline compared to our method, while the Frenet planning method performs the worst, covering only 54% of the actual trajectory distribution. Our method has improved the trajectory coverage rate by 12% compared to the potential field method.

### 4.4.2. Intent Expression Ability

When interacting with oncoming straight-line vehicles, left-turning vehicles usually adopt lateral offset behavior to clearly express their decision intent to the interactive object. Therefore, the intent expression ability of a trajectory can be characterized by analyzing the degree of lateral offset.

As shown in Fig.7, the red (blue) line segment is the center line of the virtual left-turn lane, and the gray line segments are the left and right boundary lines of the virtual left-turn lane. The offset to the left of the center line is negative, and the offset to the right is positive. The different color curves in Fig.7(a) - Fig.7(e) represent the normalized SL trajectories under different methods, and the corresponding colored areas represent the enclosed area between the normalized SL trajectories and the center line of the virtual lane. By comparing and observing Fig.7(a) - Fig.7(e), it can be seen that the offset degree exhibited by the SL trajectory of our method is closest to the real trajectory, while the

**Figure 7:** SL trajectory distribution under different planning methods
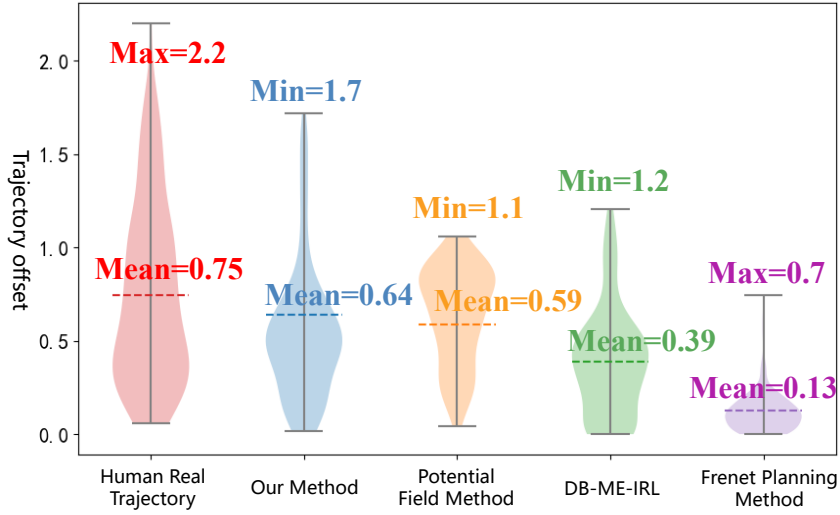
other three comparison methods all have significant differences with our method and the real trajectory. Although the potential field method can reasonably reproduce the behavior of priority trajectories quickly turning by significantly offsetting to the left of the centerline, it still cannot achieve a similar behavior to human real trajectories when planning yield trajectories, i.e., completing a longer straight line to indicate the intent to yield before turning. Both DB-ME-IRL and Frenet planning methods are inferior to our method in intent expression ability.

Further, we calculate the enclosed area $S_{SL}$ between the normalized SL trajectory and the center line of the virtual lane to represent the total lateral offset of a single trajectory. The total offset of each planned trajectory under different methods is calculated and drawn as a violin plot as shown in Fig.8. The real trajectory offset average is 0.75, our method's average is 0.64, the potential field's average is next at 0.59, and the smallest average is the Frenet planning method at 0.13. Our method's average is 85% of the real trajectory, an improvement of 8% compared to the comparison methods. In terms of maximum offset ability, the maximum offset of the real trajectory is 2.2, our method's maximum is 1.7, DB-ME-IRL's maximum is next at 1.2, and the smallest maximum is the Frenet planning method at 0.7. Our method's maximum offset is 77% of the real trajectory, an improvement of 42% compared to the comparison methods.
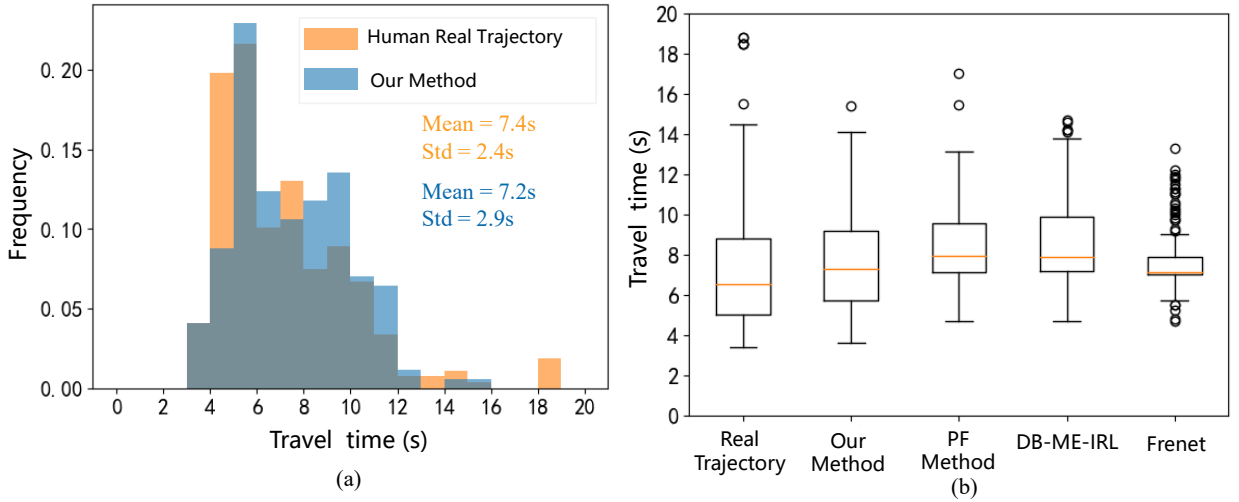
### 4.4.3. Motion Interaction Feature Analysis

We evaluate the safety and efficiency of different methods using two indicators: Post-Encroachment Time (PET) and Travel Time. We calculate the time required for the planned trajectories under different methods to enter and leave the intersection, as shown in Fig.9. Fig.9(a) compares the distribution of travel time within the intersection for the planned trajectories of our method and real trajectories. The average travel time for our method's planned trajectories is 7.4s, with a standard deviation of 2.4s. The average travel time for the real trajectories is 7.2s, with a standard deviation of 2.9s. Fig.9 (b) shows the travel time statistical results for the four different methods and the real trajectories, and we can see that the travel time of our method's planned trajectories is closest to the real trajectories.

Post-Encroachment Time (PET) has been proven to be suitable for representing the severity of conflicts during left turns. The PET distribution statistical results for different methods are shown in Fig.10. Fig.10 (a) compares the PET distribution for our method's planned trajectories and real trajectories. The average PET for our method's planned trajectories is 4.3s, with a standard deviation of 1.5s. The average travel time for real trajectories is 4.5s, with a standard deviation of 1.4s. Fig.10 (b) shows the PET statistical results for four different methods and the real trajectories. We can

**Figure 8:** Total amount of normalized SL trajectory offset from real trajectory and different methods



**Figure 9:** Comparison of travel time distribution of planning trajectories in intersections

see that our method's planned trajectories have the best PET performance. Specifically, both the DB-ME-IRL method and the Frenet planning method show a significant overall decrease in PET compared to our method, indicating an increased risk of trajectory interaction safety.

### 4.4.4. Computational Efficiency and Learning Performance

In order to validate the improvement of the trajectory planning algorithm's computational efficiency by the decision expectation space constraint, we select the average number of candidate trajectories, the average computation time (including feature computation and trajectory search), and learning performance. We conduct experimental comparisons with the DB-ME-IRL method, and the comparison results are shown in Table 2. The experimental device configuration used for comparison is an i7-8700 processor and 16GB memory. Our method reduces the average number of generated candidate trajectories by 52.5% after considering the expected trajectory space constraint and reduces the average computation time by 41.1%.
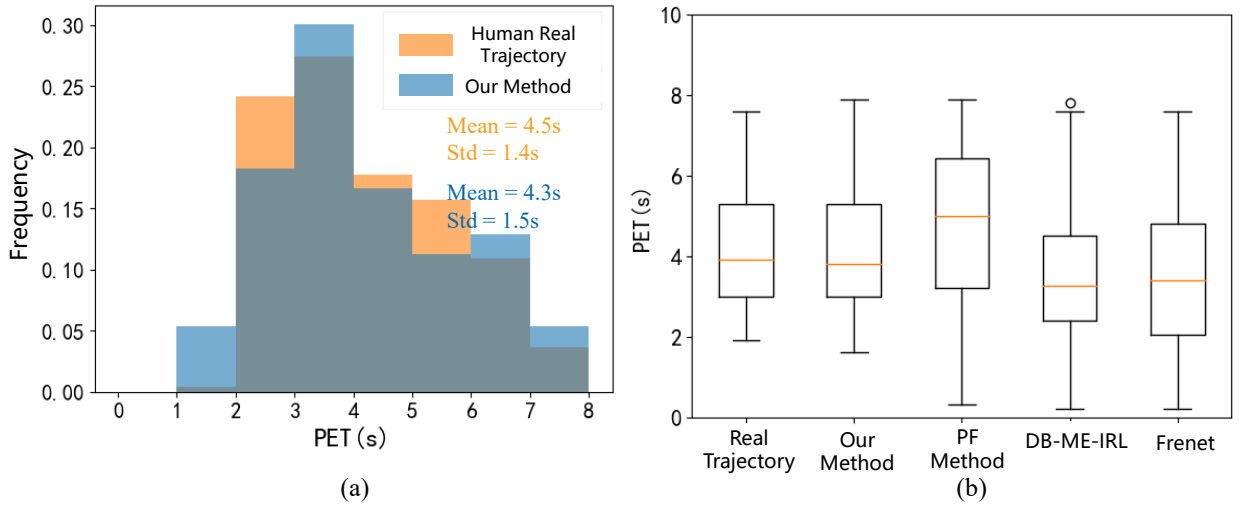
**Figure 10:** PET distribution of planning trajectory interaction process

**Table 2**
Trajectory Planning Efficiency and Effect Comparison

| Method | Average Number of Candidate Trajectories | Average Calculation Time | Learning Effect (Preceding) | Learning Effect (Yielding) |
|---|---|---|---|---|
| Our Method | 273 | 0.079s | 0.39 | 0.26 |
| Decision-based ME-IRL | 581 | 0.134s | 0.45 | 0.57 |

Further, we compared the trajectory learning performance through the Average Human Likeness (AHL) index. Our method and DB-ME-IRL were trained for 1000 generations simultaneously, and the process results are shown in Table 2 and Fig.11. After 1000 generations of simultaneous learning, the AHL of the trajectories under our method and DB-ME-IRL for priority decision are 0.39 and 0.45, respectively, and the AHL of the trajectories under our method and DB-ME-IRL for yielding decision are 0.26 and 0.57, respectively. The learning performance of our method has improved by 13% and 54% under the two types of decisions compared to the comparison method.
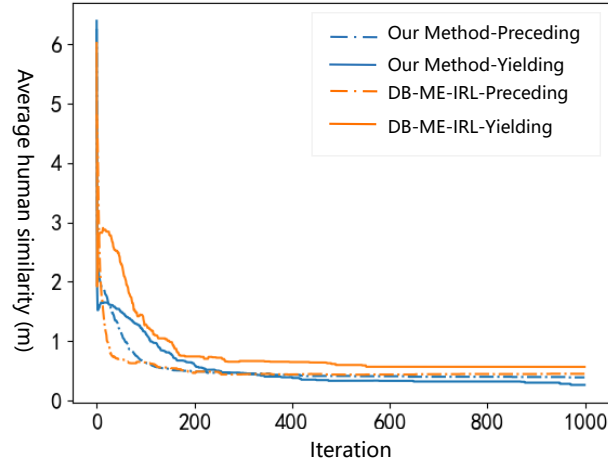
### 4.5. Case Analysis

To further substantiate whether our method resonates with human characteristics in expressing intentions during the interaction process, we have chosen two scenarios involving a left-turn vehicle yielding and proceeding first. We compared actual human trajectories, trajectories planned by our method, and trajectories planned by ME-IRL methods (Werling et al. (2010)). Animations of the two interaction cases can be accessed at the site.[2]
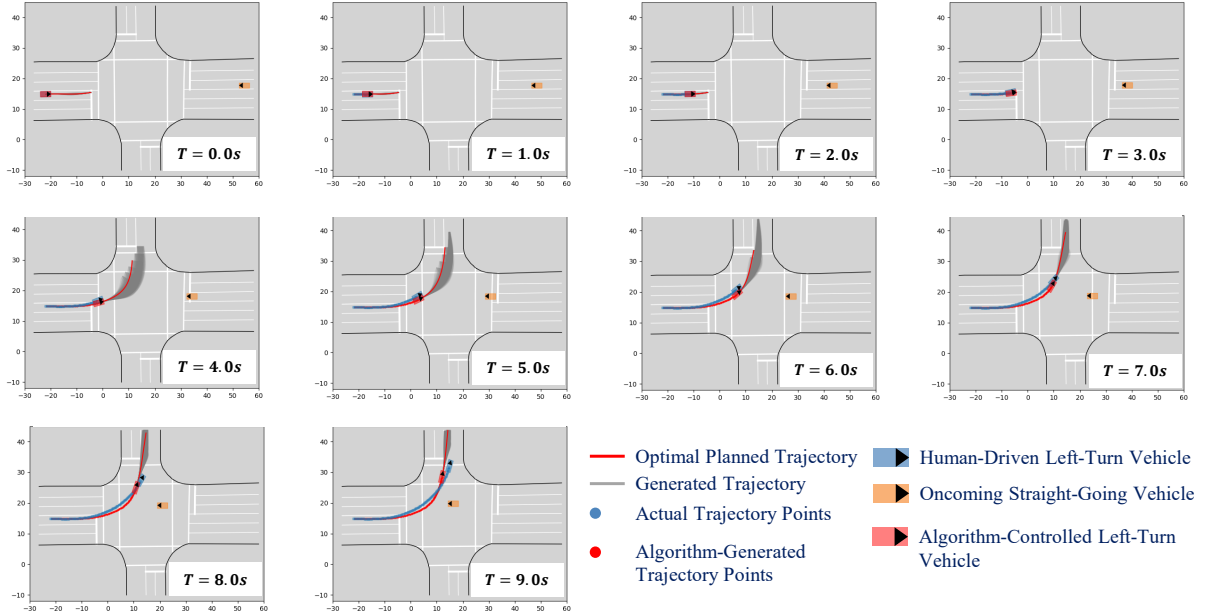
#### 4.5.1. Case 1: Left-turn Vehicle Proceeding

Fig. 12 illustrates a case where the left-turn vehicle is proceeding. At T=3s, the trajectory planned by our method begins to laterally offset prior to entering the intersection. Upon entering the intersection, the trajectory maintains a turning approach that passes through the conflict point early by significantly shifting left relative to the reference line, closely resembling the actual human trajectory. In contrast, the comparison method as shown in Figure 13 plans trajectories on both sides of the reference line, striving to minimize any lateral offset relative to the reference line. By comparison, it can be found that this method differs greatly from the actual human trajectory and fails to achieve intention communication at the trajectory behavior level.

[2]See https://shorturl.at/jqu35

**Figure 11:** Comparison of training process effects between DB-ME-IRL and our method



**Figure 12:** The trajectory of left turning from preceding case (Our Method).

Furthermore, our method reduces the trajectory generation search space significantly by adding the constraint of the expected space, whereas the comparison method uniformly generates candidate trajectories based on both sides of the reference line. Throughout the planning process, the maximum number of trajectories planned by our method is 300, while the comparison method needs to generate and search for up to 800 candidate trajectories. Our method reduces the search space by 62.5%.

We further quantitatively analyze the behavioral performance difference between the two methods by examining the changes in motion states during the interaction process, as demonstrated in Figure 14. Figure 14 (a) depicts the relationship between the heading angle of the left-turning vehicle and X. It reveals that our method exhibits a noticeable pre-turning behavior before the stop line of the lane, whereas the comparison method only commences turning after entering the intersection. Figure 14 (b) displays the relationship between the lateral offset of the left-turning vehicle
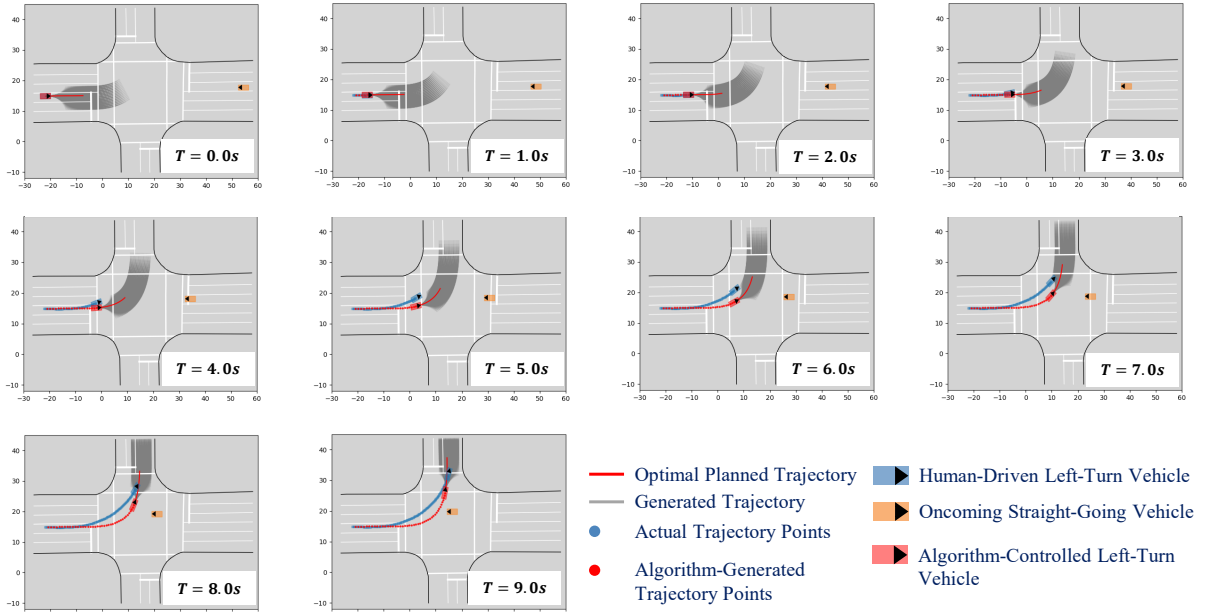
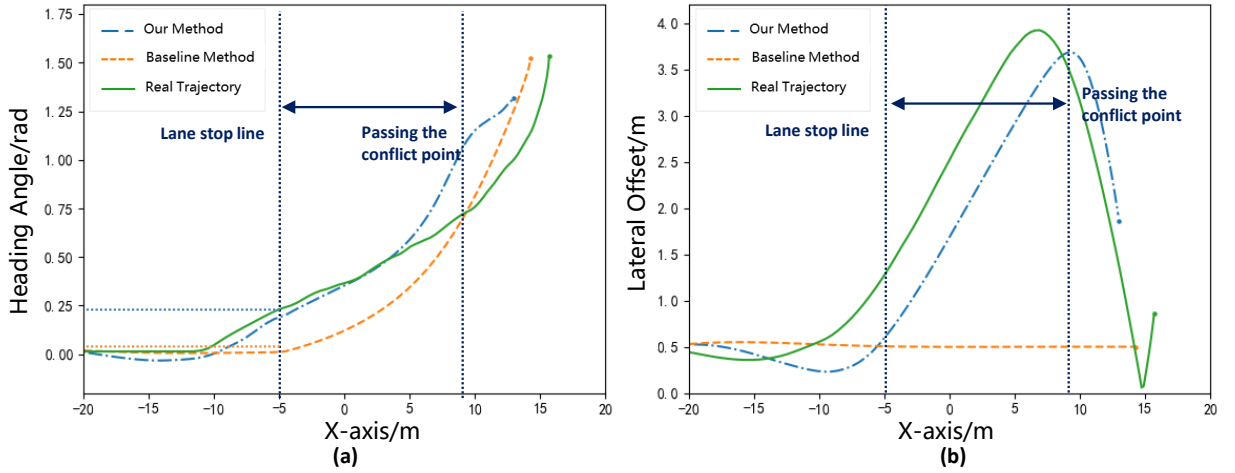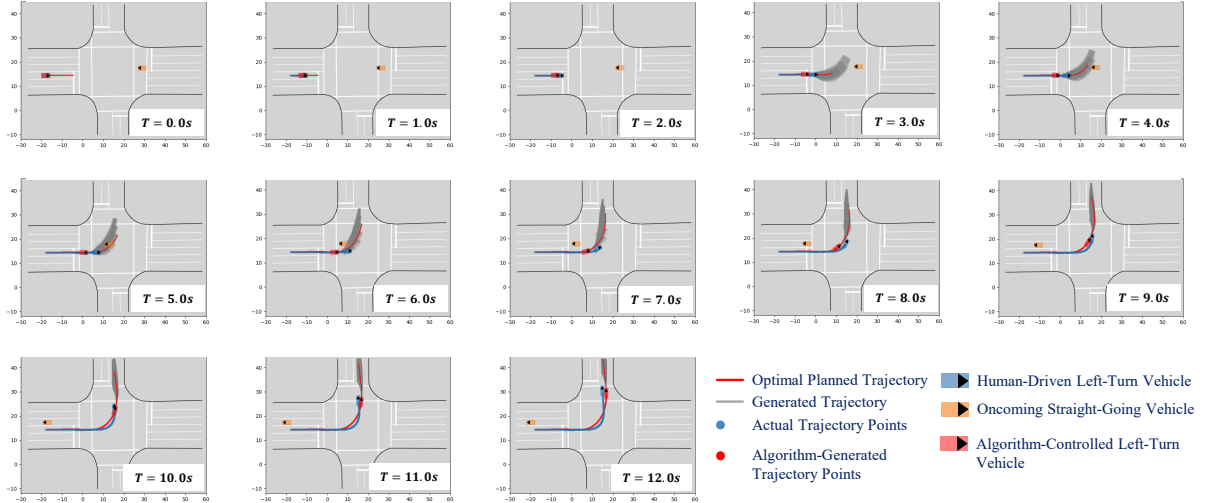**Figure 13:** The trajectory of left turning from preceding case (Baseline).



**Figure 14:** The interaction process of left turning from preceding case with different methods.

and X. It shows that our method passes through the conflict point first by significantly and swiftly shifting laterally during the turning process, whereas the comparison method adheres to the reference line, resulting in a trajectory with negligible lateral offset.

### 4.5.2. Case 2: Left-turn Vehicle Yielding

In the scenario where the left-turn vehicle yields, the trajectories planned by our method and the comparison method are illustrated as the red trajectory lines in Fig. 15 and Fig. 16, respectively. Our method, as depicted in Figure 15, can implement a two-stage turning approach of 'proceed straight then turn' under the yielding scenario when making a left turn. This both communicates its decision to yield to the interacting object and prevents the efficiency loss of stopping and restarting.

**Figure 15:** The trajectory of left turning from yielding case (Our Method).

Fig. 17(a) illustrates the relationship between the heading angle and X. If we consider a heading angle of 15° as a significant deviation, the real trajectory can maintain a straight phase up to x=10.7m, and our method can maintain a straight phase up to x=7.0m. After entering the intersection, it clearly yields and maintains a straight path. However, the comparison method can maintain a straight phase up to x=8.2m under the judgment standard of 15°. Nonetheless, the comparison method exhibits decision inconsistency; it decides to proceed first upon entering the intersection but fails, and then corrects its decision through the steering angle, resulting in a later turning time for the comparison method.

Figure 17(b) shows the relationship between the lateral offset and X. To quantify the difference between the planned trajectory and the actual trajectory during the straight phase, the area composed of the X-axis coordinates and the lateral offset is described as the degree of offset while going straight. During the straight phase of the real trajectory, our method's process offset degree is greater than that of the comparison method, thus appearing closer to the actual trajectory during the straight phase. The trajectory straight offset degree of our method is 15.6% different from the real trajectory, and the trajectory straight offset degree of the traditional method decreases by 22.1% compared to our method.

## 5. Conclusion

In an endeavor to narrow the divide between AVs and HVs and ensure that AVs can implicitly convey social intent in a manner understandable to HVs in mixed-traffic scenarios, we have introduced an innovative framework for socially-compliant trajectory planning that is robust in implicit intent expression at the unprotected left-turn scenarios.

Our proposed framework is organized into three components: trajectory generation, trajectory evaluation, and trajectory selection. The experimental results substantiate the efficacy of our framework, demonstrating a strong resemblance to actual human trajectories, considerable enhancements in intent expression, safety, and efficiency, along with improved computational efficiency and learning outcomes. Our method shows a 77% match with the actual trajectory distribution, an average offset of 85% from the real trajectory, an average travel time of 7.4 seconds within the intersection, and a decrease in the average computation time by 41.1%.

For future research, our research focus will expand the applicability of our trajectory planning methodology to include a broader range of interactive scenarios. Furthermore, we aim to confirm the effectiveness and scalability of our methodology through driving simulation experiments and real-vehicle interactions.
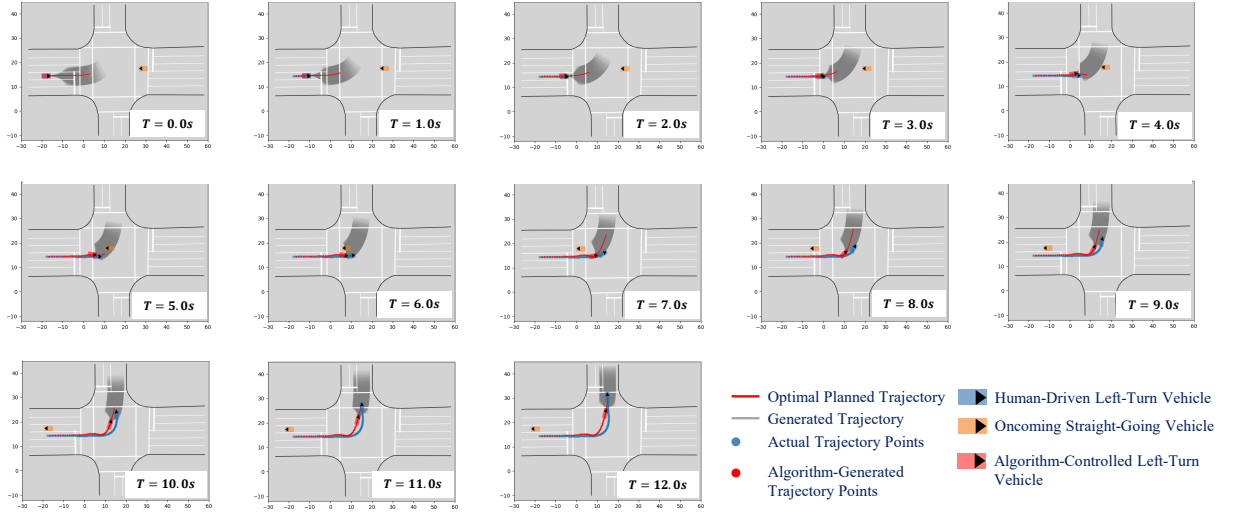
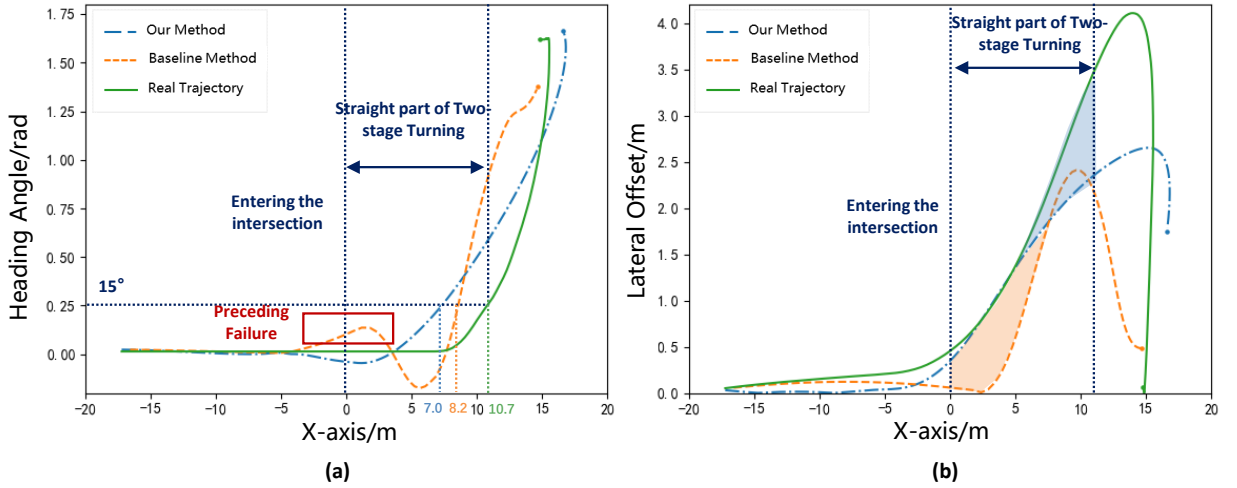**Figure 16:** The trajectory of left turning from yielding case (Baseline).



**Figure 17:** The interaction process of left turning from yielding case with different methods.

# References

Abbeel, P., Ng, A.Y., 2004. Apprenticeship learning via inverse reinforcement learning, in: Proceedings of the twenty-first international conference on Machine learning, p. 1.

Arora, S., Doshi, P., 2021. A survey of inverse reinforcement learning: Challenges, methods and progress. Artificial Intelligence 297, 103500.

De Ceunynck, T., Polders, E., Daniels, S., Hermans, E., Brijs, T., Wets, G., 2013. Road safety differences between priority-controlled intersections and right-hand priority intersections: Behavioral analysis of vehicle–vehicle interactions. Transportation research record 2365, 39–48.

Dongjian, S., Bing, Z., Jian, Z., Jiayi, H., Yanchen, L., 2022. Human-like behavior decision-making of intelligent vehicles based on driving behavior generation mechanism. Automotive Engineering 44, 1797. doi:10.19562/j.chinasae.qcgc.2022.12.001.

Hu, W., Deng, Z., Cao, D., Zhang, B., Khajepour, A., Zeng, L., Wu, Y., 2022. Probabilistic lane-change decision-making and planning for autonomous heavy vehicles. IEEE/CAA Journal of Automatica Sinica 9, 2161–2173.

Hu, X., Chen, L., Tang, B., Cao, D., He, H., 2018. Dynamic path planning for autonomous driving on various roads with avoidance of static and moving obstacles. Mechanical systems and signal processing 100, 482–500.

Huang, Z., Wu, J., Lv, C., 2021. Driving behavior modeling using naturalistic human driving data with inverse reinforcement learning. IEEE transactions on intelligent transportation systems 23, 10239–10251.

Lee, Y.M., Madigan, R., Giles, O., Garach-Morcillo, L., Markkula, G., Fox, C., Camara, F., Rothmueller, M., Vendelbo-Larsen, S.A., Rasmussen, P.H., et al., 2021. Road users rarely use explicit communication when interacting in today's traffic: implications for automated vehicles. Cognition, Technology & Work 23, 367–380.

Ma, Z., Sun, J., Wang, Y., 2017. A two-dimensional simulation model for modelling turning vehicles at mixed-flow intersections. Transportation Research Part C: Emerging Technologies 75, 103–119.

Ramachandran, D., Amir, E., 2007. Bayesian inverse reinforcement learning., in: IJCAI, pp. 2586–2591.

Wang, W., Wang, L., Zhang, C., Liu, C., Sun, L., et al., 2022. Social interactions for autonomous driving: A review and perspectives. Foundations and Trends® in Robotics 10, 198–376.

Werling, M., Ziegler, J., Kammel, S., Thrun, S., 2010. Optimal trajectory generation for dynamic street scenarios in a frenet frame, in: 2010 IEEE international conference on robotics and automation, IEEE. pp. 987–993.

Wulfmeier, M., Rao, D., Wang, D.Z., Ondruska, P., Posner, I., 2017. Large-scale cost function learning for path planning using deep inverse reinforcement learning. The International Journal of Robotics Research 36, 1073–1087.

Xu, Y., Ma, Z., Sun, J., 2019. Simulation of turning vehicles' behaviors at mixed-flow intersections based on potential field theory. Transportmetrica B: Transport Dynamics 7, 498–518.

Xu, Y., Shao, W., Li, J., Yang, K., Wang, W., Huang, H., Lv, C., Wang, H., 2022. Sind: A drone dataset at signalized intersection in china, in: 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC), IEEE. pp. 2471–2478.

Zhang, Y., Sun, H., Zhou, J., Pan, J., Hu, J., Miao, J., 2020. Optimal vehicle path planning using quadratic optimization for baidu apollo open platform, in: 2020 IEEE Intelligent Vehicles Symposium (IV), IEEE. pp. 978–984.

Zhao, J., Knoop, V.L., Sun, J., Ma, Z., Wang, M., 2023. Unprotected left-turn behavior model capturing path variations at intersections. IEEE Transactions on Intelligent Transportation Systems .

Zhou, D., Ma, Z., Zhang, X., Sun, J., 2022. Autonomous vehicles' intended cooperative motion planning for unprotected turning at intersections. IET Intelligent Transport Systems 16, 1058–1073.

Ziebart, B.D., Maas, A.L., Bagnell, J.A., Dey, A.K., et al., 2008. Maximum entropy inverse reinforcement learning., in: Aaai, Chicago, IL, USA. pp. 1433–1438.