

Article

# On the Impact of the Rules on Autonomous Drive Learning

Jacopo Talamini <sup>1</sup>, Alberto Bartoli <sup>1</sup> , Andrea De Lorenzo <sup>1</sup> , and Eric Medvet <sup>1</sup> 

<sup>1</sup> Department of Engineering and Architecture, University of Trieste, Italy; jacopo.talamini@phd.units.it, {bartoli.alberto, andrea.delorenzo, emedvet}@units.it

\* Correspondence: emedvet@units.it

Version April 1, 2020 submitted to Appl. Sci.

**Abstract:** Autonomous vehicles raise many ethical and moral issues that are not easy to deal with and that, if not addressed correctly, might be an obstacle to the advent of such a technological revolution. These issues are critical because autonomous vehicles will interact with human road users in new ways and current traffic rules might not be suitable for the resulting environment. We consider the problem of learning optimal behavior for autonomous vehicles using Reinforcement Learning in a simple road graph environment. In particular, we investigate the impact of traffic rules on the learned behaviors and consider a scenario where drivers are punished when they are not compliant with the rules, i.e., a scenario in which violation of traffic rules cannot be fully prevented. We perform an extensive experimental campaign, in a simulated environment, in which drivers are trained with and without rules, and assess the learned behaviors in terms of efficiency and safety. The results show that drivers trained with rules enforcement are willing to reduce their efficiency in exchange for being compliant to the rules, thus leading to more overall safety.

**Keywords:** Reinforcement Learning, Self-driving Vehicles, Traffic Rules

## 1. Introduction

In the recent years autonomous vehicles have attracted a lot of interest from both industrial and research groups [1,2]. The reasons for this growth are the technological advancement in the automotive field, the availability of faster computing units, and the increasing diffusion of the so-called Internet of Things. Autonomous vehicles collect a huge amount of data from the vehicle and from the outside environment, and are capable of processing these data in real-time to assist decision-making on the road. The amount of collected information and the need for real-time computing make the design of the driving algorithms a complex task to carry out with traditional techniques. Moreover the sources of information may be noisy or may provide ambiguous information that could therefore affect negatively the outcome of the driving algorithm. The combination of these factors make it very hard, if not unfeasible, to define the driver behavior by developing a set of hand-crafted rules. On the other side, the huge amount of data available can be leveraged by suitable machine learning techniques. The rise of deep learning in the last decade has proven its power in many fields, including self-driving cars development, and enabled the development of machines that take actions based on images collected by a front camera as the only source of information [3], or even using a biological inspired event-driven camera [4].

The use of simulations and synthetic data [5] for training have allowed to assess neural networks capabilities in many different realistic environments and different degrees of complexity. Many driving simulators have been designed, from the low-level ones that allow the drivers to control the hand brake of their car [6], to higher-level ones, in which the drivers can control their car acceleration and

34 lane-change [7]. Some simulators model the traffic in an urban road network [8], some others model  
35 cars intersection access [9–12], or roundabout insertion [13].

36 In a near future scenario, the first autonomous vehicles on the roads will have to make decisions in  
37 a mixed traffic environment. Autonomous vehicles will have to be able to cope with radically different  
38 road agents, i.e., agents powered by machines capable of processing information way more faster than  
39 human drivers and human drivers that could occasionally take unexpected actions. There will hardly  
40 be a single authority to control each car in a centralized fashion and thus every autonomous vehicle  
41 will have to take decisions on its own, treating all the other road agents as part of the environment. It  
42 may very well be the case that current traffic rules do not fit a scenario with self-driving cars.

43 In this work, we investigate to which extent the traffic rules affect the drivers optimization process.  
44 The problem of finding the optimal driving behavior subjected to some traffic rules is highly relevant  
45 because it provides a way to define allowed behaviors for autonomous drivers, possibly without the  
46 need to manually craft those behaviors. A first approach for solving this problem consists of defining  
47 hard constraints on driver behavior and replacing forbidden actions with fallback ones [14]. Such  
48 an approach leads to drivers which are not explicitly aware of the rules. If those hard constraints  
49 were removed, driver behavior could change in unpredictable ways. Another approach consists in  
50 punishing behaviors that are not compliant with the rules, thus discouraging drivers from taking those  
51 behaviors again. In this work we investigate this second approach based on punishing undesired  
52 behaviors. In this scenario drivers have to learn the optimal behavior that balances a trade-off between  
53 being compliant with the rules and driving fast while avoiding collisions. A scenario in which drivers  
54 have the chance of breaking the rules is particularly relevant because it could address the complex  
55 ethics issues regarding self-driving cars in a more flexible way (those issues are fully orthogonal to our  
56 work, though).

57 We perform the optimization of the self-driving controllers using Reinforcement Learning (RL),  
58 which is a powerful framework used to find the optimal policy for a given task according to a  
59 trial-and-error paradigm. In this framework, we consider the possibility of enforcing traffic rules  
60 directly into the optimization process, as part of the reward function. Experimental results show that it  
61 is therefore possible to reduce unwanted behaviors with such approach.

## 62 2. Related works

63 The rise of Reinforcement Learning (RL) [15] as an optimization framework for learning artificial  
64 agents, and the outstanding results of its combination with neural networks [16], have recently reached  
65 many new grounds becoming a promising technique for the automation of driving tasks. Deep  
66 learning advances have proved that a neural network is highly effective in automatically extracting  
67 relevant features from raw data [17], as well as allowing an autonomous vehicle to take decisions  
68 based on information provided by a camera [3,4]. These approaches may not capture the complexity  
69 of planning decisions or predicting other drivers' behavior though, and their underlying supervised  
70 learning approach could be unable to cope with multiple complex sub-problems at once, including  
71 sub-problems not relevant to the driving task itself [18]. There are thus many reasons to consider a RL  
72 self-driving framework, which can tackle driving problems by interacting with an environment and  
73 learning from experience [18].

74 An example of an autonomous driving task implementation, based on Inverse Reinforcement  
75 Learning (IRL), has been proposed by [5]. The authors of the cited paper claim that, in such a large  
76 state space task like driving, IRL can be effective in extracting the reward signal, using driving data  
77 from experts demonstrations. End-to-end low-level control through a RL driver has been done by [6],  
78 in a simulated environment, based on the racing game TORCS, in which the driver has to learn full  
79 control of its car, that is steering, brake, gas, and even hand brake to enforce drifting. Autonomous  
80 driving is a challenging task for RL because it needs to ensure functional safety and every driver  
81 has to deal with the potentially unpredictable behavior of others [14]. One of the most interesting  
82 aspects of autonomous driving is learning how to efficiently cross an intersection, which requires to

83 provide suitable information on the intersection to the RL drivers [9], as well as correctly negotiating  
84 the access with other non-learning drivers and observing their trajectory [10,11]. Safely accessing to an  
85 intersection is a challenging task for RL drivers, due to the nature of the intersection itself, which may  
86 be occluded, and possible obstacles might not be clearly visible [12]. Another interesting aspect for  
87 RL drivers is learning to overtake other cars, which can be a particularly challenging task, depending  
88 on the shape of the road section in which the cars are placed [19], but also depending on the vehicles  
89 size, as in [20], where a RL driver learns to control a truck-trailer vehicle in an highway with other  
90 regular cars. The authors of [21,22] have provided extensive classifications of the AI state-of-the-art  
91 techniques employed in autonomous driving, together with the degrees of automation that are possible  
92 for self-driving cars.

93 Despite the engineering advancements in designing self-driving cars, a lack of legal framework  
94 for these vehicles might slow down their coming [23]. There are also important ethics and social  
95 considerations. It has been proposed to address the corresponding issues as an engineering problem,  
96 by translating them into algorithms to be handled by the embedded software of a self-driving car [24].  
97 This way the solution of a moral dilemma should be calculated based on a given set of rules or  
98 other mechanisms—although the exact practical details and, most importantly, their corresponding  
99 implications, are unclear. The problem of autonomous vehicles regulation is particularly relevant in  
100 mixed-traffic scenarios, as stated by [25,26], where human drivers may behave in unpredictable ways  
101 to the machines. This problem could be mitigated by providing human drivers with more technological  
102 devices to help them drive more like robotic drivers, but mixed traffic ethics certainly introduce much  
103 deeper and more difficult problems [25].

104 A formalization of traffic rules for autonomous vehicles is provided by [27], according to which  
105 a vehicle is not responsible for a collision if satisfying all the rules while colliding. Another driving  
106 automation approach based on mixed traffic rules is proposed in [28], where the rules are inspired by  
107 current traffic regulation. Traffic rules synthesis could even be automated, as proposed by [29], where  
108 a set of rules is evolved to ensure traffic efficiency and safety. The cited paper consider rules expressed  
109 by means of a language generated from a Backus-Naur Form grammar [30], but other ways to express  
110 spatio-temporal properties have been proposed [31,32]. Given the rules, the task of automatically  
111 finding the control strategy for robotics systems with safety rules is considered in [33], where the  
112 agents have to solve the task while minimizing the number of violated rules. AI safety can be inspired  
113 by humans, who intervene on agents in order to prevent unsafe situations, and then by training an  
114 algorithm to imitate the human intervention [34], thus reducing the amount of human labour required.  
115 A different strategy is followed by [35], where the authors define a custom set of traffic rules based  
116 on the environment, the driver, and the road graph. With these rules, a RL driver learns to safely  
117 make lane-changing decisions, where the driver's decision making is combined with the formal safety  
118 verification of the rules, to ensure that only safe actions are taken by the driver A similar approach is  
119 considered in [7], where the authors replace the formal safety verification with a learnable safety belief  
120 module, as part of the driver's policy.

### 121 3. Model

122 We consider a simple road traffic scenario in the form of a directed graph where the road sections  
123 are edges, and the intersections are vertexes. Each road element is defined by continuous linear space  
124 in the direction of its length, and an integer number of lanes. In this scenario a fixed number of cars  
125 move on the road graph according to their driver decisions for a given number of discrete time steps.

#### 126 3.1. Road graph

127 A *road graph* is a directed graph  $G = (S, I)$  in which edges  $E$  represent road sections, and vertices  
128  $I$  represent road intersections. Each road element  $p \in G$  is connected to the next elements  $n(p) \subset G$ ,  
129 with  $n(p) \neq \emptyset$ . Edges are straight one-way roads with one or more lanes. For each edge  $p$  it holds  
130 that  $n(p) \subset I$ . Vertices can be either turns or crossroads, they have exactly one lane, and are used to

131 connect road sections. For each vertex  $p$  it holds that  $n(p) \subset S$ , and  $|n(p)| = 1$ . Every road element  
 132  $p \in G$  is defined by its length  $l(p) \in \mathbb{R}^+$ , and its number of lanes  $w(p) \in \mathbb{N}, w > 0$ . We do not take  
 133 into accounts traffic light nor roundabout in this scenario.

### 134 3.2. Cars

135 A *car* simulates a real vehicle that moves on the road graph  $G$ : its position can be determined  
 136 at any time of the simulation in terms of the currently occupied road element, current lane. The car  
 137 movement is determined in terms of two speeds—i.e., the *linear speed* along the road element, and the  
 138 *lane-changing speed* along the lanes of the same element. At each time step, the car state is defined by the  
 139 tuple  $(p, x, y, v_x, v_y, s)$ , where  $p \in \{S, I\}$  is the current road element,  $x \in [0, l(p)]$  is the position on the  
 140 road element,  $y \in \{1, \dots, w(p)\}$  is the current lane,  $v_x \in [0, v_{\max}]$  is the linear speed,  $v_y \in \{-1, 0, 1\}$  is  
 141 the lane-changing speed, and  $s \in \{\text{alive}, \text{dead}\}$  is the status (time reference is omitted for brevity). All  
 142 the cars have the same length  $l_{\text{car}}$  and the same maximum speed  $v_{\max}$ .

143 At the beginning of a simulation, all cars are placed uniformly among the road sections, on all  
 144 the lanes, ensuring that a minimum distance exists between cars  $i, j$  on the same road element  $p_i = p_j$ ,  
 145 i.e., such that:  $|x_i - x_j| > x_{\text{gap}}$ . The initial speeds for all the cars are  $v_x = v_y = 0$ , and the status is  
 146  $s = \text{alive}$ .

At the next time steps, if the status of a car is  $s = \text{dead}$ , the position is not updated. Otherwise,  
 if the status is  $s = \text{alive}$ , the position of a car is updated as follows. Let  $(a_x^{(t)}, a_y^{(t)}) \in \{-1, 0, 1\} \times$   
 $\{-1, 0, 1\}$  be the driver action composed respectively of  $a_x^{(t)}$  accelerating action, and  $a_y^{(t)}$  lane-changing  
 action (see details below). The linear speed and the lane-changing speed at time  $t + 1$  are updated  
 accordingly with the driver action  $(a_x^{(t)}, a_y^{(t)})$  at time  $t$  as:

$$v_x^{(t+1)} = \min \left( v_{\max}, \max \left( v_x^{(t)} + a_x^{(t)} a_{\max} \Delta t, 0 \right) \right) \quad (1)$$

$$v_y^{(t+1)} = a_y^{(t)} \quad (2)$$

where  $a_{\max}$  is the intensity of the *instant acceleration*, and  $\Delta t$  is the discrete time step duration. The car  
 linear position on the road graph at time  $t + 1$  is updated as:

$$x^{(t+1)} = \begin{cases} x^{(t)} + v_x^{(t+1)} \Delta t & \text{if } v_x^{(t+1)} \Delta t \leq x_{\text{stop}}^{(t)} \\ v_x^{(t+1)} \Delta t - x_{\text{stop}}^{(t)} & \text{otherwise} \end{cases} \quad (3)$$

where  $x_{\text{stop}}$  is the distance ahead to the next road element, and is computed as:

$$x_{\text{stop}}^{(t+1)} = l \left( p^{(t+1)} \right) - x^{(t+1)} \quad (4)$$

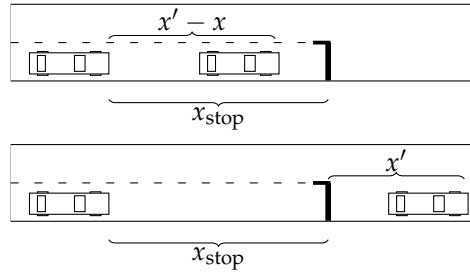
The car lane position at time  $t + 1$  is updated as:

$$y^{(t+1)} = \min \left( w(p^{(t+1)}), \max \left( y^{(t)} + v_y^{(t+1)}, 1 \right) \right) \quad (5)$$

The road element at time  $t + 1$  is computed as:

$$p^{(t+1)} = \begin{cases} p^{(t)} & \text{if } v_x^{(t)} \Delta t \leq x_{\text{stop}}^{(t)} \\ \sim U \left( n \left( p^{(t)} \right) \right) & \text{otherwise} \end{cases} \quad (6)$$

147 where  $U$  is the uniform distribution over the next road elements coming from  $p$ . In other words, when  
 148 exiting from an intersection, a car enters an intersection chosen randomly from  $n \left( p^{(t)} \right)$ .



**Figure 1.** Distance between cars in different cases.

Two cars collide, if the distance between them is smaller than the cars length  $l_{car}$ . In particular, for any cars  $(p, x, y, v_x, v_y, s)$ ,  $(p', x', y', v'_x, v'_y, s')$ , the status at time  $t + 1$  is updated as (we omit the time superscript for readability):

$$s = \begin{cases} \text{dead} & \text{if } (p = p' \wedge |x - x'| < l_{car}) \vee (p' \in n(p) \wedge x_{\text{stop}} + x' < l_{car}) \\ \text{alive} & \text{otherwise} \end{cases} \quad (7)$$

149 When a collision occurs, we simulate an impact by giving the leading car a positive acceleration of  
 150 intensity  $a_{\text{coll}}$ , while giving the following car a negative acceleration of intensity  $-a_{\text{coll}}$ , for the next  
 151  $t_{\text{coll}}$  time steps. Collided cars are kept in the simulation for the next  $t_{\text{dead}} > t_{\text{coll}}$  time steps of the  
 152 simulation, thus acting as obstacles for the alive ones.

### 153 3.3. Drivers

154 A *driver* is an algorithm that is associated to a car. Each driver is able to sense part of its car  
 155 variables and information from the road environment, and takes driving actions that affect its car state.  
 156 Every driver ability to see obstacles on the road graph is limited to the distance of view  $d_{\text{view}}$ .

#### 157 3.3.1. Observation

For the driver of a car  $(p, x, y, v_x, v_y, s)$ , the set of visible cars in the  $j$ -th relative lane, with  $j \in \{-1, 0, 1\}$ , is the union of the set  $V_{\text{same},j}$  of cars that are in the same segment and the same or adjacent lane and the set  $V_{\text{next}}$  of cars that are in one of the next segments  $p' \in n(p)$ , in both cases with a distance shorter than  $d_{\text{view}}$ :

$$V_{\text{same},j} = \left\{ (p', x', y', v'_x, v'_y, s') : p' = p \wedge 0 < x' - x \leq d_{\text{view}} \wedge y' = y + j \right\} \quad (8)$$

$$V_{\text{next}} = \left\{ (p', x', y', v'_x, v'_y, s') : p' \in n(p) \wedge x_{\text{stop}} + x' \leq d_{\text{view}} \right\} \quad (9)$$

158 We remark that the set of cars  $V_j = V_{\text{same},j} \cup V_{\text{next}}$  includes also the cars in the next segments: the  
 159 current car is hence able to perceive cars in a intersection, when in a segment, or in the connected  
 160 sections, when in an intersection, provided that they are closer than  $d_{\text{view}}$ .

The driver's observation is based on the concept of  *$j$ -th lane closest car*  $c_j^{\text{closest}}$ , based on the set  $V_j$  defined above. For each driver,  $c_j^{\text{closest}}$  is the closest one in  $V_j$ :

$$c_j^{\text{closest}} = \begin{cases} \arg \min_{(p', x', y', v'_x, v'_y, s') \in V_j} \mathbb{1}(p' = p)(x' - x) + \mathbb{1}(p' \neq p)(x_{\text{stop}} + x') & \text{if } V_j \neq \emptyset \\ \emptyset & \text{otherwise} \end{cases} \quad (10)$$

161 where  $V_j = V_{\text{same},j} \cup V_{\text{next}}$  and  $\mathbb{1} : \{\text{false}, \text{true}\} \rightarrow \{0, 1\}$  is the indicator function. Figure 1 illustrates  
 162 two different examples of  $j$ -th lane closest car, with  $j = 0$ . We can see that the  $c_j^{\text{closest}}$  might not exist for  
 163 some  $j$ , either if there is no car closer than  $d_{\text{view}}$  or if there is no such  $j$ -th lane.



Figure 2. Cars approaching intersections.

164 We define the *closeness* variables  $\delta_{x,j} \in [0, d_{\text{view}}]$ , with  $j \in \{-1, 0, 1\}$ , as the distances to the  $j$ -th  
 165 lane closest cars  $c_j^{\text{closest}}$ , if any, or  $d_{\text{view}}$ , otherwise. Similarly, we define the *relative speed* variables  
 166  $\delta_{v,j} \in [-v_{\text{max}}, v_{\text{max}}]$ , with  $j \in \{-1, 0, 1\}$ , as the speed difference of the current car w.r.t. the  $j$ -th lane  
 167 closest cars  $c_j^{\text{closest}}$ , if any, or  $v_{\text{max}}$ , otherwise.

168 At each time step of the simulation, each driver observes the distance from its car to the  
 169 next road element, indicated by  $x_{\text{stop}}$ , the current lane  $y$ , the current linear speed  $v_x$ , the status  
 170 of its vehicle  $s$ , the road element type  $e = \mathbb{1}(p \in S)$  its car is currently on, the closeness  
 171 variables  $\delta_{x,j}$  and the relative speed variable  $\delta_{v,j}$ . We define each driver *observation* as:  $o =$   
 172  $(x_{\text{stop}}, y, v_x, s, e, \delta_{x,-1}, \delta_{x,0}, \delta_{x,1}, \delta_{v,-1}, \delta_{v,0}, \delta_{v,1})$ , therefore  $o \in O = [0, l_{\text{max}}] \times \{1, w_{\text{max}}\} \times [0, v_{\text{max}}] \times$   
 173  $\{\text{alive}, \text{dead}\} \times \{0, 1\} \times [0, d_{\text{view}}]^3 \times [-v_{\text{max}}, v_{\text{max}}]^3$ .

### 174 3.3.2. Action

175 Each agent action is  $a = (a_x, a_y) \in A = \{-1, 0, 1\} \times \{-1, 0, 1\}$ . Intuitively  $a_x$  is responsible  
 176 for updating the linear speed in the following way:  $a_x = 1$  corresponds to accelerating,  $a_x = -1$   
 177 corresponds to breaking, and  $a_x = 0$  keeps the linear speed unchanged. On the other hand  $a_y$  is  
 178 responsible for updating the lane-position in the following way:  $a_y = 1$  corresponds to moving to the  
 179 left lane,  $a_y = -1$  corresponds to moving to the right lane, and  $a_y = 0$  to keeping the lane-position  
 180 unchanged.

### 181 3.4. Rules

182 A *traffic rule* is a tuple  $(b, w)$  where  $b : O \rightarrow \{\text{false}, \text{true}\}$  is the rule predicate, defined on the  
 183 drivers observation space  $O$ , and  $w \in \mathbb{R}$  is the rule weighting factor. The  $i$ -th driver *breaks* a rule at a  
 184 given time step  $t$  if the statement  $b$  that defines the rule is  $b(o_i^{(t)}) = 1$ . We define a set of three rules  
 185  $((b_1, w_1), (b_2, w_2), (b_3, w_3))$ , described in the next sections, that we use to simulate the real-world traffic  
 186 rules for the drivers. All the drivers are subjected to the rules.

#### 187 3.4.1. Intersection rule

188 In this road scenario we do not enforce any junction access negotiation protocol, nor we consider  
 189 traffic lights, and cars access interactions as in Figure 2. That is, there is no explicit reason for drivers to  
 190 slow down when approaching a junction, other than the chances of collisions with other cars crossing  
 191 the intersection at the same time. Motivated by this lack of safety at intersections, we define a traffic  
 192 rule that punishes drivers approaching or crossing an intersection at high linear speed.

In particular, the driver in road element  $p$  such that  $p \in I$  is an intersection, or equivalently  $p \in S$  and its car is in the proximity of an intersection, denoted by  $x_{\text{stop}} < 2l_{\text{car}}$ , breaks the intersection rule indicated by  $(b_1, w_1)$  if traveling at linear speed  $v_x > 10$ :

$$b_1(o) = \begin{cases} 1 & \text{if } (p \in I \vee x_{\text{stop}} < 2l_{\text{car}}) \wedge v_x > 10 \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

### 193 3.4.2. Distance rule

194 Collisions may occur when traveling with insufficient distance from the car ahead, since it is  
195 difficult to predict the leading car behavior in advance. For this reason we introduce a rule that  
196 punishes drivers that travel too close to the car ahead.

In particular, the driver observing  $c_0^{\text{closest}}$  closest car on the same lane, breaks the distance rule indicated by  $(b_2, w_2)$  if traveling at linear speed  $v_x$  such that the distance traveled before arresting the vehicle is greater than  $\delta_{x,0} - l_{\text{car}}$ , or, in other words:

$$b_2(o) = \begin{cases} 1 & \text{if } \delta_{x,0} - l_{\text{car}} < 2a_{\text{max}}v_x^2 \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

### 197 3.4.3. Right lane rule

198 In this scenario cars might occupy any lane on a road segment, without any specific constraint.  
199 This freedom might cause the drivers to unpredictably change lanes while traveling, thus endangering  
200 other drivers, who might not have the chance to avoid the oncoming collision. Motivated by this  
201 potentially dangerous behaviors, we define a rule that allows drivers to overtake when close to the car  
202 ahead, but punishes the ones leaving the right-most free lane on a road section.

In particular, the driver occupying road section  $p \in S$ , on non-rightmost lane  $y > 1$ , breaks the right lane rule indicated by  $(b_3, w_3)$  if the closest car on the right lane  $c_{-1}^{\text{closest}}$  is traveling at a distance  $\delta_{x,-1} = d_{\text{view}}$ :

$$b_3(o) = \begin{cases} 1 & \text{if } p \in S \wedge y > 1 \wedge \delta_{x,-1} = d_{\text{view}} \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

### 203 3.5. Reward

Drivers are rewarded according to their *linear speed*, thus promoting efficiency. All the cars involved in a collision, denoted by state  $s = \text{dead}$ , are then arrested after the impact, thus resulting in zero reward for the next  $t_{\text{dead}} - t_{\text{coll}}$  time steps, thus implicitly promoting safety. Each driver reward at time  $t$  is:

$$r^{(t)} = \frac{v_x^{(t)}}{v_{\text{max}}} - \sum_{i=1}^3 w_i b_i(o^{(t)}) \quad (14)$$

204 where  $w$  are the weights of the rules.

### 205 3.6. Policy learning

206 Each driver's goal is to maximize the *return* over a simulation, indicated by  $\sum_{t=0}^T \gamma^t r^{(t+1)}$ , where  
207  $\gamma \in [0, 1]$  is the discount factor, and  $T > 0$  is the number of time steps of the simulation. The driver  
208 *policy* is the function  $\pi_\theta : O \rightarrow A$  that maps observations to actions. We parameterize the drivers'  
209 policy in the form of a feed-forward neural network, where  $\theta$  is the set of parameters of the neural  
210 network. Learning the optimal policy corresponds to the problem of finding the values of  $\theta$  that  
211 maximize the return over an entire simulation. We perform policy learning by means of RL.

## 212 4. Experiments

213 Our goal is to experimentally assess the impact of the traffic rules on the optimized policies, in  
 214 terms of overall efficiency and safety. To this aim we defined 3 tuples, that are respectively the reward  
 215 tuple  $R$ , the efficiency tuple  $E$ , and the collision tuple  $C$ .

The *reward* tuple  $R \in \mathbb{R}^{n_{\text{car}}}$  is the tuple of individual rewards collected by the drivers during an episode, from  $t = 0$  to  $t = T$ , and is defined as:

$$R = \left( \sum_{t=0}^T r_1^{(t)}, \dots, \sum_{t=0}^T r_{n_{\text{cars}}}^{(t)} \right) \quad (15)$$

The *efficiency* tuple  $E \in \mathbb{R}^{n_{\text{car}}}$  is the tuple of sums of individual instant *linear speed*  $v_x$  for each driver during an episode, from  $t = 0$  to  $t = T$ , and is defined as:

$$E = \left( \sum_{t=0}^T v_{x_1}^{(t)}, \dots, \sum_{t=0}^T v_{x_{n_{\text{cars}}}}^{(t)} \right) \quad (16)$$

The *collision* tuple  $C \in \mathbb{N}^{n_{\text{car}}}$  is the tuple of individual collisions for each driver during an episode, from  $t = 0$  to  $t = T$ , and is defined as:

$$C = \left( \sum_{t=0}^T \mathbb{1}\{s_1^{(t-1)} = \text{alive} \wedge s_1^{(t)} = \text{dead}\}, \dots, \sum_{t=0}^T \mathbb{1}\{s_{n_{\text{cars}}}^{(t-1)} = \text{alive} \wedge s_{n_{\text{cars}}}^{(t)} = \text{dead}\} \right) \quad (17)$$

216 Each  $i$ -th element  $c_i$  of this tuple is defined as the number of times in which the  $i$ -th driver change its  
 217 car status  $s_i$  from  $s_i = \text{alive}$  to  $s_i = \text{dead}$  between 2 consecutive time steps  $t - 1$  and  $t$ .

218 We considered 2 different driving scenarios in which we aim at finding optimal policy parameters,  
 219 respectively “no-rules” in which traffic rules weighting factors are  $w_1 = w_2 = w_3 = 0$ , such that  
 220 drivers are not punished for breaking the rules, and “rules” in which traffic rules weighting factors are  
 221  $w_1 = w_2 = w_3 = 1$ , such that drivers are punished for breaking the rules, and all the rules have the  
 222 same relevance.

223 Moreover, we considered 2 different collision scenarios:

- 224 (a) cars are kept with status  $s = \text{dead}$  in the road graph for  $t_{\text{dead}}$  time steps, then they are removed;
- 225 (b) cars are kept with status  $s = \text{dead}$  in the road graph for  $t_{\text{dead}}$  time steps, then their status is  
 226 changed back into  $s = \text{alive}$ .

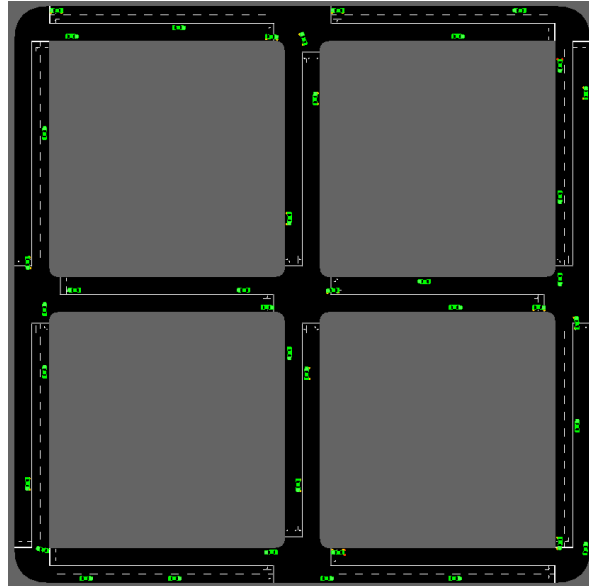
227 The rationale for considering the second option is that the condition in which we remove collided cars  
 228 after  $t_{\text{dead}}$  time steps may not be good enough for finding the optimal policy. This assumption could  
 229 ease the task of driving for the non-collided cars, when the number of collided cars grows, and, on the  
 230 other side it might provide too few collisions to learn from.

231 We simulated  $n_{\text{cars}}$  cars sharing the same driver policy parameters and moving in the simple road  
 232 graph in Figure 3 for  $T$  time steps. This road graph has 1 main intersection at the center, and 4 three-way  
 233 intersections. All the road segments  $p \in S$  have the same length  $l(p)$  and same number of lanes  $w(p)$ .  
 234 We used the model parameters shown in Table 1 and performed the simulations using Flow [36],  
 235 a microscopic discrete-time continuous-space road traffic simulator that allows implementing our  
 236 scenarios.

237 We repeated  $n_{\text{trials}}$  experiments in which we performed  $n_{\text{train}}$  training iterations in order to  
 238 optimize the initial random policy parameters  $\theta_{\text{no-rules}}$  and  $\theta_{\text{rules}}$ . We collected the values, across the  
 239  $n_{\text{trials}}$  repetitions, of  $R$ ,  $E$ , and  $C$  during the training.

240 We employed Proximal Policy Optimization (PPO) [37] as the RL policy optimization  
 241 algorithm: PPO is a state-of-the-art actor-critic algorithm that is highly effective, while being almost





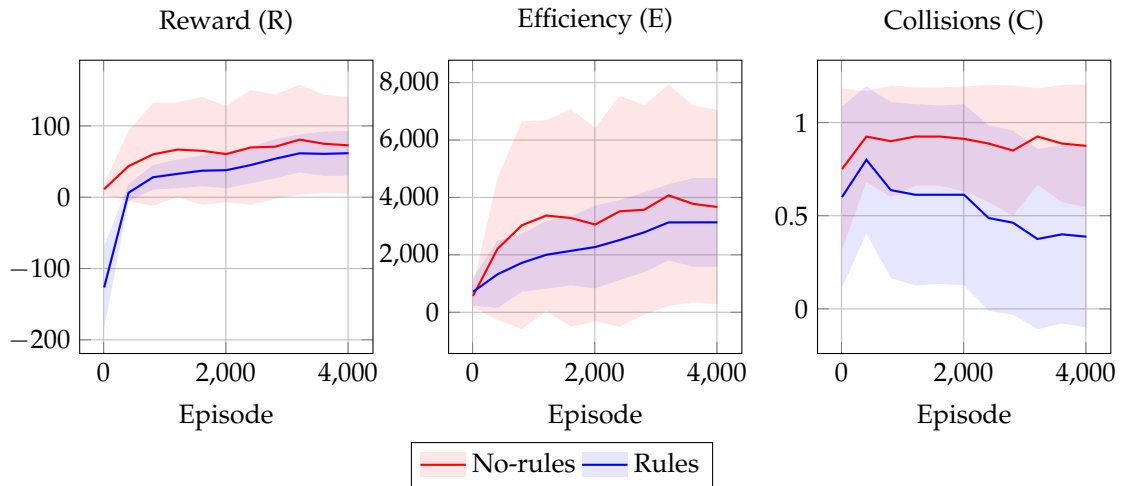
**Figure 3.** The road graph used in the experiments.

**Table 1.** Model and simulation parameters.

Param.	Meaning	Value
$l_{\text{car}}$	Car length	7
$t_{\text{coll}}$	Impact duration	10
$t_{\text{dead}}$	Collision duration	20
$d_{\text{view}}$	Driver's view distance	50
$v_{\text{max}}$	Driver's maximum speed	50
$a_{\text{max}}$	Driver's acceleration (deceleration)	2
$\Delta t$	Time step duration	0.2
$ S $	Number of road sections	12
$ I $	Number of road intersections	9
$w(p), p \in G$	Number of lanes	$\in \{1, 2\}$
$l(p), p \in S$	Section length	100
$n_{\text{car}}$	Cars in the simulation	40
$T$	Simulation time steps	500

**Table 2.** Policy learning algorithm parameters.

Param.	Meaning	Value
$n_{\text{trial}}$	Number of trials	20
$n_{\text{train}}$	Training iterations	500
$n_{\text{car}}$	Cars in the simulation	40
$\gamma$	Discount factor	0.999



**Figure 4.** Training results with cars removed after  $t_{\text{dead}}$  time steps. Here we draw the training values of  $R$ ,  $E$ , and  $C$ , at a certain training episode, averaged on  $n_{\text{trial}}$  experiments. We indicate with solid lines the mean of  $R, E$ , and  $C$  among the  $n_{\text{car}}$  vehicles, and with shaded areas their standard deviation among the  $n_{\text{car}}$  vehicles.

242 parameters-free. We used the PPO default configuration<sup>1</sup> with the parameters shown in Table 2. The  
 243 drivers policy is in the form of an actor-critic neural networks model, where each of the 2 neural  
 244 networks is made of 2 hidden layers, each one with 256 neurons and hyperbolic tangent as activation  
 245 function. The hidden layer parameters are shared between the actor and the critic networks: this is  
 246 a common practice introduced by [38] that helps to improve the overall performances of the model.  
 247 The parameters of the actor network as well as the ones of the critic network are initially distributed  
 248 according to the *Xavier initializer* [39].

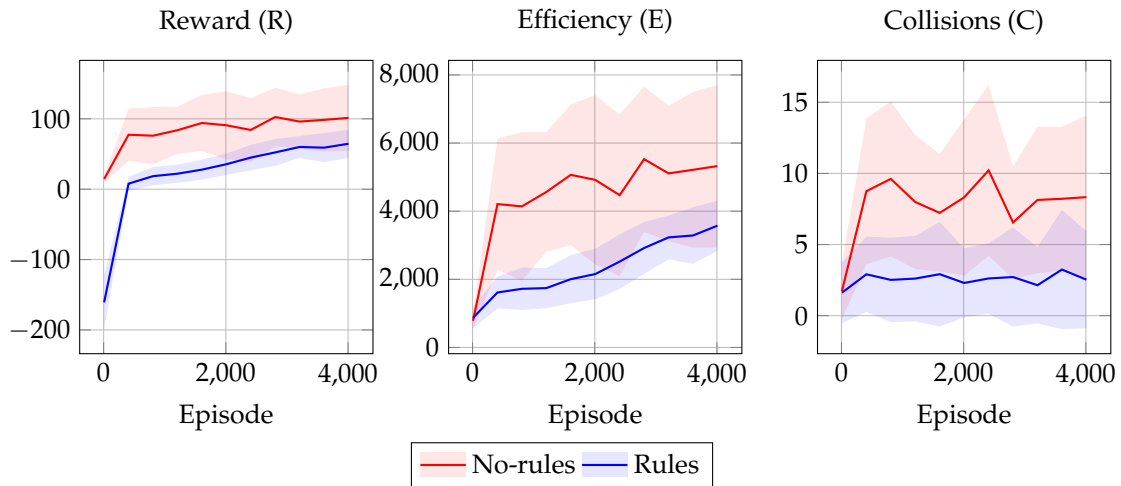
## 249 5. Results

250 Figures 4 and 5 show the training results in terms of the tuples  $R$ ,  $E$ , and  $C$  for the 2 policies  
 251  $\theta_{\text{no-rules}}$  and  $\theta_{\text{rules}}$  in the 2 collision scenarios considered.

252 In all the experimental scenarios the policy learned with rules shows driving behaviors that are  
 253 less efficient than the ones achieved by the one without rules. On the other hand, the policy learned  
 254 without rules is not even as efficient as it could theoretically be, due to the high number of collisions  
 255 that make it difficult to avoid collided cars. Moreover the values of  $E$  for the drivers employing the  
 256 rules are distributed closer to the mean efficiency value, and thus we can assume this is due to the fact  
 257 that the rules limit the space of possible behaviors to a smaller space w.r.t. the case without rules. In  
 258 other words, rules seems to favor equity among drivers.

259 On the other hand the policy learned with rules shows driving behaviors that are safer than the  
 260 ones achieved by the one without rules. This may be due to the fact that training every single driver to  
 261 avoid collisions based only on the efficiency reward is a difficult learning task, also because agents are

<sup>1</sup> <https://ray.readthedocs.io/en/latest/rllib-algorithms.html>



**Figure 5.** Training results with cars restored after  $t_{\text{dead}}$  time steps. Here we draw the training values of  $R$ ,  $E$ , and  $C$ , at a certain training episode, averaged on  $n_{\text{trial}}$  experiments. We indicate with solid lines the mean of  $R$ ,  $E$ , and  $C$  among the  $n_{\text{car}}$  vehicles, and with shaded areas their standard deviation among the  $n_{\text{car}}$  vehicles.

262 not capable of predicting the other agents' trajectories. On the other hand, we can see that the simple  
 263 traffic rules that we have designed are effective at improving the overall safety.

264 In other words these results show that, as expected, policies learned with rules are safer, but less  
 265 efficient than the ones without rules. Interestingly to us, the rules act also as a proxy for equality, as we  
 266 can see from Figures 4 and 5, in particular for the efficiency values  $E$ , where the blue shaded area is  
 267 way thinner than the red one, meaning that all the  $n_{\text{car}}$  vehicles have similar efficiency.

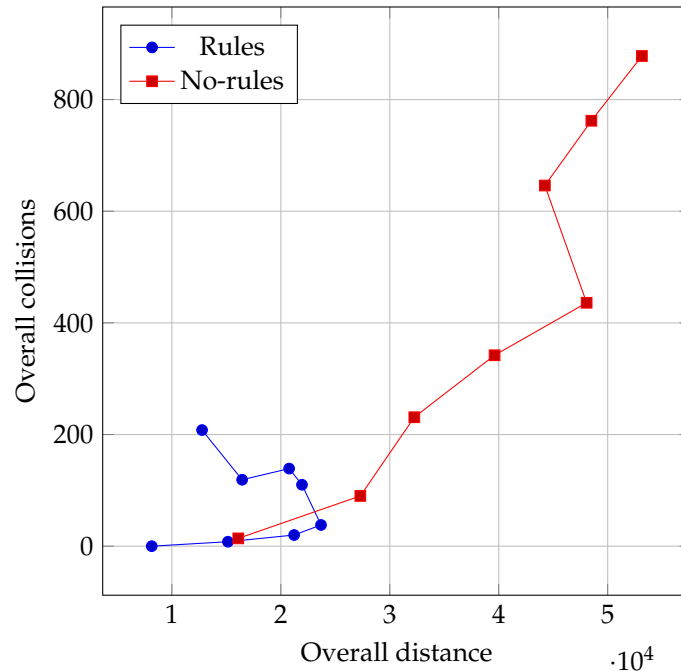
### 268 5.1. Robustness to traffic level

269 With the aim of investigating the impact of the traffic level on the behavior observed with the  
 270 learned policies (in the second learning scenario), we performed several other simulations by varying  
 271 the number of cars in the road graph. Upon each simulation, we measured the overall distance traveled  
 272  $\sum_{i=1}^{n_{\text{car}}} E_i \Delta t$  and overall collisions  $\sum_{i=1}^{n_{\text{car}}} C_i$ . We considered the overall sums, instead of the average, of  
 273 these indexes in order to investigate the impact of the variable number of cars in the graph: in principle,  
 274 the larger this number, the longer the overall distance that can be potentially traveled, and, likely, the  
 275 larger the number of collisions.

276 We show the results of this experiment in Figure 6, where each point corresponds to indexes  
 277 observed in a simulation with a given traffic level  $n_{\text{car}}$ : we considered values in 10, 20,  $\dots$ , 80. We  
 278 repeated the same procedure for both the drivers trained with and without the rules, using the same  
 279 road graph in which the drivers have been trained. For each level of traffic injected, we simulated  $T$   
 280 time steps and we measured the overall distance and overall number of collisions occurred.

281 As we can see from Figure 6, the two policies (corresponding to learning with and without rules)  
 282 exhibit very different outcomes as the injected traffic increases. In particular, the policy optimized  
 283 without rules results in an overall number of collisions that increases, apparently without any bound  
 284 in these experiment, as the traffic level increases. Conversely, the policy learned with the rules keeps  
 285 the overall number of collisions much lower also with heavy traffic. Interestingly, the limited increase  
 286 in collisions is obtained by the policy with the rules at the expense of overall traveled distance, i.e., of  
 287 traveling capacity of the traffic system.

288 From another point of view, Figure 6 shows that a traffic system where drivers learned to comply  
 289 with the rules is subjected to congestion: when the traffic level exceeds a given threshold, introducing  
 290 more cars in the system does not allow to obtain a longer traveled distance. Congestion is instead  
 291 not visible (at least not in the range of traffic level that we experimented with) with policies learned



**Figure 6.** Overall number of collisions in the simulation against the overall traveled distance in the simulation, averaged across simulations with the same  $n_{\text{car}}$ . Each dot is drawn from the sum of the values computed on the  $n_{\text{car}}$  vehicles.

292 without rules; the resulting system, however, is unsafe. All to all, congestion acts here as a mechanism,  
 293 induced by rules applied during the learning, for improving the safety of the traffic system.

## 294 6. Conclusions

295 We investigated the impact of imposing traffic rules while learning the policy for AI-powered  
 296 drivers in a simulated road traffic system. To this aim, we designed a road traffic model that allows to  
 297 analyze system-wide properties, as efficiency and safety, and, at the same time, permits learning using  
 298 a state-of-the-art RL algorithm.

299 We considered a set of rules inspired by real traffic rules and performed the learning with a  
 300 positive reward for traveled distance and a negative reward that punishes driving behaviors that are  
 301 not compliant with the rules. We performed a number of experiments and compared them with the  
 302 case where rules compliance does not impact on the reward function.

303 The experimental results show that imposing the rules during learning results in learned policies  
 304 that gives safer traffic. The increase in safety is obtained at the expense of efficiency, i.e., drivers travel,  
 305 on average, slower. Interestingly, the safer is improved also after the learning— i.e., when no reward  
 306 exists, either positive or negative—and despite the fact that, while training, rules are not enforced. The  
 307 flexible way in which rules are taken into account is relevant because it allows the drivers to learn  
 308 whether to evade a certain rule or not, depending on the current situation, and no action is prohibited  
 309 by design: rules stand hence as guidelines, rather than obligation, for the drivers. For instance, a driver  
 310 might have to overtake another vehicle in a situation in which overtaking is punished by the rules, if  
 311 this decision is the only one that allows to avoid a forthcoming collision.

312 Our work can be extended in many ways. One theme of investigation is the robustness of policies  
 313 learned with rules to the presence of other drivers, either AI-driven or human, who are not subjected  
 314 to rules or perform risky actions. It would be interesting to assess how the driving policies learned  
 315 with the approach presented in this study operate in such situations.

316 From a broader point of view, our findings may be useful in the situations where there is a trade-off  
 317 between compliance with the rules and a greater good. With the ever increasing pervasiveness of

AI-driven automation in many domains (e.g., robotics, content generation), relevance and quantity of these kind of situations will increase.

## References

1. Howard, D.; Dai, D. Public perceptions of self-driving cars: The case of Berkeley, California. Transportation research board 93rd annual meeting, 2014, Vol. 14, pp. 1–16.
2. Skrickij, V.; Sabanovic, E.; Zuraulis, V. Autonomous Road Vehicles: Recent Issues and Expectations. *IET Intelligent Transport Systems* **2020**.
3. Bojarski, M.; Del Testa, D.; Dworakowski, D.; Firner, B.; Flepp, B.; Goyal, P.; Jackel, L.D.; Monfort, M.; Muller, U.; Zhang, J.; others. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316* **2016**.
4. Maqueda, A.I.; Loquercio, A.; Gallego, G.; García, N.; Scaramuzza, D. Event-based vision meets deep learning on steering prediction for self-driving cars. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 5419–5427.
5. Sharifzadeh, S.; Chiotellis, I.; Triebel, R.; Cremers, D. Learning to drive using inverse reinforcement learning and deep q-networks. *arXiv preprint arXiv:1612.03653* **2016**.
6. Jaritz, M.; De Charette, R.; Toromanoff, M.; Perot, E.; Nashashibi, F. End-to-end race driving with deep reinforcement learning. 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 2070–2075.
7. Bouton, M.; Nakhaei, A.; Fujimura, K.; Kochenderfer, M.J. Safe reinforcement learning with scene decomposition for navigating complex urban environments. 2019 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2019, pp. 1469–1476.
8. Wang, C.; Liu, L.; Xu, C. Developing a New Spatial Unit for Macroscopic Safety Evaluation Based on Traffic Density Homogeneity. *Journal of Advanced Transportation* **2020**, 2020.
9. Qiao, Z.; Muelling, K.; Dolan, J.; Palanisamy, P.; Mudalige, P. Pomdp and hierarchical options mdp with continuous actions for autonomous driving at intersections. 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 2377–2382.
10. Tram, T.; Jansson, A.; Grönberg, R.; Ali, M.; Sjöberg, J. Learning negotiating behavior between cars in intersections using deep q-learning. 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 3169–3174.
11. Liebner, M.; Baumann, M.; Klanner, F.; Stiller, C. Driver intent inference at urban intersections using the intelligent driver model. 2012 IEEE Intelligent Vehicles Symposium. IEEE, 2012, pp. 1162–1167.
12. Isele, D.; Cosgun, A.; Subramanian, K.; Fujimura, K. Navigating intersections with autonomous vehicles using deep reinforcement learning. *arXiv preprint arXiv:1705.01196* **2017**.
13. Capasso, A.P.; Bacchiani, G.; Molinari, D. Intelligent Roundabout Insertion using Deep Reinforcement Learning. *arXiv preprint arXiv:2001.00786* **2020**.
14. Shalev-Shwartz, S.; Shammah, S.; Shashua, A. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295* **2016**.
15. Sutton, R.S.; Barto, A.G. *Reinforcement Learning: An Introduction*; A Bradford Book, 2018.
16. Mnih, V.; Kavukcuoglu, K.; Silver, D.; Graves, A.; Antonoglou, I.; Wierstra, D.; Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* **2013**.
17. Krizhevsky, A.; Sutskever, I.; Hinton, G.E. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 2012, pp. 1097–1105.
18. Sallab, A.E.; Abdou, M.; Perot, E.; Yogamani, S. Deep reinforcement learning framework for autonomous driving. *Electronic Imaging* **2017**, 2017, 70–76.
19. Loiacono, D.; Prete, A.; Lanzi, P.L.; Cardamone, L. Learning to overtake in TORCS using simple reinforcement learning. IEEE Congress on Evolutionary Computation. IEEE, 2010, pp. 1–8.
20. Hoel, C.J.; Wolff, K.; Laine, L. Automated speed and lane change decision making using deep reinforcement learning. 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 2148–2155.
21. Grigorescu, S.; Trasnea, B.; Cocias, T.; Macesanu, G. A survey of deep learning techniques for autonomous driving. *Journal of Field Robotics* **2019**.

- 368 22. Kiran, B.R.; Sobh, I.; Talpaert, V.; Mannion, P.; Sallab, A.A.A.; Yogamani, S.; Pérez, P. Deep Reinforcement  
369 Learning for Autonomous Driving: A Survey. *arXiv preprint arXiv:2002.00444* **2020**.
- 370 23. Brodsky, J.S. Autonomous vehicle regulation: How an uncertain legal landscape may hit the brakes on  
371 self-driving cars. *Berkeley Technology Law Journal* **2016**, *31*, 851–878.
- 372 24. Holstein, T.; Dodig-Crnkovic, G.; Pelliccione, P. Ethical and social aspects of self-driving cars. *arXiv preprint*  
373 *arXiv:1802.04103* **2018**.
- 374 25. Nyholm, S.; Smids, J. Automated cars meet human drivers: Responsible human-robot coordination and  
375 the ethics of mixed traffic. *Ethics and Information Technology* **2018**, pp. 1–10.
- 376 26. Kirkpatrick, K. The Moral Challenges of Driverless Cars. *Commun. ACM* **2015**, *58*, 19–20.  
377 doi:10.1145/2788477.
- 378 27. Rizaldi, A.; Althoff, M. Formalising traffic rules for accountability of autonomous vehicles. 2015 IEEE 18th  
379 International Conference on Intelligent Transportation Systems. IEEE, 2015, pp. 1658–1665.
- 380 28. Vanholme, B.; Gruyer, D.; Lusetti, B.; Glaser, S.; Mammari, S. Highly automated driving on highways based  
381 on legal safety. *IEEE Transactions on Intelligent Transportation Systems* **2013**, *14*, 333–347.
- 382 29. Medvet, E.; Bartoli, A.; Talamini, J. Road traffic rules synthesis using grammatical evolution. European  
383 Conference on the Applications of Evolutionary Computation. Springer, 2017, pp. 173–188.
- 384 30. O'Neill, M.; Ryan, C. Grammatical evolution. *IEEE Transactions on Evolutionary Computation* **2001**,  
385 *5*, 349–358.
- 386 31. Nenzi, L.; Bortolussi, L.; Ciancia, V.; Loreti, M.; Massink, M. Qualitative and quantitative monitoring of  
387 spatio-temporal properties. *Runtime Verification*. Springer, 2015, pp. 21–37.
- 388 32. Bartocci, E.; Bortolussi, L.; Loreti, M.; Nenzi, L. Monitoring mobile and spatially distributed cyber-physical  
389 systems. Proceedings of the 15th ACM-IEEE International Conference on Formal Methods and Models for  
390 System Design, 2017, pp. 146–155.
- 391 33. Tumova, J.; Hall, G.C.; Karaman, S.; Frazzoli, E.; Rus, D. Least-violating control strategy synthesis with  
392 safety rules. Proceedings of the 16th international conference on Hybrid systems: computation and control,  
393 2013, pp. 1–10.
- 394 34. Saunders, W.; Sastry, G.; Stuhlmüller, A.; Evans, O. Trial without error: Towards safe reinforcement  
395 learning via human intervention. Proceedings of the 17th International Conference on Autonomous Agents  
396 and MultiAgent Systems. International Foundation for Autonomous Agents and Multiagent Systems, 2018,  
397 pp. 2067–2069.
- 398 35. Mirchevska, B.; Pek, C.; Werling, M.; Althoff, M.; Boedecker, J. High-level decision making for safe and  
399 reasonable autonomous lane changing using reinforcement learning. 2018 21st International Conference  
400 on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 2156–2162.
- 401 36. Wu, C.; Kreidieh, A.; Parvate, K.; Vinitsky, E.; Bayen, A.M. Flow: A Modular Learning Framework for  
402 Autonomy in Traffic. *arXiv preprint arXiv:1710.05465* **2017**.
- 403 37. Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; Klimov, O. Proximal policy optimization algorithms.  
404 *arXiv preprint arXiv:1707.06347* **2017**.
- 405 38. Mnih, V.; Badia, A.P.; Mirza, M.; Graves, A.; Lillicrap, T.; Harley, T.; Silver, D.; Kavukcuoglu, K.  
406 Asynchronous methods for deep reinforcement learning. International conference on machine learning,  
407 2016, pp. 1928–1937.
- 408 39. Glorot, X.; Bengio, Y. Understanding the difficulty of training deep feedforward neural networks.  
409 Proceedings of the thirteenth international conference on artificial intelligence and statistics, 2010, pp.  
410 249–256.