

Collaborative Planning for Mixed-Autonomy Lane Merging

Shray Bansal¹, Akansel Cosgun², Alireza Nakhaei³ and Kikuo Fujimura³

Abstract—Driving is a social activity: drivers often indicate their intent to change lanes via motion cues. We consider mixed-autonomy traffic where a Human-driven Vehicle (HV) and an Autonomous Vehicle (AV) drive together. We propose a planning framework where the degree to which the AV considers the other agent's reward is controlled by a selfishness factor. We test our approach on a simulated two-lane highway where the AV and HV merge into each other's lanes. In a user study with 21 subjects and 6 different selfishness factors, we found that our planning approach was sound and that both agents had less merging times when a factor that balances the rewards for the two agents was chosen. Our results on double lane merging suggest it to be a non-zero-sum game and encourage further investigation on collaborative decision making algorithms for mixed-autonomy traffic.

I. INTRODUCTION

Driving is a social activity: drivers indicate their willingness to change lanes by subtle cues such as eye contact, or by not-so-subtle cues such as adjusting their speed and position [1]. There has been impressive demonstrations of Autonomous Vehicle (AV) technology [2]–[4], however one of the remaining challenges in this area is reading those cues to estimate the intentions of other agents as well as using cues to communicate the intentions of the AV. As AVs become commonplace, the situations where AV's and Human-driven Vehicles (HV) interact will increase. A number of issues in mixed-autonomy traffic need to be addressed before wide deployment, many posing interesting technical challenges.

Prior work focused on designing robust controllers for low-level tasks such as lane following and lane changing [3]–[7] either considered the other drivers as obstacles to avoid [5], [7] or did not consider their presence while modeling [6]. Close interactions with other agents in high-level tasks like deciding when to pass or change lanes [8], [9] require a more sophisticated model. For instance, a highway merge where there is a short distance to the next exit creates a situation where cars entering and exiting have to negotiate with one another to merge safely into their desired lanes. In this paper, we study a planning approach for navigating an autonomous vehicle for this challenging double lane merge in the presence of a human-driven car, see Fig. 1.

Recent progress in Reinforcement Learning (RL) has led to work exploring its application to autonomous driving [10], [11]. However, most of this progress was for the single-agent setting but driving is inherently multi-agent. [10] solve the

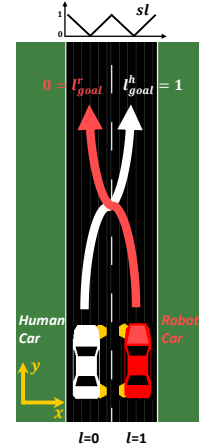


Fig. 1: Double lane merging problem. The autonomous and human-driven vehicles start together and aim to switch lanes. Successful merging requires signaling the agent's own intent while estimating the other's intention and then acting on it.

issue by only considering the other agent's actions through statistics of the traffic and although, [11] study the double lane merging problem in a multi-agent framework, they simplify the problem using expert knowledge. Similar to our approach of online planning, RL with Monte Carlo tree search [12] has been applied to autonomous driving in [13].

The idea of an autonomous agent acting deliberately to be predictable was formalized by [14]. Game theoretic models have been used to model human's adaptation to the robot [15] in the human-robot interaction domain. We are also influenced by research on collaborative planning with pedestrians for navigating a mobile robot in a social environment [16], [17]. Our work is most related to [18], which leverages the adaptability of the human by being cognizant of the effect the AV's actions have on the human driver. They show that the AV can influence a human driver such as making them slow down by moving in front of them.

However, while their approach maximizes the AV's own reward, it is not uncommon to see people slow down and let others into their lane while driving. Inspired by this idea of consideration and of leveraging the adaptability of HVs, we frame mixed-autonomy driving as a collaboration. In particular, we consider the double lane merging problem where the AV and HV start in adjacent lanes and must merge into each others' lanes in a limited road length. The actions of the AV are considered explicitly via collaborative planning with those of the HV to optimize a collective reward that combines the rewards of both of the agents. We validate our planning approach in simulation and conduct a user study where subjects can engage in interactions with our

¹S. Bansal is with Georgia Institute of Technology, Atlanta, GA, USA sbansal34@gatech.edu

²A. Cosgun is with Monash University, Clayton, VIC, Australia

³A. Nakhaei and K. Fujimura are with Honda Research Institute, Mountain View, CA, USA

autonomous agents. We study the effect of the AV having different levels of selfishness with users.

We make the following two contributions:

- A mixed-autonomy merge planning algorithm that is able to vary its consideration for other drivers with a single tunable selfishness factor.
- Demonstration that an AV with balanced consideration leads to better performance for both agents, suggesting that lane merging is a non-zero sum game.

II. PROBLEM

This paper assumes a discrete-time Markov Decision Process defined by,

$$M = \langle S, U, T, R \rangle, \quad (1)$$

where $S \subseteq \mathbb{R}^n$ is a n -dimensional state space, $U = \{1, \dots, K\}$ is the set of discrete actions the system can perform, $T : S \times U \rightarrow S$ is a deterministic transition function, and $R : S \rightarrow \mathbb{R}$ is a reward function.

A. State Representation

The state representation s provides full information of the vehicles present in the system. It contains the continuous lateral x and longitudinal positions y , as well as the longitudinal speed v of all cars including the AV.

$$s = (y_i, x_i, v_i)_{i \in \{1, \dots, k\}} \quad (2)$$

Where k is the number of cars in the environment. Since the focus of this study is on decision-making pertaining to intelligent interaction we limit the scope to the fully-observable setting.

B. Action Representation

The action representation u contains actions for each car $u = (u_i)_{i \in \{1, \dots, k\}}$. Each car has five control actions available to them at every time-step:

- **accelerate**: increases the speed by a constant factor.
- **decelerate**: decreases the speed by a constant factor.
- **stay**: maintains the speed.
- **turn-right**: moves the car in the positive latitudinal direction while reducing its longitudinal motion.
- **turn-left**: moves the car in the negative latitudinal direction while reducing its longitudinal motion.

The **right** and **left** actions are not permitted if the car is at the rightmost and leftmost edges of the road respectively. These discrete high-level actions were chosen because the article aims to study the high-level decision-making of the cars and so a controller is assumed.

C. Dynamics

The state at the next time-step is determined by applying a control action to each car in the system $s_{t+1} = T(s_t, u)$. These actions enable longitudinal and lateral control of the car and are applied instantaneously. The dynamics follow a simple deterministic model for each car. The accelerate,

decelerate and stay actions change the state of the corresponding vehicle in the following way.

$$v_{t+1} = v_t + a\Delta t \quad (3)$$

$$y_{t+1} = y_t + v_t\Delta t \quad (4)$$

$$x_{t+1} = x_t \quad (5)$$

Here, Δt is the length of the time-step and a is the acceleration. For the turning actions we introduce a fixed lateral velocity (v_{lat}),

$$v_t^x = \min\{v_t, v_{lat}\} \quad (6)$$

$$v_t^y = \sqrt{v_t^2 - v_t^{x2}} \quad (7)$$

$$y_{t+1} = y_t + v_t^y\Delta t \quad (8)$$

$$x_{t+1} = x_t + v_t^x\Delta t \quad (9)$$

$$v_{t+1} = v_t \quad (10)$$

D. Problem Domain

We use the two-lane scenario depicted in Fig. 1 where the HV and AV start in adjacent lanes and have the goal of merging into each other's lanes before the road ends. Here, if the cars change lanes independent (without regard) of each other it can cause a collision, which leads to the interesting challenge of planning while cognizant of the human driver's goals. The investigation is restricted to a two car scenario, with one human driver (H) and the other an AV, which will be referred to as the robot (R) car. We do so to reduce the effect of other interactions like HV-HV and AV-AV, and make the effects of the focus of our study - interactions between HVs and AVs - more clearly observable.

Even though our approach is extensible to multiple cars, we limit the system to two cars in this discussion and will apply a subscript of R or H to the AV and human, respectively. Thus, the state s is composed of the state for the human s_H and that of the robot s_R , where $s_H = (y_H, x_H, v_H)$ and $s_R = (y_R, x_R, v_R)$. We assume that the AV has noise-free access to the state at the current time s^t and the human's goal.

E. Reward

The goal of the AV is to reach its target lane (l_{goal}) while avoiding collisions and preferring to drive in the center of the lane. We encoded this high-level goal into a simple reward function that gets applied at every time-step. It has a high negative reward for collisions and a positive one for being in the goal-lane. To keep the car near the center of the lane, we add a small reward that decays exponentially with distance to the center.

$$r(s) = \begin{cases} -10, & \text{if collision} \\ \gamma e^{-sl} + (1 - \gamma), & \text{if } l = l_{goal} \\ 0, & \text{if } l \neq l_{goal} \end{cases} \quad (11)$$

Here, $l = \{0, 1\}$ and $sl = [0, 1]$ are the lane and sublane positions of the robot, explained in Fig. 1, and, $\gamma = 0.3$ is a constant used to make the agent favor the middle of the lane after it reaches the target lane. The positions l and sl can be easily determined from the latitudinal position x .

III. APPROACH

The goal of the AV is to perform a set of actions that optimize the reward it receives over time. To do so, it will plan for a finite sequence of actions, perform the first one and re-plan after receiving an observation. Let N be the length of the sequence, $\mathbf{u}_R = (u_R^t)_{t=0}^{N-1}$ be a sequence of the AV's actions and $\mathbf{u}_H = (u_H^t)_{t=0}^{N-1}$ be those of the human. Then the accumulated reward for the AV over the sequence starting at a state s^t is given by

$$R_R(s^t, \mathbf{u}_R, \mathbf{u}_H) = \sum_{i=0}^{N-1} r_R(T(s^{t+i}, u_H^{t+i}, u_R^{t+i})). \quad (12)$$

Here, r_R is the instantaneous reward defined in (11) and the next state $s^{t+1} = T(s^t, u_H^t, u_R^t)$ is computed using the dynamics T described in Section II-C. From a state s^t , the optimal set of actions for the AV, \mathbf{u}_R^* , can be found through the following optimization.

$$\mathbf{u}_R^* = \arg \max_{\mathbf{u}_R} R_R(s^t, \mathbf{u}_R, \mathbf{u}_H) \quad (13)$$

In the ideal scenario, \mathbf{u}_H would hold the HV's actual future actions, however, these are unknown. Some previous work presumed that \mathbf{u}_H depends only on the state, for instance, [19] uses a constant velocity model, but a more accurate model for a rational driver would also be influenced by the AV's future plan. We model this influence by assuming that the human driver plans in an action space that includes the AV as well and chooses to optimize a reward which includes the goals of both. This is based on the hypothesis that driving is not a zero-sum game and people drive with some consideration of each other's goals. In this system, we are indirectly influencing the human agent's actions by choosing the reward function for this joint planning system.

A. Planning and Prediction for Collaboration

Our collaborative planning approach determines the optimal set of actions by maximizing the joint reward R_J which is a weighted combination of the human and robot rewards,

$$R_J = \alpha R_R + (1 - \alpha) R_H. \quad (14)$$

Here, α is the selfishness factor which determines the relative importance of each reward in the collaboration. For instance, $\alpha = 1$ is an AV with no consideration for the human and $\alpha = 0$ is the opposite, while $\alpha = 0.5$ equally considers the goal of both agents. We do not assume apriori knowledge of α and plan to study its effect with the user study. The optimal plan for both agents under this model is determined by the following optimization.

$$\mathbf{u}_R^*, \mathbf{u}_H^* = \arg \max_{\mathbf{u}_R, \mathbf{u}_H} R_J(s, \mathbf{u}_R, \mathbf{u}_H) \quad (15)$$

Since we do not control the human's car, the set of actions \mathbf{u}_H^* is a prediction of their plan under the influence of the AV's plan \mathbf{u}_R^* . This influence is implicit in our model through choice of R_J . The AV is controlled by executing the first action from \mathbf{u}_R^* and the planning is repeated after every time-step. This planning is performed by an online search procedure described next.

B. Limited Horizon Tree Search

The optimization from (15) is implemented as an online tree-search. Here, the nodes are states containing both human and AV positions and velocities, from (2). The search starts at the current state and then simulates taking actions, terminating when either all the unpruned nodes, until maximum depth, are explored or if the time limit is reached. The time limit is set to be the length of the simulation time-step for real-time planning. The action leading to the maximally rewarding node is chosen and ties are broken randomly.

The algorithm is similar to A^* search [20] while also handling multiple unknown goal nodes and a simplified version is provided in `FIND_OPTIMAL_ACTION()` of Algorithm 1. The current state s^{init} and maximum depth t_{max} are used as input, a priority queue *open* keeps track of the states to explore and a set *closed* of those that have already been explored. States are ordered by f , which is the sum of their accumulated reward R and the heuristic. `GET_REWARD()` returns the reward at a state (11), while `rewards[s]` maps an explored state to its accumulated reward. `IS_TERMINAL()` checks if a given state is a terminal one, and `GET_HEURISTIC()` provides an upper bound by returning the maximum reward in the remaining horizon (or depth), keeping the search optimal. The `TIMER()` tracks the amount of time spent in the search and outputs true in case the time budget runs out. We also keep track of parent states and actions (not shown) to efficiently trace-back optimal actions. In case search is terminated prematurely, the state with maximal accumulated reward is chosen.

C. Selfish Baseline

Inspiring work from [18] had a similar formulation for mixed-autonomy driving. However, they determine the AV's plan \mathbf{u}_R^* from (13) by optimizing R_R only and the human's plan, \mathbf{u}_H , is determined by optimizing R_H , leading to the following nested optimization.

$$\mathbf{u}_R^* = \arg \max_{\mathbf{u}_R} R_R(s, \mathbf{u}_R, (\arg \max_{\mathbf{u}_H} R_H(s, \tilde{\mathbf{u}}_R, \mathbf{u}_H))), \quad (16)$$

To solve it they make some simplifying assumptions. They assume a turn-taking scenario where the robot plans first, also they give the human model access to \mathbf{u}_R , i.e. $\tilde{\mathbf{u}}_R = \mathbf{u}_R$. This leads to the human considering the robot's plan as fixed when planning her response. The reward R_R in our task depends only upon the state of the AV apart from avoiding collisions, reducing (16) to the selfish $\alpha = 1$ condition of our approach because the human-model here considers the plan \mathbf{u}_R as fixed. So, the only case that can cause a \mathbf{u}_R that optimizes $R_J = R_R$ in (15) to be suboptimal in (16) is when the optimal \mathbf{u}_H causes a collision. This, however, is prevented by the high negative reward of collisions for R_H .

This formulation works well in scenarios where R_R is designed to encode desired human behavior explicitly. For instance, they included the negative squared velocity of the human into R_R to reward the AV for slowing down the human. However, reward functions encoding the agent's own goals, rather than those of other drivers, are more intuitive,

Algorithm 1 Finite horizon action search

```

1: procedure FIND_OPTIMAL_ACTION( $s^{init}, t_{max}$ )
2:   Initialize  $open, closed, rewards$  as  $\emptyset$ 
3:    $t \leftarrow 0$ 
4:   Push  $((s^{init}, 0), GET\_HEURISTIC(t_{max} - 0))$  into
    $open$ 
5:    $R_{max} \leftarrow -\infty$ ,
6:   while  $open$  not empty do
7:     break if TIMER()
8:      $(s, t), f \leftarrow POP(open)$ 
9:     break if  $f < R_{max}$  ▷ Pruning
10:    add  $s$  to  $closed$ 
11:    if IS_TERMINAL( $s, t$ ) then
12:      if  $rewards[s] > R_{max}$  then
13:         $R_{max} \leftarrow rewards[s]$ 
14:      continue
15:     $t \leftarrow t + 1$ 
16:    for all  $s_+ \in GET\_NEIGHBORS(s)$  do
17:       $R \leftarrow rewards[s] + GET\_REWARD(s_+)$ 
18:      if  $(s_+ \in closed) \wedge (R < rewards[s_+])$  then
19:        continue
20:       $f_+ \leftarrow R + GET\_HEURISTIC(t_{max} - (t + 1))$ 
21:       $rewards[s_+] \leftarrow R$ 
22:      Push  $((s_+, t + 1), f_+)$  into  $open$ 
23:  Traceback from state with reward  $R_{max}$  to find
  optimal action  $(u_H^*, u_R^*)$ 
24:  return  $u_R^*$ 
25: procedure GET_HEURISTIC( $t_{remaining}$ )
26:  return  $t_{remaining} \cdot r_{max}$ 

```

and there, this formulation might lead to selfish behaviors as shown in our experiments.

D. Implementation Details

The human reward R_H was chosen to have the same simple form as the robot (11). Our results indicate that the model performs well even with this approximation of the human's reward. However, it can easily be replaced by a learned model in our framework in case sequences for learning this reward are available, like in [18]. The search algorithm 1, was implemented in Python. For the simulation we used the open source Simulation of Urban MObility (SUMO) [21] software and its in-built visualizer to render the cars. We used a discrete-time simulation with a time-step of 0.2 seconds, which also served as the time-limit for the search algorithm to return an action. The cars were allowed a maximum longitudinal speed of 30m/s and a lateral speed of 3m/s, with a lane width of 4m. If the cars start together with similar speeds it can take multiple seconds even for the optimal plan, to separate the cars enough so that a lane merge can commence safely. Thus, the time horizon of the search was kept at 6 seconds for interesting plans to develop. However, with a 0.2s time-step this would mean searching a tree of depth 30 and branching factor of 5^2 this leads to exploring 25^{30} states which is infeasible in 0.2s. To make

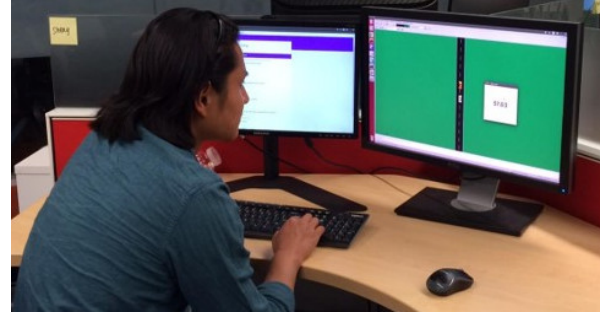


Fig. 2: A participant in our user study drives the car using the keyboard.

this search faster we increased the time-step to 1s in the planner, which is equivalent to repeating each action for five time-steps. This, combined with our choice of data-structures enabled us to typically complete the search in less than 0.2s.

IV. USER STUDY

A human-subject study was conducted to evaluate our method on the double lane merge scenario described earlier.

A. Independent Variables

We manipulated two aspects of this experiment:

- 1) **Road Length.** We used two lengths of the road $\{100, 200\}m$ referred to as *short-road* and *long-road*. This creates two levels of difficulty for the experiment where the more challenging *short road* scenario gives the HV and AV less time to complete their merges.
- 2) **Selfishness Factor.** The selfishness factor (α), defined earlier, was given six different values,

$$\alpha \in \{0.0, 0.2, 0.4, 0.6, 0.8, 1.0\}$$

which will be referred to as *0R*, *.2R*, *.4R*, *.6R*, *.8R*, *1R* respectively, where the level of the AV's selfishness increases with increasing α . We consider the *.4R*, *.6R* agents as fair and *1R* agent as selfish in the discussion.

Through these we explored a total of $2 \times 6 = 12$ conditions.

B. Hypotheses

The following hypotheses are tested.

- H1** *The selfishness factor will have an impact on the human's performance on the task.* In particular, we expect the less selfish versions of our algorithm to help improve the human driver's performance.
- H2** *Decreasing selfishness to a fair balance of rewards will not adversely effect the AV's performance.* We take *.4R* and *.6R* to be fair versions of our approach and expect that their performance is not significantly worse than the selfish *1R* agent.

C. Experimental Design

We recruited 21 participants, out of which 4 were female and 17 male, 19 had a valid driver's license (and were kept), and 20 were between the ages of 21 and 40 with one older.

The subjects had a top-down view of the environment and controlled the car using a keyboard, see Fig. 2. Each participant was instructed to drive their car to the goal

lane, in a safe and natural manner, before the end of the road. A small window displayed the distance to the end of the road. There was one other agent (AV) present in the scene that used indicator-lights to communicate its goal lane. We did not inform the participants that this agent was autonomously controlled to avoid biasing the perception of it. An unrecorded practice scenario was used to familiarize them with the simulation environment and the controls of their car. A within-subject design was used to mitigate the effects of inter-subject variability. Each subject performed 18 trials, each of which was uniformly sampled from one of the 12 conditions mentioned before. After every trial, a questionnaire gauged their interaction with the robot on a three-point scale. The sequences and responses of these trials were recorded and analyzed. To avoid making the task repetitive and losing the user's interest, we varied start lanes and the color of the robot's car through uniform sampling. The robot's initial speed was sampled from a normal distribution with a mean $15m/s$ and $3m/s$ standard deviation, while the human car's was kept constant at $15m/s$. The subjects were not informed of the variables being manipulated.

V. RESULTS AND ANALYSIS

A. Objective Measures

We studied the effect of α , using our 325 recorded trials on two measures, the average reward and the average merge-time for both agents. The results are shown in Fig. 3. By reward we refer to the individual rewards R_R and R_H as defined in Section 2 and merge-time refers to the amount of time it takes for an agent to reach their goal lane. Reward is an obvious choice for a measure of performance because the planner aims to optimize it for the AV and assumes that the human driver has a similar goal. We choose merge time here as a measure of task performance because the reward function was designed to encode it but also, because it might be a better measure for what the human drivers are optimizing. We believe that people do not prefer to take the risk of waiting to the end of the road to make a lane merge and will do it early if the opportunity presents itself. So, if the lane merging goes smoothly it should lead to lower average merge times.

Human. In Fig. 3 (a) human performance peaks when the $.4R$ condition is used and in (b) $.6R$ leads to lowest merge-time. The variation of these measures with the selfishness factor supports our hypothesis **H1** that the human performance is affected by the robot's collaboration. To test whether the cooperative condition $.6R$ outperforms the selfish $1R$ baseline, we compared their merge-times using an unpaired Student's t-test and found that the merge-time was significantly lower for the human ($p < 0.05$) in the $.6R$ condition.

Robot. A one-way ANOVA was conducted to test the variation of the AV's merge-time with different α and was found to be statistically significant ($F(5, 302) = 20.96$, $p < 0.0001$). This affirms that different levels of cooperation affect the performance of the AV. From Fig. 3 we observe that the two cooperative settings of $.4R$ and $.6R$ perform

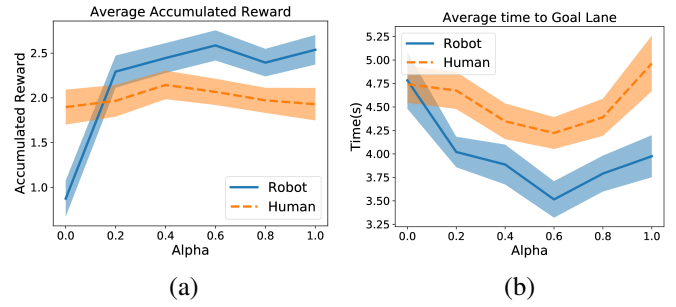


Fig. 3: Quantitative effect of α . (a) Plots the average reward accumulated by both agents and (b) shows the time it takes for a car to merge into its goal lane under all six conditions of the cooperation factor α . Fig. shows the mean as the line and the standard errors as the region. We observe that the fair conditions $.4R$, $.6R$ fare better than others including the selfish $1R$ for both agents. Note that higher rewards and lower merge-times relate to better performance.

equally-well or even better when compared to the selfish $1R$ baseline for both measures, supporting our hypothesis **H2**. The observation that $.6R$ improves the humans' ability to reach the goal is not surprising since it is more considerate towards the human's goal. However, it is surprising that the $.6R$ condition performs better even for the robot than optimizing its selfish goal. We attribute this to its possession of a more accurate model of the user's behavior. During planning, the $1R$ assumes that the human actively plans for the robot's goal and not her own. However, this is false and probably leads to it modifying its plan often making it suboptimal.

In Fig. 3, we also observe that the human is generally slower at lane merging than the robot. This is because the AV's reward favors the fastest lane merge but the subjects were not instructed to merge as quickly as possible, they were merely asked to do it before the road ends.

TABLE I: Failure Rate to Reach Goal Lane

	0R	.2R	.4R	.6R	.8R	1R
HV	15.2%	10.3%	7.5%	4.3%	9.5%	8.3%
AV	52.2%	10.3%	7.5%	2.1%	6.3%	6.7%

Failure Rate. Merge time was only computed for successful trials, i.e. trials where the car was able to successfully merge into its goal lane before the end of the road was reached. Table I shows the percentage to reach the goal for the human and the AV. Again, we found that the cooperative $.6R$ condition outperformed all others, including the $1R$, which means that both the human and AV were able to reach their goal lanes more often when the robot had a fair reward function. This lends more support to the hypotheses **H1** and **H2**. Every condition barring $0R$ achieves failures of $< 11\%$, supporting the argument that our framework is able to perform the task. We explain the reason for this high incidence of failures for $0R$ in section V-C.

B. Subjective Measures

After every trial, the participants were asked to answer the following two questions on a three-point scale.

Q1. Was the other car considerate of your goals?

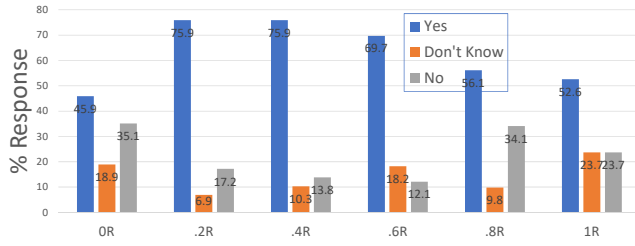


Fig. 4: Q1. Was it considerate of your goals?

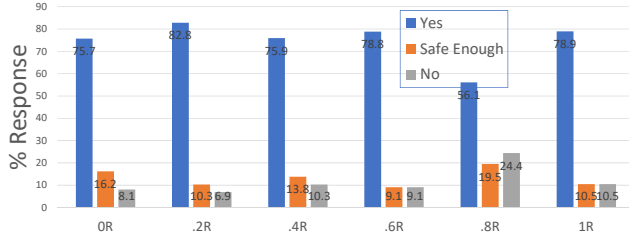


Fig. 5: Q2. Did the other car drive safely?

Q2. Did the other car drive safely?

For **Q1**, we explained that a considerate driver is one that changes their behavior for the benefit of other drivers. For instance, in the case of a car waiting to turn left at an intersection (with no traffic light on a U.S. road), being considerate at an intersection could mean reducing speed to give them time to make the turn or it might be better to speed up and pass them quickly. This notion of consideration is strongly influenced by the context, e.g. depending on relative speed, traffic condition, *etc.*, and can not be easily determined by a static rule.

Fig. 4 shows the response. Here, the three-point scale goes from "No", to "Don't Know", to "Yes". Excluding *OR*, the percentage of considerate trials decreases with increasing selfishness (α). Supporting the argument that less selfish agents are perceived more favorably in terms of consideration. A surprising observation is that the agent considered the least considerate is *OR*, the AV which solely optimizes the human's reward. We hypothesize that this agent is not perceived as rational and we believe that people reserve the considerateness attribute for rationally acting agents. The *OR* agent might be considered irrational because it does not act in its own self-interest when given the chance, for instance, even after the human has made a lane change and the agent is free to merge into its target lane it does not do so and seemingly takes random actions. This is illustrated in Fig. 6 and the next subsection gives more details on why this occurs. We believe that rationality is a prerequisite for being perceived considerate.

Q2 gauges whether the subject felt safe during the interaction on a three-point scale, from "No", to "Safe Enough", to a "Yes". We wanted the participants to use the notion of safety that they have developed by real-world driving and so did not provide examples of safety for the scenario. The results are grouped by the selfishness factor (α) in Fig. 5. Apart from *.8R*, every condition seems safe and there is no clear trend. This leads us to believe that safety is difficult to

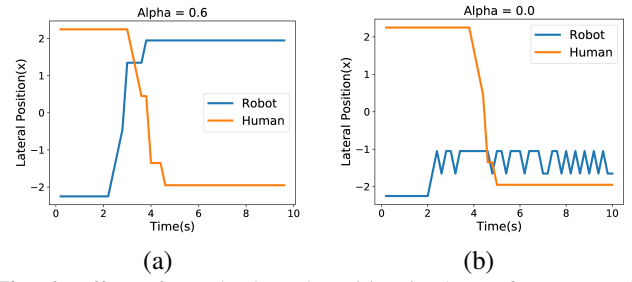


Fig. 6: Effect of α . The lateral position is shown for agents *.6R* and *OR* for two different trajectories in (a) and (b) respectively. (a) is a typical lane change sequence where both agents successfully steer to their goal lane. While in (b), the human again completes their merge successfully, but the AV takes random steering decisions making it appear irrational.

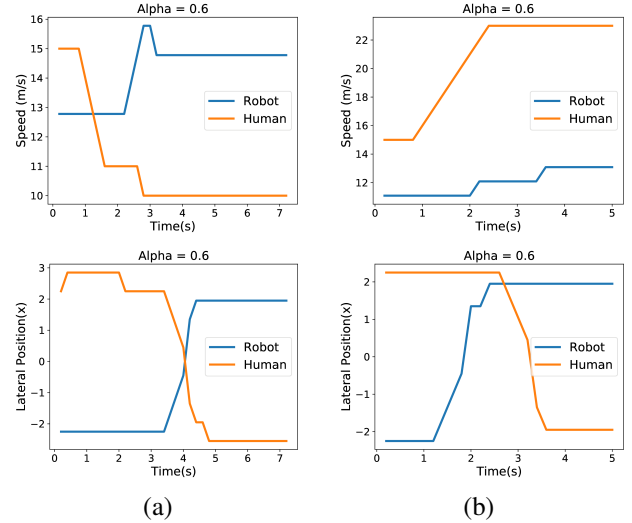


Fig. 7: Adaptation to human actions. The first and second rows plot longitudinal speed and lateral positions respectively. (a) and (b) are two trajectories. In (a), the human brakes and slows down so the AV accelerates to get ahead performing the merge afterwards. In (b), the human accelerates getting ahead of the AV allowing it to start steering early because it does not have to use the accelerate action first. Since the initial speeds and selfishness factor were the similar across columns, the difference in the AV's behavior can be attributed to the human's actions.

judge in this environment due to the detachment caused by the top-down view as opposed to a first-person view and in the future, a more involved simulation set-up might be able to recover it.

C. Qualitative Analysis

In this section we illustrate the behavior of our method through select examples. In particular, we explore the effects of the selfishness factor α and of the human driver's behavior in Figs. 6 and 7 respectively. Fig. 6 plots the lateral position (x) of both vehicles with respect to time for two trajectories. In (a), a typical scenario of a successful merge is shown, where both the human and AV start steering towards their goals and after a short time period, reach them. In (b), the unselfish *OR* steers towards its goal lane initially but then starts steering back and appears to be taking random actions. This is because *OR* optimizes only for R_H , i.e. for this AV's reward, while the human driver reaching their goal lane is ideal, their own lane does not matter. This leads the planner

to oscillate between choosing equally rewarding plans that may lead to either lane. The human in this scenario is aware of the AV's goal lane, but not of its peculiar reward function, and so might conclude from this behavior that the agent is irrational. We used this argument in Section V-B to explain why participants in our study did not find the *OR* agent to be considerate. This behavior explains the high failure rate for the AV, and its departure from the human's expectation might also explain the higher failure of the HV, for *OR* in Table I.

Fig. 7 highlights the capability of our approach to adapt to the human driver's behavior. Columns (a) and (b) are two trajectories that were performed with the same *.6R* selfishness factor and similar initial conditions but show different behavior due to the difference in human actions.

VI. CONCLUSIONS AND FUTURE WORK

We proposed a collaborative method for mixed-autonomy driving that optimizes a collective reward and showed that it can adapt to the human's behavior and find feasible solutions. To do this, we were able to use an approximate reward model for the humans due to their adaptability and did not require a large corpus of demonstrations to learn from or hand-code specific behaviors. Manipulating a single parameter in our framework gave rise to different levels of cooperation and we found that optimizing for a fair-share of the reward led to decreased merge times and failure rates. Providing evidence for the utility of collaboration in this domain and affirming our assumption that driving was a non-zero-sum game.

Future. We found that subjective notions of human safety are difficult to judge in this setup and hope to see future work utilizing Virtual Reality to create realistic environments to measure it. An interesting avenue would be to personalize the selfishness factor and estimate it in an on-line manner for every interaction. The idea of collaborative planning can be applied to other domains where agents interact to achieve personalized goals, for instance, factory robots sharing tools with a human worker. We think the idea of controlling the meta-behavior of an AV through variables like a selfishness factor can be used to personalize an autonomous car to the comfort of its rider. For example, some people might prefer an AV that yields easily to other drivers and is generally more considerate, but the same person may choose to change this behavior in situations where they are pressed for time.

ACKNOWLEDGMENT

We would like to thank colleagues at HRI - Priyam, Mustafa, Trevor, Reza and Sasha for insightful discussions. We would also like to acknowledge Amirreza, Ching-An, Brian, Byron, Himanshu, Kalesha, Sonia, Sid, Sangeeta, and Shubh for their valuable suggestions and comments.

REFERENCES

- [1] Y. Ba, W. Zhang, B. Reimer, Y. Yang, and G. Salvendy, "The effect of communicational signals on drivers' subjective appraisal and visual attention during interactive driving scenarios," *Behaviour & Information Technology*, vol. 34, no. 11, pp. 1107–1118, 2015.
- [2] A. Cosgun, L. Ma, J. Chiu, J. Huang, M. Demir, A. M. Anon, T. Lian, H. Tafish, and S. Al-Stouhi, "Towards full automated drive in urban environments: A demonstration in gomentum station, california," in *Intelligent Vehicles Symposium (IV)*, 2017 IEEE. IEEE, 2017.
- [3] C. Urmsen, J. Anhalt, D. Bagnell, C. Baker, R. Bittner, M. Clark, J. Dolan, D. Duggins, T. Galatali, C. Geyer *et al.*, "Autonomous driving in urban environments: Boss and the urban challenge," *Journal of Field Robotics*, vol. 25, no. 8, pp. 425–466, 2008.
- [4] J. Leonard, J. How, S. Teller, M. Berger, S. Campbell, G. Fiore, L. Fletcher, E. Frazzoli, A. Huang, S. Karaman *et al.*, "A perception-driven autonomous urban vehicle," *Journal of Field Robotics*, vol. 25, no. 10, pp. 727–774, 2008.
- [5] P. Petrov and F. Nashashibi, "Adaptive steering control for autonomous lane change maneuver," in *Intelligent Vehicles Symposium (IV)*, 2013 IEEE. IEEE, 2013, pp. 835–840.
- [6] C. Hatipoglu, U. Ozguner, and K. A. Redmill, "Automated lane change controller design," *IEEE transactions on intelligent transportation systems*, vol. 4, no. 1, pp. 13–22, 2003.
- [7] Z. Shiller and S. Sundar, "Emergency lane-change maneuvers of autonomous vehicles," *Journal of Dynamic Systems, Measurement, and Control*, vol. 120, no. 1, pp. 37–44, 1998.
- [8] M. Matthews, D. Bryant, R. Webb, and J. Harbluk, "Model for situation awareness and driving: Application to analysis and research for intelligent transportation systems," *Transportation Research Record: Journal of the Transportation Research Board*, 2001.
- [9] R. Sukthankar, S. Baluja, and J. Hancock, "Multiple adaptive agents for tactical driving," *Applied Intelligence*, vol. 9, no. 1, pp. 7–23, 1998.
- [10] D. Isele, R. Rahimi, A. Cosgun, K. Subramanian, and K. Fujimura, "Navigating occluded intersections with autonomous vehicles using deep reinforcement learning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018.
- [11] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "Safe, multi-agent, reinforcement learning for autonomous driving," *arXiv preprint arXiv:1610.03295*, 2016.
- [12] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, p. 484, 2016.
- [13] M. Bouton, A. Cosgun, and M. J. Kochenderfer, "Belief state planning for autonomously navigating urban intersections," in *Intelligent Vehicles Symposium (IV)*, 2017 IEEE. IEEE, 2017.
- [14] A. D. Dragan, K. C. Lee, and S. S. Srinivasa, "Legibility and predictability of robot motion," in *Proceedings of the 8th ACM/IEEE international conference on Human-robot interaction*. IEEE, 2013.
- [15] S. Nikolaidis, S. Nath, A. D. Procaccia, and S. Srinivasa, "Game-theoretic modeling of human adaptation in human-robot collaboration," in *Proceedings of the 2017 ACM/IEEE international conference on human-robot interaction*. ACM, 2017, pp. 323–331.
- [16] P. Trautman and A. Krause, "Unfreezing the robot: Navigation in dense, interacting crowds," in *Intelligent Robots and Systems (IROS)*, 2010 IEEE/RSJ International Conference on. IEEE, 2010.
- [17] A. Cosgun, E. A. Sisbot, and H. I. Christensen, "Anticipatory robot path planning in human environments," in *Robot and Human Interactive Communication (RO-MAN)*, 2016 25th IEEE International Symposium on. IEEE, 2016, pp. 562–569.
- [18] D. Sadigh, S. Sastry, S. A. Seshia, and A. D. Dragan, "Planning for autonomous cars that leverage effects on human actions," in *Robotics: Science and Systems*, 2016.
- [19] M. P. Vitus and C. J. Tomlin, "A probabilistic approach to planning and control in autonomous urban driving," in *Decision and Control (CDC)*, 2013 IEEE 52nd Annual Conference on. IEEE, 2013.
- [20] S. J. Russell, P. Norvig, J. F. Canny, J. M. Malik, and D. D. Edwards, *Artificial intelligence: a modern approach*. Prentice hall Upper Saddle River, 2003, vol. 2, no. 9.
- [21] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, "Recent development and applications of sumo-simulation of urban mobility," *International Journal On Advances in Systems and Measurements*, vol. 5, no. 3&4, pp. 128–138, 2012.