

Modeling and Testing Multi-Agent Traffic Rules within Interactive Behavior Planning

Klemens Esterle¹, Luis Gressenbuch², and Alois Knoll²

Abstract—Autonomous vehicles need to abide by the same rules that humans follow. Some of these traffic rules may depend on multiple agents or time. Especially in situations with traffic participants that interact densely, the interactions with other agents need to be accounted for during planning. To study how multi-agent and time-dependent traffic rules shall be modeled, a framework is needed that restricts the behavior to rule-conformant actions during planning, and that can eventually evaluate the satisfaction of these rules. This work presents a method to model the conformance to traffic rules for interactive behavior planning and to test the ramifications of the traffic rule formulations on metrics such as collision, progress, or rule violations. The interactive behavior planning problem is formulated as a dynamic game and solved using Monte Carlo Tree Search, for which we contribute a new method to integrate history-dependent traffic rules into a decision tree. To study the effect of the rules, we treat it as a multi-objective problem and apply a relaxed lexicographical ordering to the vectorized rewards. We demonstrate our approach in a merging scenario. We evaluate the effect of modeling and combining traffic rules to the eventual compliance in simulation. We show that with our approach, interactive behavior planning while satisfying even complex traffic rules can be achieved. Moving forward, this gives us a generic framework to formalize traffic rules for autonomous vehicles.

I. INTRODUCTION

Traffic rules have been created to help humans manage the otherwise chaotic traffic environment. When sharing the road with human drivers, autonomous vehicles will have to obey the same rules humans do, which often depend on the actions of other agents or the past. Previous work has proposed a method for combining model checking techniques for traffic rule satisfaction with motion planning in static environments [1]. However, dense scenarios have shown to be difficult for such motion planning approaches, as they do not model the interactions with human drivers and are thus unable to correctly anticipate human reactions. The research line of *interactive behavior planning* addresses this problem but has mostly ignored the aspect to obey traffic rules other than collision prevention and speed compliance. The problem setting that is closest to ours is [2]. However, the approach is limited to rules that depend on only one agent and consequently cannot incorporate rules such as keeping a safe distance or merging in a zipper fashion.

Merging scenarios have proven to be challenging due to the dense interaction with others, combined with the

eventual lane ending, prompting the need for the driver to make a decision. Game-theoretic approaches offer an elegant way to model such interactions. We propose a game-based planning approach that monitors traffic rules, which depend on multiple agents and past information, at runtime. We refer to these rules by *multi-agent, time-dependent traffic rules*. To solve this formulation, we use Monte Carlo Tree Search (MCTS). Evaluating time-dependent traffic rules violates the Markov property, i.e., it does not rely only on the current state but also on previous states. This raises the question how to integrate those rules in MCTS, which depends on the Markov property to be computationally efficient.

Specifically, we contribute a game-based planning method monitoring multi-agent and time-dependent traffic rules, and a method to model non-Markovian traffic rules within MCTS. To study the effect of modeling a certain traffic rule within a set of rules, we treat the set of rules to be priority-ordered, and incorporate it to MCTS by leveraging methods from multi-objective optimization theory. In simulations based on real-world data, we use the same runtime monitors for both planning and evaluating a scenario run. We provide a comparable study of the ramifications on collision, progress, and the violation of rules in the simulated scenarios when modeling the applicable traffic rules within MCTS.

II. RELATED WORK

To safely drive in mixed traffic, multiple goals, such as collision and safety metrics, various traffic rules, and comfort metrics should be considered. Some of these goals are strictly preferred over others. Expressing the preference using a weighted sum scalarization is tedious, sometimes even impossible. Therefore, Reyes Castro *et al.* [1] define the preference relation in lexicographical order. Compared to a weighted sum scalarization, their method does not require to adapt the weights of all the cost terms.

Reyes Castro *et al.* [1] propose a sampling-based method to generate a trajectory that minimizes a set of formalized traffic rules in a finite fragment of Linear Temporal Logic (LTL). They only check collisions with static obstacles, as their use case is limited to a static environment while traveling at a constant speed. Their state space consists of position and orientation, which yields an analytic steering function, and thus enables the efficient connection of samples in the tree. However, there is no straight forward way to include dynamic obstacles from a prediction module.

When LTL is employed to formalize the traffic rules within motion planning [1, 3], violating such a formula only yields a binary satisfaction signal. With these planning approaches,

¹Klemens Esterle is with fortiss GmbH, Research Institute of the Free State of Bavaria, Munich, Germany, esterle@fortiss.org

²Luis Gressenbuch and Alois Knoll are with Robotics, Artificial Intelligence and Real-time Systems, Technische Universität München, Munich, Germany

neither the reactions of others are taken into account, nor does the violation penalty incorporate any prediction uncertainty, that circumvents this lack of interaction.

Satisfying multi-agent traffic rules has been realized by establishing contracts between vehicles [4]. However, their approach is used for deriving requirements during the design phase, not for penalizing or restricting actions during planning.

Chaudhari *et al.* [2] define a two-player non-zero-sum non-cooperative game. Both the ego agent as well as the environment (the other agents) try not to violate any traffic rules, which are expressed in LTL. Each agent builds up a tree as in [1]. That prohibits them from incorporating rules that depend on other agents than the ego agent, such as the safe distance or zipper merge rule.

Lee *et al.* [5] perform runtime verification on full traces of MCTS. They formalize simple single-agent traffic rules in LTL, and demonstrate it for an intersection. Instead, our approach can incorporate multi-agent and time-dependent traffic rules in an efficient way, by exploiting automata-based model checking in the search tree.

To summarize, no work currently exists that allows to study multi-agent traffic rules in dense scenarios by modeling it as part of the interactive planning problem.

III. PROBLEM STATEMENT

In a lane-merging scenario populated with multiple agents, this work aims to plan the behavior of a single agent, i.e. the generation of a sequence of desired future dynamic states $\{x(t = t_1), x(t = t_2), \dots, x(t = t_K)\}$ while obeying the applying traffic rules. We assume perfect observation of the dynamic state of the other agents, of the map and our localization within that.

We treat this setting of interactive behavior planning as a dynamic game, that formally consists of:

- A set of N agents, each having a dynamic state $x_i \in X_i$,
- An environment state $s \in S = \times X_i$ with state space S ,
- Agent 1 (referred to as “ego agent”) having an action space A_1 with discrete actions,
- Agents 2... N (referred to as “other agents”) following a behavior model B_o , that determines the agent’s action $a_i^t \sim B_o(s^t)$.

Based on the joint action of the agents, the environment described by joint state s transitions to the next state s' . The ego agent gets a reward $r(s, s', \mathbf{a})$ after joint action $\mathbf{a} \in \mathcal{A} = \times A_i$ is applied. The goal is to find an optimal action a for the ego vehicle that maximizes its cumulative reward $\sum_0^K \gamma^k r_k$ along a planning horizon of K steps, where the reward incorporates penalties for discomfort, traffic rule violation, and collision. The discount factor is denoted by γ . As an optimal solution of the decision problem is infeasible due to the size of the state space, we will use MCTS to obtain an approximation of the optimal solution using sampling.

We will use the formalized traffic rules from [6], which are formulated in LTL on finite traces (LTL_f). Some of these traffic rules, such as the zipper merge or overtaking, do not rely only on the current state s but also on previous

information. To be used in an interactive planning framework as described above, we aim to find a Markovian formulation for evaluating the history-dependent rules for each agent.

IV. PRELIMINARIES

A. Linear Temporal Logic on Finite Traces

Linear Temporal Logic is a discrete formal logic to reason not just about an absolute truth but about truths which might hold only at some points in time. It can thus be used to represent non-Markovian properties.

Let Π be a set of atomic propositions. The powerset of Π , i.e., the set of all subsets of Π , is denoted by 2^Π . A labeling function $\mathcal{L} : S \rightarrow 2^\Pi$ obtains the labels from state s .

Formally, the language φ of LTL formulas is defined as

$$\varphi ::= \pi \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \vee \varphi_2 \mid \varphi_1 \Rightarrow \varphi_2 \mid \bigcirc\varphi \mid \varphi_1 \cup \varphi_2 \mid \square\varphi \mid \diamond\varphi,$$

where $\pi \in \Pi$ denotes an *atomic* proposition, \neg (resp. \wedge , \vee , \Rightarrow) denote the boolean operators “not”, “and”, “or” and “implies”, and \bigcirc , (resp. \cup , \square , \diamond) denote the temporal operators “next”, “until”, “globally” (or “always”), “finally” (or “eventually”). See [7] for a definition of the semantics.

In the context of continuously replanning over a receding horizon, we use LTL_f, as it provides a formalism to reason over bounded periods of time. It uses the same syntax as LTL. We refer to the work of De Giacomo and Vardi [8] for definitions of the semantics of LTL_f.

Following the categorization of Manna and Pnueli [9], we will restrict ourselves to formulas with obligation properties, which include safety properties as well as guarantee properties. Safety formulas can be represented as $\square p$, whereas guarantees are generally captured by $\diamond p$, where p is a finite LTL formula only consisting of atomic propositions and the operators $\neg, \vee, \wedge, \Rightarrow, \bigcirc$, and past-LTL operators.

B. Automaton-based Verification of LTL_f Formulas

To verify if a trace satisfies the LTL_f formula, we utilize automata-based model checking. Given a formula in LTL_f, a deterministic finite automaton (DFA) can be constructed, that recognizes words satisfying the formula. For the LTL_f formula to be translated into a DFA, we follow the concept of [8]. An additional atomic proposition *alive* is introduced to the formula before translation to the automaton. This symbol will be set to true as long as the search horizon has been reached, otherwise false. The automaton is defined over the alphabet $\Sigma = \Pi$, i.e., the set of atomic propositions of the LTL_f formula.

Definition 4.1: A deterministic finite automaton is a tuple $\Lambda = (Q, q_0, \Sigma, \delta, F)$, where

- Q is a set of states
- $q_0 \in Q$ is the initial state
- Σ is a finite alphabet
- $\delta : Q \times \Sigma \rightarrow Q$ is a transition function
- $F \subseteq Q$ denotes a set of accepting (or final) states.

For a word w , given as a sequence of symbols $\sigma \in \Sigma$ and a DFA Λ , the automaton state is sequentially updated from the

initial automaton state with respect to δ and w . Therefore, if the automaton halts in an accepting state, the trajectory satisfies the formula.

C. Monte Carlo Tree Search for Behavior Planning

MCTS is an online method to approximate the solution of sequential decision problems via sampling. Throughout the search, the estimate of the state-action-value function $Q(s, a)$, which maps the expected cumulative reward of performing action a in state s , is updated iteratively. Each iteration consists of *selection*, *expansion*, *rollout*, and *backup*.

For *selecting* a new node, the tree is traversed following the tree policy of best actions until a state with untried actions is reached. A commonly used selection strategy is called Upper Confidence Bounds for Trees (UCT), which balances exploitation and exploration. During *expansion*, an untried actions is applied, and the child node is added to the tree. From the newly expanded node, a value estimate is obtained using a heuristic *rollout* (also called *simulation*): Based on a default policy, actions are applied until a terminal state is reached. The simplest case is to make uniform random moves. The observed value is *backed up* to the root node and used to update $Q(s, a)$. The search is repeated until some predefined computation budget (time or iterations) is reached. Finally, the best performing root action is returned.

Lenz *et al.* [10] applied this concept to cooperative driving, as they introduce cooperative costs that include the costs of all agents. When sampling the actions of all agents, the size of the search tree grows exponentially with the number of agents. To mitigate that, only a subset of all agents are included in the joint action, while the actions of the remaining agents are based on a predefined model. We will base our work on a special case of this, called *Single-Agent MCTS*, where all other agents are predicted using this model. However, the approach can be easily extended to multiple agents.

V. APPROACH

In this section, we first present our method for monitoring traffic rules within interactive behavior planning based on MCTS. We describe how we obtain a reward from the rule violation. We then describe how to extend MCTS to rewards in a lexicographical ordering.

A. Runtime Monitoring within Monte Carlo Tree Search

To allow for a valuation of the traffic rules given in LTL_f , we label the environment state s according to whether an atomic proposition such as “agent i is in lane” or “agent i is in front of j ” is *true* or *false*. A description of the necessary labels to express the traffic rules is available in [6].

To monitor temporal rules within MCTS, the straightforward way is to evaluate the run from the root node to the current node for each expansion and simulation step. However, this will significantly limit the performance of the MCTS. Instead, we exploit the structure of the decision tree, by augmenting the tree nodes of the MCTS by the automaton state. Thus, non-Markovian properties captured in

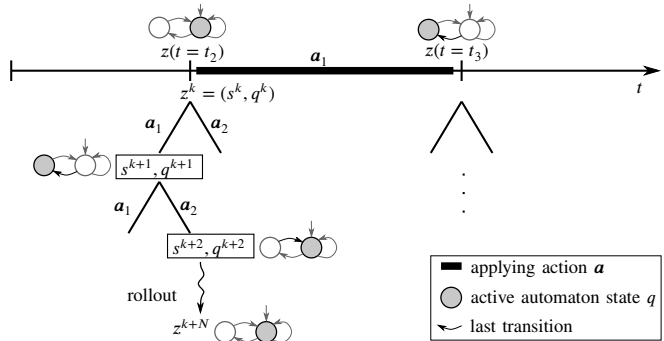


Fig. 1: The evolution of the product state z for a single rule monitor. An MCTS planning step is depicted exemplarily at t_2 .

LTL can be efficiently verified during the MCTS by encoding necessary historic information into the automaton’s state.

We first translate each φ_i formula into their corresponding DFA representation Λ_i . Instead of combining them to a single product automaton, we consider them as independent automata, which is advantageous for the time complexity of finding a valid transition. The states of m automata form the automaton state vector \bar{q}

$$\bar{q} = (q_1 \quad \dots \quad q_m)^T. \quad (1)$$

With the joint state s and the automaton state vector \bar{q} , we can define the combined state z

$$z^k = (s^k \quad \bar{q}^k)^T. \quad (2)$$

The successor state is then defined by

$$z^{k+1} = \begin{pmatrix} s^{k+1} \\ \delta_1(q_1^k, \mathcal{L}(s^{k+1})) \\ \vdots \\ \delta_m(q_m^k, \mathcal{L}(s^{k+1})) \end{pmatrix} \quad (3)$$

where s^{k+1} is defined by the joint action \mathbf{a} . Fig. 1 illustrates the evolution of the product state z during MCTS and over time for a single rule monitor.

B. Costs for Rule Violations

Previous work [1, 11] introduced a weighted transition for rule violation to the automaton. During planning, each rule automaton is evaluated over the word w , and the weight equals the penalty for the respective rule. Specifically, a self-looping automaton is used in [1], i.e., a weighted transition to the same state is added for a violating transition. However, this formulation does not allow to model violations that cause a penalty once vs. violations that accumulate penalties over time. To account for this, Schluter *et al.* [11] extended the work of [1], defining an explicit violation symbol.

Although violations of safety properties can be detected on partial traces, a safety property $\Box p$ will never be satisfied again after being violated. Consequently, violating a safety property leads to a non-accepting state, that only has a self-loop. The automaton thus cannot be brought back to an accepting state anymore, once it has registered a violation.

To obtain penalties for consecutive rule violations, such as multiple times overtaking on the right [3], we reset the automaton to its initial state, if a violation occurred. Formally, we modify $\Lambda = (Q, q_0, \Sigma, \delta, F)$ to $\bar{\Lambda} = (Q, q_0, \Sigma, \bar{\delta}, F)$, where

$$\bar{\delta}(q, \sigma) := \begin{cases} q_0 & \text{if } \forall \sigma' \in \Sigma : \delta(\delta(q, \sigma), \sigma') = \delta(q, \sigma) \wedge \\ & \delta(q, \sigma) \notin F \\ q' \in F & \text{if } q = q_0 \wedge \text{alive} \notin \sigma \\ \delta(q, \sigma) & \text{else} \end{cases} \quad (4)$$

Extending our formulation to be invariant to the step size will be the subject of future work.

Let a given rule formalized in LTL φ be represented by $\bar{\Lambda}_\varphi$. The weighting function $W : Q \times \Sigma \rightarrow \mathbb{R}$ defines the penalty for violating that rule, assigning each transition a scalar weight:

$$W(q, \sigma) := \begin{cases} \omega & \text{if } \forall \sigma' \in \Sigma : \delta(\delta(q, \sigma), \sigma') = \delta(q, \sigma) \wedge \\ & \delta(q, \sigma) \notin F // \text{safety} \\ \omega & \text{if } \text{alive} \notin \Sigma \wedge q \notin F // \text{guarantee} \\ 0 & \text{else} \end{cases} \quad (5)$$

A penalty ω is returned, if the modified automaton $\bar{\Lambda}_\varphi$ gets reset to its initial state (violation of a safety property $\square p$) or the automaton does not halt in an accepting state (violation of a guarantee property $\diamond p$). Finally, the reward for violating φ is then defined as

$$r(s^k, s^{k+1}) = W(q, \mathcal{L}(s^{k+1})) \quad (6)$$

C. Multi-Objective Reward Function with Priorities

To model the multi-objective reward (i.e., safety, legal, comfort) and to prevent weight-tuning for a scalar reward, we treat the reward elements as a vector:

$$\mathbf{r} = (r_1 \quad \dots \quad r_n)^\top \quad (7)$$

If goals are meant to be traded off among each other, we perform weighted sum scalarization for these. While Wang and Sebag [12] try to find elements on the Pareto-optimal front, the multi-objective formulation for us simplifies, as our objective vector is ordered according to the priority levels. An entry at index i in the reward vector denotes a higher priority than at index $i + 1$.

1) *UCT for vectorized rewards*: To incorporate a vectorized reward to the MCTS, we modify the MCTS selection strategy by computing the UCT value per reward vector element. We then compare the vector elements according to the lexicographical order.

2) *Relaxed lexicographical order*: If employing strict lexicographical order, the selection would favor tree branches that are suboptimal in terms of all other criteria over tree branches, where one single outcome was a collision (if that is the top-level priority). To account for the inaccuracy of the sampling-based approximation of the optimal solution, we implement a relaxation of the lexicographical order,

called *Thresholded Lexicographical Ordering* (TLO) [13]. Comparing two reward vectors \mathbf{r}, \mathbf{r}' according to TLO, is defined as

$$\mathbf{r} \preceq_{\text{TLO}} \mathbf{r}' \iff \exists i : r_i \leq r'_i \wedge \forall j < i : (r_j > \tau_j \wedge r'_j > \tau_j) \vee r_j = r'_j, \quad (8)$$

where \preceq_{TLO} denotes the thresholded lexicographic comparison. TLO uses a threshold vector $\boldsymbol{\tau}$ to determine if two goals are sufficiently close so that the next lower priority level can be considered.

VI. EVALUATION

A. Experimental Setup

We study our approach using the open-source benchmarking and development framework BARK proposed in [14]. We use Spot [15], a library for model checking, to translate the formalized LTL formula to a DFA, and to manipulate the automata. We implement the runtime monitoring and the multi-objective reward function on top of a template-based MCTS library [16].

1) *Benchmarking Framework*: BARK is a multi-agent environment tailored to develop interactive behavior models. It allows to easily exchange behavior models for planning, prediction and simulation. BARK already offers a range of behavior models, which we can use for predicting and simulating the other agents. It provides efficient collision checking and realistic maps. A kinematic bicycle model is employed to simulate the vehicles.

2) *Variants*: We will study

- **SA** as a scalar single agent variant of MCTS with penalties for collision and comfort,
- **SA-Lex** as a lexicographic baseline implementing the same rewards as SA, but as a vectorized reward,
- **SA-Lex (Zip)** extending SA-Lex by including a penalty for the zipper merge,
- **SA-Lex (SD)** extending SA-Lex by including a penalty for the safe distance (SD) rule.
- **SA-Lex (Zip > SD)** and **SA-Lex (SD > Zip)** extending SA-Lex by including a penalty for for the zipper merge and the safe distance rule.

We define a base reward

$$r_{\text{base}} = r_a + r_{\text{lat}} + r_{\Delta v} + r_\phi, \quad (9)$$

which consists of penalties for

- Longitudinal acceleration $r_a = -w_a a^2 \Delta t$
- Lateral acceleration $r_{\text{lat}} = -w_{\text{lat}} |\dot{\theta}| v^2 \Delta t$
- Difference to desired velocity $r_{\Delta v} = -w_v |v - v_r| \Delta t$

with respective weights w_\square , velocity v , reference velocity v_r , acceleration a , orientation rate $\dot{\theta}$ and time increment Δt . The convergence of MCTS can be accelerated by incorporating additional domain knowledge. We define a potential function $\phi(s) = -w_\phi |v - v_r| \Delta t$ and the potential-based shaping function as

$$r_\phi = \gamma \phi(s_{k+1}) - \phi(s_k) \quad (10)$$

r_{col} denotes the penalty for not colliding. r_{sd} and r_{zip} denote the penalty for violating the safe distance and the

TABLE I: Variants of \mathcal{B}_e showing the reward vector \mathbf{r} and the size of the reward vector.

\mathcal{B}_e	Reward vector \mathbf{r}	size(\mathbf{r})
SA	$\mathbf{r} = (r_{\text{col}} + r_{\text{base}})^T$	1
SA-Lex	$\mathbf{r} = (r_{\text{col}} \quad r_{\text{base}})^T$	2
SA-Lex (Zip)	$\mathbf{r} = (r_{\text{col}} \quad r_{\text{zip}} \quad r_{\text{base}})^T$	3
SA-Lex (SD)	$\mathbf{r} = (r_{\text{col}} \quad r_{\text{sd}} \quad r_{\text{base}})^T$	3
SA-Lex (Zip > SD)	$\mathbf{r} = (r_{\text{col}} \quad r_{\text{zip}} \quad r_{\text{sd}} \quad r_{\text{base}})^T$	4
SA-Lex (SD > Zip)	$\mathbf{r} = (r_{\text{col}} \quad r_{\text{sd}} \quad r_{\text{zip}} \quad r_{\text{base}})^T$	4

zipper merge rules, respectively. The “zipper merge” rule requires vehicles, that are on a continuing lane, to let vehicles on an ending lane merge in a zipper fashion. The “safe distance” rule requires to leave a safe distance to the vehicle in front. Both rules are defined in [6]. Table I shows the reward vectors of these variants.

3) *Action Space*: We model the discrete action space for the ego agent as “lane keeping at constant acceleration” for $a \in \{0 \frac{\text{m}}{\text{s}^2}, 1 \frac{\text{m}}{\text{s}^2}, -2 \frac{\text{m}}{\text{s}^2}, -8 \frac{\text{m}}{\text{s}^2}\}$, “lane changing at constant velocity”, and “gap keeping”³ based on the Intelligent Driver Model (IDM) [17]. The parameters for the “gap keeping” primitive are shown in Table II with the exception of the desired speed, which we set to $14 \frac{\text{m}}{\text{s}}$.

4) *Scenarios*: We want to study our approach using real-world data from the INTERACTION dataset [18]. However, we cannot just replace the agent with our model and replay the other agents, as the other agents would not react to our model under test anymore. We thus preserve the initial configuration of the vehicles from the dataset but simulate them according to \mathcal{B}_o . Each vehicle is simulated as the ego agent in one scenario. The ego agent is controlled by the behavior model \mathcal{B}_e . The scenario is passed successfully, if the ego agent reaches its goal region, which we create from the last pose of the agent in the dataset.

5) *Behavior Model for Others*: The actions of the other vehicles are calculated using the behavior model \mathcal{B}_o , for which we employ a rule-based model BehaviorMobilRule-Based, with IDM as a longitudinal and MOBIL [19] as a lateral model, and a lane filtering mechanism on top of MOBIL. The model is available as open-source in BARK. Table II shows the parameters. To cause challenging situations for the ego vehicle, we let the other vehicles travel at a desired speed of $10 \frac{\text{m}}{\text{s}}$, and the ego vehicle at $14 \frac{\text{m}}{\text{s}}$.

6) *Rule Evaluators*: BARK provides an abstract evaluator class, that calculates a given metric such as collision or step count based on the simulated world state. We have extended this evaluator concept to evaluate arbitrary LTL formulas on finite traces, and made it available within BARK. Each rule is captured in a rule monitor, which we can use to monitor compliance throughout the simulation. Both the behavior model and the evaluator employ the same rule monitor, which we have contributed as open-source¹. This allows us to study whether \mathcal{B}_e truly satisfies the modeled rules. Fig. 2 shows the resulting evaluation framework.

¹<https://github.com/bark-simulator/rule-monitoring>

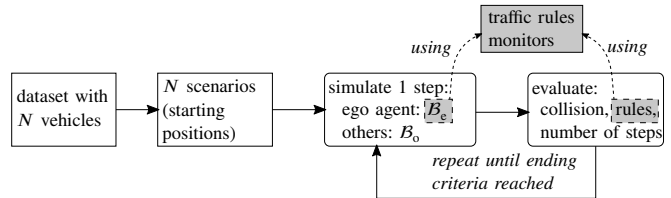


Fig. 2: Framework to evaluate the behavior of the ego agent \mathcal{B}_e in a closed-loop simulation. The initial starting positions for the vehicles are taken from the dataset. The other vehicles are simulated using a behavior model \mathcal{B}_o . The traffic rule monitors are used within \mathcal{B}_e and to evaluate the simulation.

B. Quantitative Evaluation

We evaluate \mathcal{B}_e for 200, 500 and 1000 search iterations. Fig. 3 shows the share of controlled agents to collide, successfully reach the goal region (within 30 s) and to violate the zipper merge or safe distance rule.

We observe nearly no collisions, which indicates that MCTS with high level actions can approximate the solution well, if predicted and simulated behavior are matching. *SA*, *SA-Lex* and *SA-Lex (SD)* do violate the zipper merge at about 25%, which shows the this rule will not be satisfied implicitly and motivates its explicit modeling. *SA-Lex (SD)* violates it to some less extent, as keeping a safe distance all the time sometimes leaves enough space for the merging vehicle to fit in. The variants not implementing the safe distance rule usually do not leave much space to the front vehicle, as the correct prediction model yields an accurate estimate of what will happen given a certain action, as long as the tree is explored enough. However, in order to obey to the safe distance rule, and to be prone to unexpected actions of the human drivers, it must be modeled within the planner. The remaining safety distance violations stem from unsafe initial states at the beginning of the scenario. The variants *SA-Lex (Zip > SD)* and *SA-Lex (SD > Zip)* do not cause any violations to the zipping rule, but also keep a safe distance for most scenarios. The results for both are similar, which indicates that there exists no trade-off in our scenarios between the safe distance and the zipper rule.

When varying the number of search iterations, the results remain mostly unchanged for the variants *SA*, *SA-Lex* and *SA-Lex (SD)*. With more search iterations, \mathcal{B}_e becomes more certain about which actions will prevent a zipper merge violation. As \mathcal{B}_o is not necessarily collision-free (it does not employ any prediction), this creates more dense situations, and thus explains the slightly increasing number of collisions.

VII. CONCLUSION AND FUTURE WORK

In this work, we investigated the problem of how interactive behavior planning for an autonomous vehicle can be modeled to obey the traffic rules, and how this can be evaluated and tested. We modeled the interaction as a dynamic game and used MCTS to solve the problem, while incorporating ideas from model checking and multi-objective optimization.

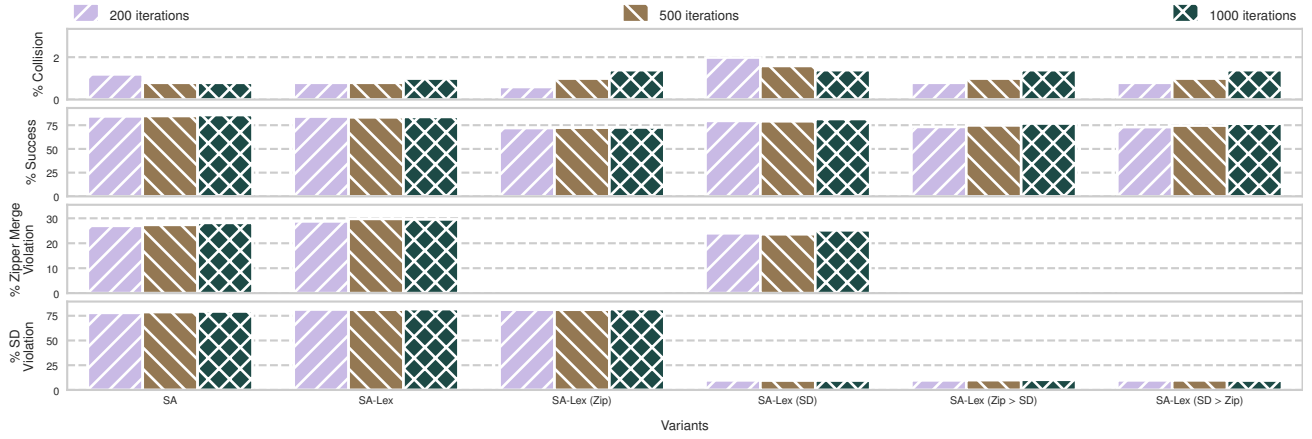


Fig. 3: Benchmark of the variants for \mathcal{B}_c . The comparison of *SA-Lex (Zip > SD)* and *SA-Lex (SD > Zip)* with reversed priorities shows that both modeled rules can be satisfied at the same time and thus no tradeoff between those rules is required for our scenario.

TABLE II: Parameters of BehaviorMobilRuleBased for \mathcal{B}_o . The lane change rules serve as an additional filter of the lanes, to which MOBIL can choose to change to.

	Parameter	Unit	Value
IDM	Desired velocity	[m/s]	10
	Maximum acceleration	[m/s ²]	1.7
	Desired time headway	[s]	2.5
	Comfortable deceleration	[m/s ²]	2
	Minimum distance	[m]	2
Lane change rules	Min. rear distance	[m]	0.5
	Min. front distance	[m]	1
	Time Gap	[s]	0.5
MOBIL	Politeness Factor	[1]	0
	Safe deceleration	[m/s ²]	4
	Acceleration threshold	[m/s ²]	0.2

In our evaluation, we demonstrated the capabilities of our approach in a closed-loop simulation for a merging scenario at dense traffic in a systematic way. For this, we selected two rules that apply to this situation, namely to keep a safe distance and to merge in a zipper fashion.

Our new method allows us to formalize multi-agent rules, that apply to more than one agent, whose future motion is modeled interactively. With this, multi-agent time-dependent traffic rules can be studied on whether they have been correctly formalized, and on what the ramifications of the rules are on safety and progress. Future work should investigate the probabilistic interpretation of predicates.

ACKNOWLEDGMENT

This research was funded by the Bavarian Ministry of Economic Affairs, Regional Development and Energy.

REFERENCES

- [1] L. I. Reyes Castro, P. Chaudhari, J. Tumova, *et al.*, “Incremental sampling-based algorithm for minimum-violation motion planning,” in *52nd IEEE Conference on Decision and Control*, 2013.
- [2] P. Chaudhari, T. Wongpiromsarny, and E. Frazzoli, “Incremental minimum-violation control synthesis for robots interacting with external agents,” in *2014 American Control Conference*, 2014.
- [3] K. Esterle, V. Aravantinos, and A. Knoll, “From specifications to behavior: Maneuver verification in a semantic state space,” in *IEEE Intelligent Vehicles Symposium, Proceedings*, 2019.
- [4] J. Decastro, L. Liebenwein, C.-i. Vasile, *et al.*, “Counterexample-Guided Safety Contracts for Autonomous Driving,” in *Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2018.
- [5] J. Lee, A. Balakrishnan, A. Gaurav, *et al.*, “WiseMove: A Framework for Safe Deep Reinforcement Learning for Autonomous Driving,” 2019. arXiv: 1902.04118.
- [6] K. Esterle, L. Gressenbuch, and A. Knoll, “Formalizing Traffic Rules for Machine Interpretability,” in *2020 IEEE 3rd Connected and Automated Vehicles Symposium (CAVS)*, 2020.
- [7] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT Press, 2008.
- [8] G. De Giacomo and M. Y. Vardi, “Linear Temporal Logic and Linear Dynamic Logic on Finite Traces,” in *Twenty-Third International Joint Conference on Artificial Intelligence*, 2013.
- [9] Z. Manna and A. Pnueli, “A hierarchy of temporal properties,” in *Proceedings of the ninth annual ACM symposium on Principles of distributed computing*, 1990.
- [10] D. Lenz, T. Kessler, and A. Knoll, “Tactical Cooperative Planning for Autonomous Vehicles using MCTS,” in *IEEE Intelligent Vehicles Symposium (IV)*, IEEE, Jun. 2016, pp. 447–453.
- [11] H. Schluter, P. Schillinger, and M. Burger, “On the Design of Penalty Structures for Minimum-Violation LTL Motion Planning,” in *2018 IEEE Conference on Decision and Control (CDC)*, 2018.
- [12] W. Wang and M. Sebag, “Multi-objective monte-carlo tree search,” *Journal of Machine Learning Research*, vol. 25, 2012.
- [13] P. Vamplew, R. Dazeley, A. Berry, *et al.*, “Empirical evaluation methods for multiobjective reinforcement learning algorithms,” *Machine Learning*, vol. 84, no. 1-2, 2011.
- [14] J. Bernhard, K. Esterle, P. Hart, *et al.*, “BARK: Open Behavior Benchmarking in Multi-Agent Environments,” in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [15] A. Duret-Lutz, A. Lewkowicz, A. Fauchille, *et al.*, “Spot 2.0 — a framework for LTL and ω -automata manipulation,” *International Symposium on Automated Technology for Verification and Analysis*, 2016.
- [16] J. Bernhard, *MA-MCTS: A configurable library for Multi-Agent Monte Carlo Tree Search in C++*, 2019.
- [17] M. Treiber, A. Hennecke, and D. Helbing, “Congested traffic states in empirical observations and microscopic simulations,” *Physical Review E - Statistical Physics, Plasmas, Fluids, and Related Interdisciplinary Topics*, vol. 62, no. 2, pp. 1805–1824, 2000.
- [18] W. Zhan, L. Sun, D. Wang, *et al.*, “INTERACTION Dataset: An INTERNATIONAL, Adversarial and Cooperative moTION Dataset in Interactive Driving Scenarios with Semantic Maps,” 2019. arXiv: 1910.03088.
- [19] A. Kesting, M. Treiber, and D. Helbing, “General lane-changing model MOBIL for car-following models,” *Journal of the Transportation Research Board*, 2007.