# Cooperative Behavior Planning for Automated Driving Using Graph Neural Networks

Marvin Klimke[1,2], Benjamin Völz[1], and Michael Buchholz[2]

*Abstract*— Urban intersections are prone to delays and inefficiencies due to static precedence rules and occlusions limiting the view on prioritized traffic. Existing approaches to improve traffic flow, widely known as automatic intersection management systems, are mostly based on non-learning reservation schemes or optimization algorithms. Machine learning-based techniques show promising results in planning for a single ego vehicle. This work proposes to leverage machine learning algorithms to optimize traffic flow at urban intersections by jointly planning for multiple vehicles. Learning-based behavior planning poses several challenges, demanding for a suited input and output representation as well as large amounts of ground-truth data. We address the former issue by using a flexible graph-based input representation accompanied by a graph neural network. This allows to efficiently encode the scene and inherently provide individual outputs for all involved vehicles. To learn a sensible policy, without relying on the imitation of expert demonstrations, the cooperative planning task is considered as a reinforcement learning problem. We train and evaluate the proposed method in an open-source simulation environment for decision making in automated driving. Compared to a first-in-first-out scheme and traffic governed by static priority rules, the learned planner shows a significant gain in flow rate, while reducing the number of induced stops. In addition to synthetic simulations, the approach is also evaluated based on real-world traffic data taken from the publicly available inD dataset.

## I. INTRODUCTION

Urban traffic regularly exhibits disturbances and inefficiencies caused by simple traffic management schemes faced with large volumes of vehicles. Especially smaller intersections are typically handled by static priority rules, resulting in vehicles approaching from a minor road having to yield. Moreover, occlusions through buildings or other objects are highly prevalent in urban areas, limiting the view for both human drivers and vehicle-bound sensory systems.

The increasing use of connected vehicles (CVs) and connected automated vehicles (CAVs) opens up new opportunities to increase the traffic efficiency. Those vehicles can announce their presence and possibly share perception data with surrounding road users via a communication link. Moreover, with edge computing resources becoming available in
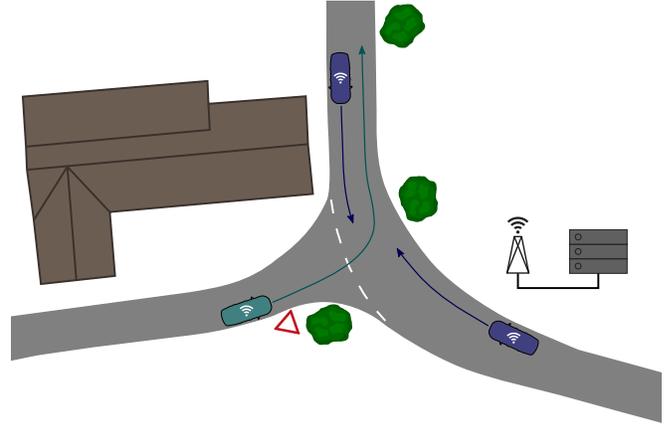


Fig. 1. Cooperative maneuver at an urban intersection. The planning module on the edge server requests the blue vehicles on the main road to slow down. Hence, the turquoise vehicle can merge without having to stop.

urban areas, it is viable to build and maintain a local environment model of, e.g., an intersection and its surroundings. Such an edge server can distribute the environment model to connected vehicles in the operational area. CAVs, for instance, can make use of the information by incorporating it into their planning algorithms.

In the publicly funded project MEC-View, research on connected automated driving was conducted using a testing site at a suburban three-way intersection in the city of Ulm in Germany [1]. Due to buildings occluding the view onto the priority road, an automated vehicle merging from the side road has to decelerate strongly, before being able to safely enter the intersection based solely on its own perception system. With the support of the environment model provided by the edge server, the automated vehicle can transition smoothly onto the main road, given that appropriate space is available. We build upon this approach and discuss the potential of multi-agent planning schemes that are executed on the server. Based on the fused environment model, a cooperative plan for handling intersection traffic is derived that can be proposed to the connected vehicles as behavioral instructions. Explicit deviations from static priority rules become possible. For instance, vehicles on the main road can be requested to slow down and thus letting a vehicle from the side road merge into the emerging gap, as depicted in Fig. 1.

Prior research on automatic intersection management (AIM) primarily focuses on non-learning algorithms, like

[1]The authors are with the Robert Bosch GmbH, Corporate Research, D-71272 Renningen, Germany. E-Mail: {marvin.klimke, benjamin.voelz}@de.bosch.com
[2]The authors are with the Institute of Measurement, Control and Microtechnology, Ulm University, D-89081 Ulm, Germany. E-Mail: michael.buchholz@uni-ulm.de

reservation-based or optimization-based programs. At the same time, machine learning-based approaches show remarkable results on prediction tasks in automated driving as well as planning for a single ego vehicle. The lack of fitting ground-truth data prevents the application of supervised learning for cooperative behavior planning. To bridge this gap, the present work proposes to train a reinforcement learning (RL) policy for multi-agent planning in a simulated environment, resulting in the following contributions:

- Leveraging machine learning to perform cooperative multi-agent planning for urban automated driving,
- To the best of our knowledge, we propose the first AIM system exploiting graph neural networks,
- Evaluation based on real-world traffic data taken from a publicly available urban driving dataset.

The remainder of the paper is structured as follows: Section II discusses related work in the field of AIM and machine learning-based planning for automated driving. The proposed behavioral planning scheme is comprehensively introduced in Section III. Afterwards, evaluation results obtained in synthetic simulations and based on real-world traffic data are given in Section IV. Section V concludes the paper and gives an outlook on future work.

## II. RELATED WORK

The analysis on the state of the art first considers existing approaches to AIM. Because machine learning is seldom used for AIM, we subsequently investigate learning-based behavioral planning methods. Due to the large body of existing literature, we present a selection of commonly used approaches and refer the reader to surveys for a more extensive overview.

### A. Automatic Intersection Management

Past research brought forth a variety of AIM schemes, surveyed for instance by [2]. The authors identify centralization as a crucial feature for distinguishing AIM schemes. Thereby, a fully centralized scheme exhibits a single coordination unit that is in charge of planning the intersection traversal and acts as the communication partner for all vehicles. In a fully distributed AIM, a cooperative plan is negotiated by the vehicles on their own.

A centralized, reservation-based AIM system is proposed in [3], which employs a first-in-first-out policy for assigning clearance to cross the intersection. A driver agent places a request when its vehicle is about to enter the monitored intersection area, covering possible conflict points with other paths. The intersection manager maintains tile-based reservations and confirms the request if the affected tiles are free. This approach can be combined with a traffic light to enable the co-usage of the intersection by human drivers and automated vehicles.

Optimization-based intersection management systems are published, for instance, in [4] and [5]. Those works assume full penetration of CAVs that laterally follow predefined lanes on urban intersections. The distributed energy-optimizing approach [4] further disallows turning maneuvers and the utilization of two conflicting paths at the same time. In [5], the longitudinal control of vehicles is performed by a centralized intersection coordination unit employing a model predictive control (MPC) scheme. Both works demonstrate efficiency gains in time and fuel consumption by comparison to a traditional signalized intersection. A prevalent issue with optimization-based approaches is the unfavorable scaling of computational demand with increasing traffic density. In [6], a novel AIM scheme is presented, that is also capable of handling mixed traffic, i.e., simultaneous usage by automated and human-driven vehicles.

The concept of platooning [7] can also be used for AIM. Based on a so-called virtual inter-vehicle distance, a pair of vehicles can adapt their velocities to cross their conflict point with sufficient clearance. The authors acknowledge that significant adaptions would have to be made for managing an intersection under mixed traffic.

### B. Machine Learning-Based Planning

Machine learning-based approaches to automated driving experience rising interest of researchers during the last years. A survey of recent deep reinforcement and imitation learning planning methods for a single ego vehicle can be found in [8]. The authors categorize published works by the type of input data (e.g. sensor measurements or object detections) and output representation (e.g. behavioral planning or direct control outputs). Because individual sensor measurements are not suited for cooperative planning over multiple vehicles, we limit our analysis to methods that require a prior perception system to be in place. Readers interested in machine learning-based prediction for automated driving are referred to comprehensive surveys on the topic, like [9]. On the output side, multi-agent planning requires an intermediate representation that can be passed to various vehicles in the scene, laying the focus on high-level behavioral planning approaches.

Imitation learning describes the application of supervised learning techniques to automated driving by training on expert drivers' demonstrations, which can be obtained from datasets or accordingly equipped testing vehicles. Based on object detections from a dedicated perception system and high-definition map information, a typical approach is to render the surroundings of the ego vehicle in a raster image that is subsequently processed by a convolutional neural network (CNN) [10], [11]. To address the problem of distributional shift between training data and closed-loop test conditions, various improvements have been proposed, like perturbing a random subset of training trajectories to teach the model to recover from atypical states [11]. Being based on supervised learning, those techniques share the large needs for high-quality training data. This limits their prospective transfer to cooperative multi-agent planning because ground-truth data showing cooperative maneuvers is virtually not available. Urban traffic datasets (e.g. the inD dataset [12]) show road users obeying to static priority rules or traffic lights; both entities that shall become obsolete with cooperative planning.

In contrast to supervised learning, RL approaches evade the requirement for large datasets by instead exploring possible actions in a simulated environment and exploiting a reward signal to learn the desired behavior. In [13], a driving policy for controlling the acceleration and steering angle is trained through RL that is applied to multiple vehicles in a common simulated environment. As there is no explicit communication between the different vehicles' policies, no cooperation is shown in traffic. Based on a raster image representation, this approach shares the unfavorable scaling of computational load with the number of participants in the scene, because each vehicle requires an individual image, centered on its pose for sensible inference. An alternative RL-based approach to coordinated driving on an urban intersection was published in [14]. By maintaining a tile-based reservation of the intersection, the decentralized policies can choose from the set of actions that do not cause a collision. Apart from this limitation of the action space, there is no further inter-agent communication that could enable cooperative maneuvers.

When encoding the semantic environment of a vehicle in urban traffic, the number of potentially relevant entities (e.g. other vehicles) is highly dynamic. This makes fixed-size network architectures and input representations often used in RL unsuitable for the task at hand. In [15], it is proposed to encode input features per vehicle using a multilayer perceptron followed by a permutation invariant operation for pooling the resulting features. The aggregated feature vector is then propagated through another fully connected network to finally infer actions for a single ego vehicle. The authors extend their work to encode whole traffic scenes including lanes and traffic signs and compare it to using a graph convolutional network for the same task in [16]. Similarly, [17] proposes to encode the vehicles being present in the scene as graph vertices. However, none of the described works can handle multi-agent planning.

## III. PROPOSED APPROACH

In this section, our proposed approach is presented, beginning with a discussion on learning paradigms for multi-agent usage. Afterwards, the graph-based input representation is introduced, followed by details on the network architecture and reward engineering.

### A. Learning Paradigm

In RL algorithms, an agent typically interacts with an environment in discrete time steps. The agent observes the current state of the environment and subsequently chooses an action, whose effect is evaluated by a reward signal. With multiple entities to be controlled, one can pursue different learning paradigms depending on the degree of centralization in multi-agent reinforcement learning [18]. Instead of having the various agents interact individually with the environment and learn independent policies, cooperative planning is modeled best by the centralized training centralized execution paradigm. Because different agents, which shall take part in a cooperative maneuver, have to be

able to communicate explicitly. This holds not only during training, but also at inference time. In paradigms relying on decentralized execution, the agents would have to learn an implicit communication scheme through their behavior. With the fused environment model being available to the server-side planner, considering the planning problem over all vehicles in the scene using a joint RL agent allows for explicit communication and hence better cooperation.

Like many RL problems, the cooperative planning problem can be denoted as a Markov decision process (MDP), defined as the tuple $(S, A, T, R)$. It consists of a set of states $S$ that fully describe the traffic scene at a given time. $A$ denotes the set of actions the RL agent can choose from while interacting with the environment. The transition function $T(s, a, s')$ describes the probability of changing from state $s \in S$ to $s' \in S$ when applying action $a \in A$, whereas the reward signal is determined by the function $R(s, a)$. Since the multi-agent planning problem contains different vehicles in the scene, the dimensionalities of the state space and the action space depend on the number of vehicles currently present and may vary over time.

### B. Input Representation

A well-suited input representation is crucial for applying artificial neural networks successfully. We identify three major requirements for an input scene representation to be used in cooperative multi-agent planning:

- Invariance on the number of vehicles in the scene,
- Permutation invariance of the input nodes,
- Permutation equivariance regarding the output nodes.

Simple tabular representations already lack the invariance properties. The limitations of fixed-sized inputs for behavior planning are elaborated more extensively in [15]. A rendered raster image of the scene, as often used for CNNs, fulfills the invariance requirements, but typically requires a target agent to be centered around [10]. This process must be repeated to produce individual outputs for each agent, making the application to a large number of vehicles computationally infeasible. Permutation equivariance means that the inferred outputs of given agents in the scene are independent of their ordering in the input vector. Hence, our work proposes to use a lean and flexible graph-based scene representation, shown in Fig. 2, which fulfills all above requirements. The current state of the environment is thus defined as $S = (V, E, U)$, with $V$ being the set of vertices corresponding to the vehicles in the scene and $E$ denoting edges depending on the pairwise relation between vehicles. For each vehicle, one vertex in $V$ stores the corresponding input features. Each of the directed edges is assigned one of two edge types in $U$, either *same lane*, or *crossing*:

$$(v_1, r, v_2) \in E = V \times U \times V. \tag{1}$$

Two vehicles in front of the intersection whose paths cross or merge are being connected bidirectionally with crossing edges, like $v_1$, $v_2$, and $v_3$ in the figure. The same lane edge, in contrast, is used to connect two vertices of vehicles on the same path, pointing from the predecessor to the following
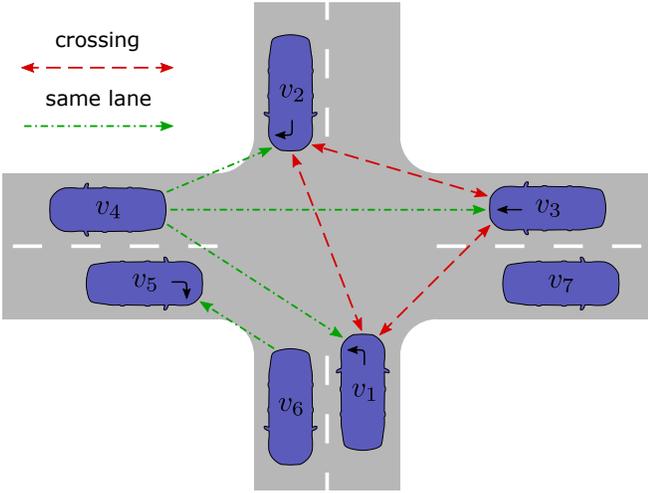
Fig. 2. The graph-based input representation illustrated on an arbitrary traffic scene at a four-way intersection. The vehicles' turning intentions are denoted by arrows on their hood.
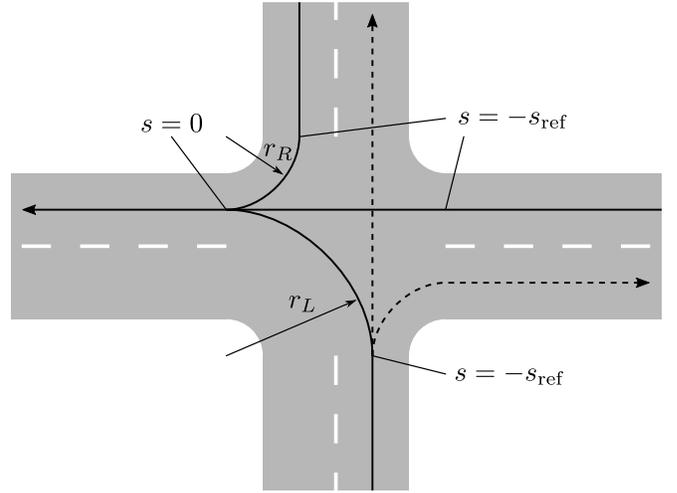
Fig. 3. Path parameterization by a common reference length. Each path enters the intersection area at $s = -s_{\text{ref}}$ and leaves it at $s = 0$.

one (e.g. $v_6$ and $v_5$). This is motivated by the observation that vehicles should adapt their behavior to the preceding vehicle and not vice-versa. Note that the graph does not need to be connected. Some vehicles may form a disjoint sub-graph, if they share no conflicts with the remaining vehicles, as it is the case for $v_5$ and $v_6$ or $v_7$ in Fig. 2.

The input feature vector for each vehicle consists of three values and is denoted as $h^{(0)} = [s, v, d]^T$, where the upper index denotes the layer number. The longitudinal position of the vehicle along its path is denoted by $s$, with $s = 0$ defined as the point where the path leaves the intersection area. This corresponds to the longitudinal coordinate of a Frenet coordinate pair. Because different maneuvers (e.g. straight driving and right turns) cause a difference in the path length on the intersection, the effective path length is scaled to a common reference length $s_{\text{ref}}$, as depicted in Fig. 3. Thereby, the entry point to the intersection area is located at $s = -s_{\text{ref}}$ consistently, ensuring that the localization on the incoming lanes is independent of the maneuver to be driven. Moreover, this normalization makes the scene representation robust to slight changes in intersection geometry.

The second input feature $v$ denotes the scalar velocity of the corresponding vehicle normalized over the speed limit of the lane it is currently driving on. To allow the network to sense immediate proximity of other vehicles, the input features are complemented by a distance measure $d$ based on the Mahalanobis distance [19]. The distance measured from vehicle $i$ to vehicle $j$ is calculated as

$$d_{ij} = \sqrt{(\boldsymbol{p}_j - \boldsymbol{p}_i)^T \boldsymbol{\Sigma}_i^{-1} (\boldsymbol{p}_j - \boldsymbol{p}_i)}, \quad (2)$$

where $\boldsymbol{p}_i$ denotes the position of vehicle $i$ in cartesian coordinates. The covariance matrix is given as

$$\boldsymbol{\Sigma}_i = \boldsymbol{R}_{\psi_i} \begin{bmatrix} l^2/4 & 0 \\ 0 & w^2/4 \end{bmatrix} \boldsymbol{R}_{\psi_i}^T, \quad (3)$$

with $l = 5\,\text{m}$ and $w = 2\,\text{m}$ describing the standardized length and width of a vehicle. $\boldsymbol{R}_{\psi_i}$ denotes the 2D rotation matrix

using the heading angle $\psi_i$. To determine the input feature for a particular vehicle, the distance to each other vehicle is computed according to (2), and the inverse of the minimum distance value is passed to the network. Using the inverse instead of the plain distance value proved to yield better model convergence.

### C. Network Architecture

In the present work, the behavioral control of vehicles is performed by applying a commanded longitudinal acceleration within the range $[a_{\min}, a_{\max}]$, requiring an RL algorithm suited for continuous control. We propose to use the twin delayed deep deterministic policy gradient (TD3) algorithm [20], an extension of the deep deterministic policy gradient (DDPG) [21]. Both methods are actor-critic RL algorithms for actions in continuous space. TD3 consists of two function approximators, namely the *actor* and the *critic*. Based on a given state and action input, the critic network is trained to predict the discounted reward as $Q$ value estimates. The actor gets the current environment state as an input and outputs an action to be performed in the particular time step, optimized by using the critic output as the loss.

The proposed graph neural network (GNN) architecture is depicted in Fig. 4. For each vertex, the low-dimensional input features are first processed by a dense layer `enc`. Note that this operation is performed individually for each vertex using shared parameters, disregarding the graph structure defined by the edges. The encoded vertex features are then propagated through two relational graph convolution layers, `conv_1` and `conv_2`. In contrast to simple graph convolution layers, multiple weight matrices corresponding to the different edge types are used for message passing [22]. During a forward pass, the hidden features of the vertices are propagated along the outgoing edges, while being multiplied by the respective weight matrix. Hence, each node receives a variable amount of such messages that have to be integrated into its own feature vector. This is done using a
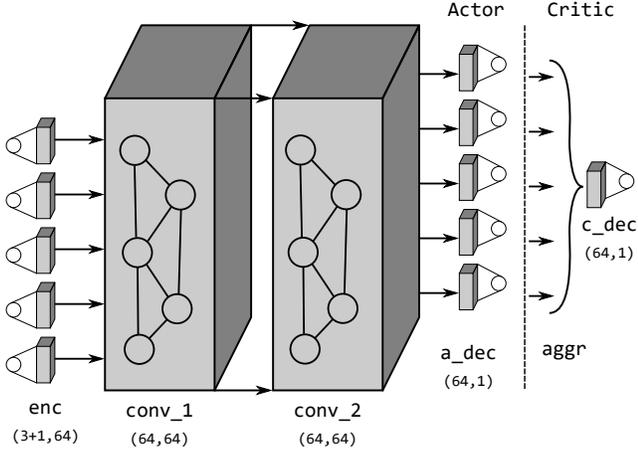
Fig. 4. The graph neural network architecture is depicted, consisting of one dense encoder layer, two graph convolution layers, and one dense output layer. Below the layer identifiers, their dimensionalities are shown. The encoder has three input channels for the actor and four channels for the critic network.

permutation invariant aggregation function like the element-wise maximum, mean, or sum. In the present work, the maximum operation delivered the best results. The update for the hidden feature vector of node $i$ is thus given as

$$h_i^{(l+1)} = \sigma \left( \Sigma_{r \in U} \max_{j \in \mathcal{N}_i^r} \boldsymbol{W}_r^{(l)} h_j^{(l)} + \boldsymbol{W}_0^{(l)} h_i^{(l)} \right), \quad (4)$$

where the set of neighbor nodes connected to the target node $i$ by incoming edges is denoted by $\mathcal{N}_i^r$. The weight matrices for each edge type $r \in U$ are called $\boldsymbol{W}_r$, while the previous target node vector is multiplied by $\boldsymbol{W}_0$. In case a vertex has no incoming edges (like $v_4$ or $v_7$ in Fig. 2), the node update coincides with a single dense layer on the node's own state. Finally, the resulting feature vector is passed through a non-linear activation function, empirically chosen as rectified linear unit (ReLU).

While the actor network and the critic network are constructed analogously up to the point described above, they differ in the output layer, depicted on the right side of Fig. 4. The actor network is responsible for deriving an action for each entity in the scene to be executed for a given time horizon, while the critic is in charge of estimating the Q value for the entire graph. With each vehicle being represented by a vertex in the graph, there is one regression target per vertex that describes the commanded acceleration for the corresponding vehicle in the actor network. The latent feature output by the GNN is reduced to a single unit using a final dense layer a_dec, whose weights are shared across nodes. To limit the action output to a defined range, a tangens hyperbolicus activation function is used on the output layer. The normalized value range is subsequently mapped to an acceleration between $-5\,\frac{\mathrm{m}}{\mathrm{s}^2}$ and $5\,\frac{\mathrm{m}}{\mathrm{s}^2}$. With the critic network's output being a performance measure in form of a single Q value estimate, an aggregation function for the latent feature vector of all nodes is required. The Q values'

range is not limited, hence a final dense layer a_dec with linear activation is used as the output layer. Because the critic network requires the chosen action in addition to the state representation, the action values are concatenated with all vertex input features.

### D. Reward Engineering

Apart from the network architecture described above, the RL algorithm requires a reward scheme to learn a reasonable behavior within the simulation environment. The reward signal is composed of a weighted sum of reward components

$$R = \sum_{k \in \mathcal{R}} w_k R_k, \quad (5)$$

where the set of reward components is given as $\mathcal{R} = \{\text{velocity, action, idle, proximity, collision}\}$. The velocity reward is the main driver for learning a non-trivial solution through rewarding large velocities and is defined as

$$R_{\text{velocity}} = \begin{cases} 1.25\frac{v}{v_{\text{lim}}} & \frac{v}{v_{\text{lim}}} \leq 0.8 \\ 1.0 & 0.8 < \frac{v}{v_{\text{lim}}} \leq 1.0 \\ 6.0 - 5.0\frac{v}{v_{\text{lim}}} & 1.0 < \frac{v}{v_{\text{lim}}}, \end{cases} \quad (6)$$

where $v_{\text{lim}}$ describes the vehicle's lane speed limit. Regularizing the model against applying large acceleration magnitudes is done by the action penalty that is defined as the negative absolute commanded acceleration. When striving to avoid collisions, the simplest solution is to stop the whole traffic, which is not desirable. Therefore, the idle penalty is set to $R_{\text{idle}} = 1$ in case all vehicles are standing still. To teach the model to keep suitable safety distances to nearby vehicles, the proximity component is used to penalize actions that cause two vehicles to get dangerously close. This penalty is calculated based on the aforementioned modified Mahalanobis distance measure (cf. (2)), which takes the relative direction of the obstacle into account. In the case that two vehicles collide, the collision penalty $R_{\text{collision}} = 1$ is used to let the model implicitly learn collision avoidance, while aborting the episode on the spot preventing further positive rewards from being accumulated.

## IV. EXPERIMENTS

Training and evaluation of RL algorithms require a suited simulation environment. For the task of behavioral planning in automated driving, the simulator should at least provide a kinematic vehicle model and reasonable interaction between vehicles. In the present study, the open-source environment Highway-env [23] is used and slightly adapted to be employed for centralized multi-agent planning. The simulation of vehicle kinematics is done according to the kinematic bicycle model [24], which suffices for behavioral planning. The graph-based scene representation and graph neural network layers are based on the PyTorch Geometric API [25].

The choice of reward weights is based on a grid search that was conducted on a reduced variant of the simulation environment resembling the key behavior, while being much less computationally demanding. The reward weights used

| Reward | velocity | action | idle | proximity | collision |
|--------|----------|--------|------|-----------|-----------|
| **Weight** | 0.03 | 0.01 | 0.01 | 0.2 | 1.0 |



Fig. 5. Flow rate during various evaluation runs by using static priority rules (PR), the FIFO policy, and the RL planner.

throughout this study are given in Table I. During training, the latest model is evaluated on a separate validation environment for ten episodes every 5000 time steps. The validation environment is constructed the same way as the training environment but initialized with a different seed. Each time the validation shows a new best validation reward, the current model parameters are saved to disk.

We benchmark our approach against two baselines in synthetic simulation: static priority rules (PR) and a first-in-first-out scheme (FIFO), resembling the currently prevalent approach in real-world and a seminal AIM scheme. Traffic obeying to static priority rules is simulated using the driver models provided by Highway-env, which rely on the intelligent driver model (IDM) [26] for longitudinal control and additional logic for handling intersections. We applied minor tweaks to the driver models to obtain reasonable results for more dense traffic:

- The derivation of the commanded acceleration is modified to correctly handle a target speed of zero.
- Scheduling at an intersection is based on vehicle priorities that are inferred from their intended maneuver instead of the current lane priority.
- The yielding logic is extended to respect a specific stop point in front of the intersection.

The FIFO scheme, on the other hand, prioritizes the incoming vehicles based on their distance to the intersection. Thereby, non-conflicting paths can be used at the same time, possibly allowing multiple vehicles on the intersection at a given time. Note that this policy does not enforce a strict FIFO ordering on the whole intersection, but rather groups of conflicting paths, which leads to a considerable performance increase.

The way of generating simulated traffic differs between training and evaluation runs. During training, vehicles are being spawned on all incoming lanes featuring enough space at a certain probability with their destination also being chosen randomly. In the course of training, the spawn probability is continuously increased until it saturates at $5\%$ per time step. Thereby, the intersection is kept busy and allows the RL algorithm to obtain meaningful data samples to learn from. The evaluation runs are based on scenario definitions that are generated using a slightly different scheme. In that case, the time between vehicles appearing on a particular lane is governed by a shifted exponential distribution. This resembles a Poisson process, where the shift on the distribution of the spawn period ensures a minimum distance between vehicles. If a traffic jam has formed so that there is no space for a vehicle to be spawned, the generation is suspended to prevent immediate collisions. For each scenario to be generated, the desired vehicle rate is chosen randomly from a uniform distribution over the interval $[0.2, 0.4]$ vehicles per
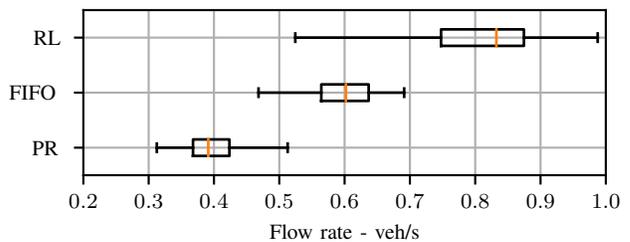
second and major road lane. The vehicle rate on the minor road lanes is set to half of that value. In both training and evaluation, the initial vehicle velocity is chosen uniformly between $60\%$ and $100\%$ of the corresponding lane speed limit.

*A. Synthetic Simulation*

All experiments within this subsection were performed on a four-way intersection, whose layout is analogue to the one depicted in Fig. 2. The *flow rate* describes the number of vehicles that cross an intersection (or other road infrastructure) during a given time frame. Figure 5 shows the flow rate distribution over 100 evaluation runs (each of $100\,\mathrm{s}$ length) of varying traffic density. It can be observed that already the rather simple FIFO scheme achieves a benefit over static priority rules, while the learned RL planner outperforms both baselines regarding the median values. To further investigate the performance gain, we analyze the ratio of vehicles that had to stop during the maneuver. A vehicle trajectory is considered to contain a stop, if the velocity falls below $0.3\,\frac{\mathrm{m}}{\mathrm{s}}$ for at least one time step. This threshold was chosen due to numeric reasons. By categorizing the vehicles by their incoming road priority, the effect on traffic approaching from a minor road becomes apparent, as depicted in Fig. 6. Clearly, the static priority rules induce a significant traffic buildup on the minor road that forces nearly all vehicles to stop. The FIFO policy manages to let more vehicles from the minor road pass the intersection, but in turn causes a large proportion of stops also on the major road. In contrast, the RL planner succeeds to get a large amount of vehicles across the intersection, while keeping the traffic flow mostly intact. This behavior might be explained by the planner's learned ability to adapt the vehicles' velocities early to fit into an emergent gap. It should be mentioned that the RL planner cannot completely eliminate collisions, as denoted in the first row of Table II. However, the occurrence is extremely rare, making it very challenging to further reduce them, given the RL algorithm has to implicitly learn it via the reward signal. In practice, the remaining failure cases are not an issue, because the cooperative maneuver will only be advertised to the connected vehicles, if it fulfills sanity checks like being collision-free. In the case no viable cooperative plan was found, the vehicles simply resort to local planning.
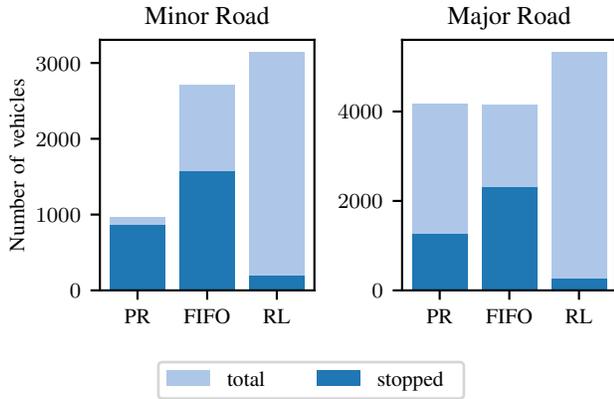
Fig. 6. The number of vehicles approaching the intersection from the major and minor road as well as the ratio of those that had to stop when governed by static priority rules (PR), the FIFO policy, and the RL planner.



Fig. 8. The number of vehicles that perform a certain maneuver and ratio of those that had to stop in real-world data (inD), using the FIFO policy, and the RL planner.

| Intersection | Priority rules | FIFO scheme | RL planner |
|---|---|---|---|
| Synthetic 4-way | 0.0 % | 0.0 % | 0.028 % |
| inD | 0.0 % | 1.918 % | 0.584 % |

## B. Simulation Based on Real-World Traffic Data

Apart from the simulation based on synthetic data, we also evaluated the cooperative planning scheme on real-world urban traffic data taken from the inD dataset [12]. The dataset contains tracks of vehicles and vulnerable road users that were recorded at four urban intersections in Germany. We selected a four-way intersection connecting a priority road with a minor road that is managed by static priority rules. The major road also features isolated lanes for turning left, as depicted in Fig. 7. Simulating the traffic according to the cooperative planning approach on this intersection makes the following assumptions inevitable. Firstly, the intersection geometry is only approximated in simulation. However, this is not an issue for behavioral planning, which is mostly independent of road geometry. Considering the road curvature, vehicles may not traverse it with arbitrary speed, which is
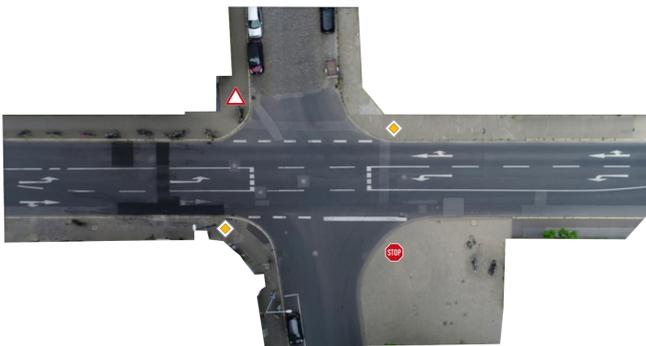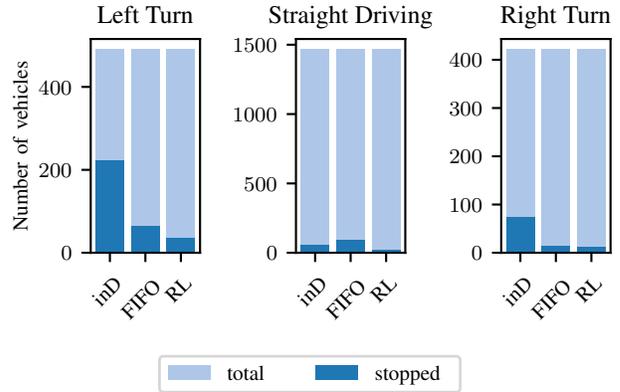


Fig. 7. Bird's-eye view of the intersection, where the real-world vehicle tracks taken from the inD dataset were collected (image adopted from [12]).

ensured by defining lane-dependent velocity limits. As the dynamic properties of the vehicles shown in the dataset are unknown, a default parameter set is used in simulation. Moreover, the real-world tracks deviate from the lane center lines that are used for guiding the simulated vehicles. The recorded intersection is also used by vulnerable road users like pedestrians and bicyclists that cannot be modeled in the simulation as of now.

Compared to the synthetic simulation, not all metrics are viable for evaluation when using the dataset as the baseline. The flow rate, for instance, cannot be improved by any intersection management system, because the number of vehicles in the scene is specified by the dataset. In total, the used dataset excerpt provides 2446 vehicle tracks over a recording time of 3.08 hours, resulting in an average traffic density of $794 \frac{\text{veh}}{\text{h}}$ or $0.221 \frac{\text{veh}}{\text{s}}$. Compared to the flow rates obtained in the synthetic simulations, those numbers are rather small, which might indicate that there are not many interesting situations during most of the recording. Evaluating the RL planner and FIFO policy on real-world data is performed by spawning vehicles in simulation according to the appearance time in the dataset and subsequently simulating their motion based on the vehicle models.

Figure 8 compares the number of vehicles that had to come to a complete stop categorized into the maneuvers left turn, straight driving, and right turn. The total number of vehicles being managed by the FIFO policy and the RL planner has to be identical to the amount given by the dataset recording. It is clearly visible from the real-world data that many left turning vehicles have to come to a stop before being able to safely pass the intersection. Note that although one minor road access is governed by a stop sign, the real-world recordings show that by far not all road users obey to it. The FIFO scheme naturally distributes the stops to all maneuvers including vehicles driving straight on. Meanwhile, the RL planner is able to avoid the vast majority of stops and maintains a smooth flow of traffic.

As it can be seen in the last row of Table II, both the

FIFO and the RL planner suffer from an increased collision rate. This can, at least in parts, be attributed to the way the real-world tracks are mapped to simulation. Especially the minor road is only depicted for a very short distance in front of the intersection (cf. Fig. 7), which makes it difficult for the planner to influence the incoming vehicles before they enter the intersection area. In case one vehicle waits at the stop point while a second vehicle enters the scene at high speed, a collision might be inevitable. This is merely an issue of the evaluation and does not diminish the remarkable improvement in traffic efficiency that is raised by the RL planner.

## V. Conclusion

In this work, a novel multi-agent behavioral planning scheme for connected automated vehicles at urban intersections was presented. We chose a reinforcement learning algorithm to leverage recent advances in machine learning while evading the need for ground truth data that is virtually unavailable for cooperative maneuvers. The developed graph-based input representation effectively encodes the semantic environment at the operational area. By employing graph neural networks, our approach confidently handles the varying number of vehicles in the scene. The proposed approach was evaluated in synthetic simulation and additionally based on real-world traffic data. Compared to static priority rules and a FIFO scheme as baselines, the learned planner increases the vehicle throughput significantly. In addition, the number of induced stops is reduced which indicates better traffic flow.

The proposed behavioral planning framework can serve as a sound foundation for solving more sophisticated planning problems. In the future, we plan to extend this work to be applicable to intersection layouts that were not seen during training. Moreover, cooperative planning in mixed traffic, i.e. human drivers and automated vehicles sharing the road, shall be addressed.

## References

[1] M. Buchholz, J. C. Müller, M. Herrmann, J. Strohbeck, B. Völz, M. Maier, J. Paczia, O. Stein, H. Rehborn, and R.-W. Henn, "Handling Occlusions in Automated Driving Using a Multiaccess Edge Computing Server-Based Environment Model From Infrastructure Sensors," *IEEE Intelligent Transportation Systems Magazine*, to appear, doi: 10.1109/MITS.2021.3089743.

[2] Z. Zhong, M. Nejad, and E. E. Lee, "Autonomous and Semi-Autonomous Intersection Management: A Survey," *IEEE Intelligent Transportation Systems Magazine*, 2020.

[3] K. Dresner and P. Stone, "A Multiagent Approach to Autonomous Intersection Management," *Journal of Artificial Intelligence Research*, vol. 31, pp. 591–656, Mar. 2008.

[4] A. A. Malikopoulos, C. G. Cassandras, and Y. J. Zhang, "A Decentralized Energy-Optimal Control Framework for Connected Automated Vehicles at Signal-Free Intersections," *Automatica*, vol. 93, pp. 244–256, Jul. 2018.

[5] M. A. S. Kamal, J.-i. Imura, T. Hayakawa, A. Ohata, and K. Aihara, "A Vehicle-Intersection Coordination Scheme for Smooth Flows of Traffic Without Using Traffic Lights," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 3, pp. 1136–1147, Jun. 2015.

[6] M. B. Mertens, J. Müller, and M. Buchholz, "Cooperative Maneuver Planning for Mixed Traffic at Unsignalized Intersections Using Probabilistic Predictions," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, 2022, to appear.

[7] A. I. Morales Medina, N. van de Wouw, and H. Nijmeijer, "Cooperative Intersection Control Based on Virtual Platooning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 6, pp. 1727–1740, Jun. 2018.

[8] Z. Zhu and H. Zhao, "A Survey of Deep RL and IL for Autonomous Driving Policy Learning," *IEEE Transactions on Intelligent Transportation Systems*, to be published, doi: 10.1109/TITS.2021.3134702.

[9] S. Lefèvre, D. Vasquez, and C. Laugier, "A Survey on Motion Prediction and Risk Assessment for Intelligent Vehicles," *ROBOMECH Journal*, vol. 1, no. 1, Dec. 2014.

[10] J. Chen, B. Yuan, and M. Tomizuka, "Deep Imitation Learning for Autonomous Driving in Generic Urban Scenarios with Enhanced Safety," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, Nov. 2019, pp. 2884–2890.

[11] M. Bansal, A. Krizhevsky, and A. Ogale, "ChauffeurNet: Learning to Drive by Imitating the Best and Synthesizing the Worst," in *Robotics: Science and Systems*. Robotics: Science and Systems Foundation, 2019.

[12] J. Bock, R. Krajewski, T. Moers, S. Runde, L. Vater, and L. Eckstein, "The inD Dataset: A Drone Dataset of Naturalistic Road User Trajectories at German Intersections," in *2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, Oct. 2020, pp. 1929–1934.

[13] A. P. Capasso, P. Maramotti, A. Dell'Eva, and A. Broggi, "End-to-End Intersection Handling using Multi-Agent Deep Reinforcement Learning," in *2021 IEEE Intelligent Vehicles Symposium (IV)*, 2021, pp. 443–450.

[14] Y. Wu, H. Chen, and F. Zhu, "DCL-AIM: Decentralized Coordination Learning of Autonomous Intersection Management for Connected and Automated Vehicles," *Transportation Research Part C: Emerging Technologies*, vol. 103, pp. 246–260, Jun. 2019.

[15] M. Huegle, G. Kalweit, B. Mirchevska, M. Werling, and J. Boedecker, "Dynamic Input for Deep Reinforcement Learning in Autonomous Driving," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, Nov. 2019, pp. 7566–7573.

[16] M. Huegle, G. Kalweit, M. Werling, and J. Boedecker, "Dynamic Interaction-Aware Scene Understanding for Reinforcement Learning in Autonomous Driving," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, May 2020, pp. 4329–4335.

[17] P. Hart and A. Knoll, "Graph Neural Networks and Reinforcement Learning for Behavior Generation in Semantic Environments," in *2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, Oct. 2020, pp. 1589–1594.

[18] S. Gronauer and K. Diepold, "Multi-Agent Deep Reinforcement Learning: A Survey," *Artificial Intelligence Review*, Apr. 2021.

[19] R. De Maesschalck, D. Jouan-Rimbaud, and D. Massart, "The Mahalanobis Distance," *Chemometrics and Intelligent Laboratory Systems*, vol. 50, no. 1, pp. 1–18, Jan. 2000.

[20] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing Function Approximation Error in Actor-Critic Methods," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, Jul. 2018, pp. 1587–1596.

[21] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous Control with Deep Reinforcement Learning," in *4th International Conference on Learning Representations, ICLR 2016*, San Juan, Puerto Rico, 2016.

[22] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling, "Modeling Relational Data with Graph Convolutional Networks," in *The Semantic Web*, A. Gangemi, R. Navigli, M.-E. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai, and M. Alam, Eds. Cham: Springer International Publishing, 2018, vol. 10843, pp. 593–607.

[23] E. Leurent, "An Environment for Autonomous Driving Decision-Making," May 2018. [Online]. Available: https://github.com/eleurent/highway-env

[24] J. Kong, M. Pfeiffer, G. Schildbach, and F. Borrelli, "Kinematic and Dynamic Vehicle Models for Autonomous Driving Control Design," in *2015 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, Jun. 2015, pp. 1094–1099.

[25] M. Fey and J. E. Lenssen, "Fast Graph Representation Learning with PyTorch Geometric," May 2019. [Online]. Available: https://github.com/pyg-team/pytorch_geometric

[26] M. Treiber, A. Hennecke, and D. Helbing, "Congested Traffic States in Empirical Observations and Microscopic Simulations," *Physical Review E*, vol. 62, no. 2, pp. 1805–1824, Aug. 2000.