# Dynamic Interaction-Aware Scene Understanding for Reinforcement Learning in Autonomous Driving

Maria Huegle[1], Gabriel Kalweit[1], Moritz Werling[2] and Joschka Boedecker[1,3]

*Abstract*— The common pipeline in autonomous driving systems is highly modular and includes a perception component which extracts lists of surrounding objects and passes these lists to a high-level decision component. In this case, leveraging the benefits of deep reinforcement learning for high-level decision making requires special architectures to deal with multiple variable-length sequences of different object types, such as vehicles, lanes or traffic signs. At the same time, the architecture has to be able to cover interactions between traffic participants in order to find the optimal action to be taken. In this work, we propose the novel Deep Scenes architecture, that can learn complex interaction-aware scene representations based on extensions of either 1) Deep Sets or 2) Graph Convolutional Networks. We present the Graph-Q and DeepScene-Q off-policy reinforcement learning algorithms, both outperforming state-of-the-art methods in evaluations with the publicly available traffic simulator SUMO.

## I. INTRODUCTION

In autonomous driving scenarios, the number of traffic participants and lanes surrounding the agent can vary considerably over time. Common autonomous driving systems use modular pipelines, where a perception component extracts a list of surrounding objects and passes this list to other modules, including localization, mapping, motion planning and high-level decision making components. Classical rule-based decision-making systems are able to deal with variable-sized object lists, but are limited in terms of generalization to unseen situations or are unable to cover all interactions in dense traffic. Since Deep Reinforcement Learning (DRL) methods can learn decision policies from data and off-policy methods can improve from previous experience, they offer a promising alternative to rule-based systems. In the past years, DRL has shown promising results in various domains [1], [2], [3], [4], [5]. However, classical DRL architectures like fully-connected or convolutional neural networks (CNNs) are limited in their ability to deal with variable-sized, structured inputs or to model interactions between objects.

Prior works on reinforcement learning for autonomous driving that used fully-connected network architectures and fixed sized inputs [6], [7], [5], [8], [9] are limited in the number of vehicles that can be considered. CNNs using occupancy grids [10], [11] are limited to their initial grid size. Recurrent neural networks are useful to cover temporal

[1,3]Dept. of Computer Science, University of Freiburg, Germany.
{hueglem,kalweitg,jboedeck}@cs.uni-freiburg.de
[2]BMWGroup, Unterschleissheim, Germany.
Moritz.Werling@bmw.de
[3]Cluster of Excellence BrainLinks-BrainTools, Freiburg, Germany.

context, but are not able to handle a variable number of objects permutation-invariant w.r.t to the input order for a fixed time step. In [12], limitations of these architectures are shown and a more flexible architecture based on Deep Sets [13] is proposed for off-policy reinforcement learning of lane-change maneuvers, outperforming traditional approaches in evaluations with the open-source simulator SUMO.

In this paper, we propose to use Graph Networks [14] as an interaction-aware input module in reinforcement learning for autonomous driving. We employ the structure of Graphs in off-policy DRL and formalize the Graph-Q algorithm. In addition, to cope with multiple object classes of different feature representations, such as different vehicle types, traffic signs or lanes, we introduce the formalism of Deep Scenes, that can extend Deep Sets and Graph Networks to fuse multiple variable-sized input sets of different feature representations. Both of these can be used in our novel DeepScene-Q algorithm for off-policy DRL. Our main contributions are:

1) Using Graph Convolutional Networks to model interactions between vehicles in DRL for autonomous driving.
2) Extending existing set input architectures for DRL to deal with multiple lists of different object types.

## II. RELATED WORK

Graph Networks are a class of neural networks that can learn functions on graphs as input [15], [16], [17], [18], [19] and can reason about how objects in complex systems interact. They can be used in DRL to learn state representations [20], [21], [22], [17], e.g. for inference and control of physical systems with bodies (objects) and joints (relations). In the application for autonomous driving, Graph Networks were used for supervised traffic prediction while modeling traffic participant interactions [23], where vehicles were modeled as objects and interactions between them as relations. Another type of interaction-aware network architectures, Interaction Networks, were proposed to reason about how objects in complex systems interact [18]. A vehicle behavior interaction network that captures vehicle interactions was presented in [24]. In [25], a convolutional social pooling component was proposed using a CNN to model spatial connections between vehicles for vehicle trajectory prediction.

## III. PRELIMINARIES

We model the task of high-level decision making for autonomous driving as a Markov Decision Process (MDP), where the agent is following a policy $\pi$ in an environment
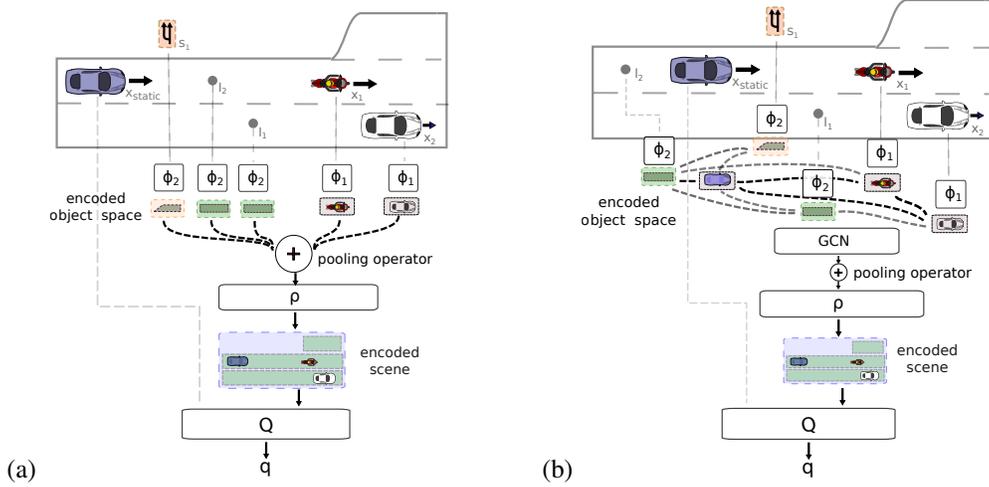
Fig. 1. Scheme of DeepScene-Q, using (a) Deep Sets and (b) Graphs. Both architectures combine multiple variable-length object lists in a scene, here a traffic sign $s_1$, lanes $l_1, l_2$ and vehicles $x_1, x_2$. The modules $\phi_i$, $\rho$ and $Q$ are fully-connected networks. As permutation invariant pooling operator, we use the *sum*. The vector $x^{\text{static}}$ includes static features and $q$ the action value output.

in a state $s_t$, applying a discrete action $a_t \sim \pi$ to reach a successor state $s_{t+1} \sim \mathcal{M}$ according to a transition model $\mathcal{M}$. In every time step $t$, the agent receives a reward $r_t$, e.g. for driving as close as possible to a desired velocity. The agent tries to maximize the discounted long-term return $R(s_t) = \sum_{i>=t} \gamma^{i-t} r_i$, where $\gamma \in [0,1]$ is the discount factor. In this work, we use Q-learning [26]. The Q-function $Q^\pi(s_t, a_t) = \mathbf{E}_{a_{i>t} \sim \pi}[R(s_t)|a_t]$ represents the value of following a policy $\pi$ after applying action $a_t$. The optimal policy can be inferred from the optimal action-value function $Q^*$ by maximization over actions.

### A. Q-Function Approximation

We use DQN [1] to estimate the optimal $Q$-function by function approximator $Q$, parameterized by $\theta^Q$. It is trained in an offline fashion on minibatches sampled from a fixed replay buffer $\mathcal{R}$ with transitions collected by a driver policy $\hat{\pi}$. As loss, we use $L(\theta^Q) = \frac{1}{b} \sum_i \left(y_i - Q(s_i, a_i|\theta^Q)\right)^2$ with targets $y_i = r_i + \gamma \max_a Q'(s_{i+1}, a|\theta^{Q'})$, where $Q'$ is a target network, parameterized by $\theta^{Q'}$, and $(s_i, a_i, s_{i+1}, r_i)|_{0 \leq i \leq b}$ is a randomly sampled minibatch from $\mathcal{R}$. For the target network, we use a soft update, i.e. $\theta^{Q'} \leftarrow \tau \theta^Q + (1-\tau)\theta^{Q'}$ with update step-size $\tau \in [0,1]$. Further, we use a variant of Double-$Q$-learning [27] which is based on two $Q$-network pairs and uses the minimum of the predictions for the target calculation, similar as in [28].

### B. Deep Sets

A network $Q_{\mathcal{DS}}$ can be trained to estimate the $Q$-function for a state representation $s = (X^{\text{dyn}}, x^{\text{static}})$ and action $a$. The representation consists of a static input $x^{\text{static}}$ and a dynamic, variable-length input set $X^{\text{dyn}} = [x^1, .., x^{\text{seq len}}]^\top$, where $x^j|_{1 \leq j \leq \text{seq len}}$ are feature vectors for surrounding vehicles in sensor range. In [12], it was proposed to use Deep Sets to handle this input representation, where the Q-network consists of three network modules $\phi, \rho$ and $Q$. The representation of the dynamic input set is computed by

$\Psi(X^{\text{dyn}}) = \rho\left(\sum_{x \in X^{\text{dyn}}} \phi(x)\right)$, which makes the Q-function permutation invariant w.r.t. the order of the dynamic input [13]. Static feature representations $x^{\text{static}}$ are fed directly to the Q-module, and the Q-values can be computed by $Q_{\mathcal{DS}} = Q(\Psi(X^{\text{dyn}})||x^{\text{static}})$, where $||$ denotes a concatenation of two vectors. The Q-learning algorithm is called DeepSet-Q [12].

## IV. METHODS

### A. Deep Scene-Sets

To overcome the limitation of DeepSet-Q to one variable-sized list of the same object type, we propose a novel architecture, Deep Scene-Sets, that are able to deal with $K$ input sets $X^{\text{dyn}_1}, ..., X^{\text{dyn}_K}$, where every set has variable length. A combined, permutation invariant representation of all sets can be computed by

$$\Psi(X^{\text{dyn}_1}, ..., X^{\text{dyn}_K}) = \rho\left(\sum_k \sum_{x \in X^{\text{dyn}_k}} \phi^k(x)\right),$$

where $1 \leq k \leq K$. The output vectors $\phi^k(\cdot) \in \mathbb{R}^F$ of the neural network modules $\phi^k$ have the same length $F$. We additionally propose to share the parameters of the last layer for the different $\phi$ networks. Then, $\phi^k(\cdot)$ can be seen as a projection of all input objects to the same encoded *object space*. We combine the encoded objects of different types by the *sum* (or other permutation invariant pooling operators, such as *max*) and use the network module $\rho$ to create an encoded *scene*, which is a fixed-sized vector. The encoded scene is concatenated to $x^{\text{static}}$ and the Q-values can be computed by $Q_{\mathcal{D}} = Q(\Psi(X^{\text{dyn}_1}, ..., X^{\text{dyn}_K})||x^{\text{static}})$. We call the corresponding Q-learning algorithm DeepScene-Q, shown in Algorithm 2 (Option 1) and Figure 1 (a).

### B. Graphs

In the Deep Set architecture, relations between vehicles are not explicitly modeled and have to be inferred in $\rho$. We extend this approach by using Graph Networks, considering

**Algorithm 1:** Graph-Q

1  initialize $Q_{\mathcal{G}} = (\phi, \rho, H, Q)$ and $Q'_{\mathcal{G}} = (\phi', \rho', H', Q')$, set replay buffer $\mathcal{R}$
2  **for** *optimization step o=1,2,...* **do**
3     get minibatch $(s_i, a_i, (X^{\text{dyn}}_{i+1}, x^{\text{static}}_{i+1}), r_{i+1})$ from $\mathcal{R}$
4     **foreach** *transition* **do**
5        **foreach** *object* $x^j_{i+1}$ *in* $X^{dyn}_{i+1}$ **do**
6           $(\phi'_{i+1})^j = \phi'\left(x^j_{i+1}\right)$
7        compute $H'^{(L)}_{i+1}$ by GCN with $H'^{(0)}_{i+1} = [(\phi'_{i+1})^1, ..., (\phi'_{i+1})^{\text{seq len}}]^\top$
8        get $\rho'_{i+1} = \rho'\left(\sum_k \sum_j H'^{(L)}_{i+1}\right)$
9     $y_i = r_{i+1} + \gamma \max_a Q'(\rho'_{i+1}, x^{\text{static}}_{i+1}, a)$
10    perform a gradient step on loss: $\frac{1}{b}\sum_i (Q_{\mathcal{G}}(s_i, a_i) - y_i)^2$
11    update target network by: $\theta^{Q_{\mathcal{G}}} \leftarrow \tau\theta^{Q_{\mathcal{G}}} + (1-\tau)\theta^{Q'_{\mathcal{G}}}$

---

**Algorithm 2:** DeepScene-Q

1  initialize $Q_{\mathcal{D}} = (\phi^1, ..., \phi^K, \rho, H, Q)$ and $Q'_{\mathcal{D}} = (\phi^{1'}, ..., \phi^{K'}, \rho', H', Q')$, set replay buffer $\mathcal{R}$
2  **for** *optimization step o=1,2,...* **do**
3     get minibatch $(s_i, a_i, (X^{\text{dyn}_1}_{i+1}, ..., X^{\text{dyn}_K}_{i+1}, x^{\text{static}}_{i+1}), r_{i+1})$ from $\mathcal{R}$
4     **foreach** *transition* **do**
5        **foreach** *object type* $k \in (1, ..., K)$ **do**
6           **foreach** *object* $x^j_{i+1}$ *in* $X^{dyn_k}_{i+1}$ **do**
7              $(\phi^{k}_{i+1}{'})^j = \phi^{k'}\left(x^j_{i+1}\right)$
8        *Set (Option 1):*
9        get $\rho'_{i+1} = \rho'\left(\sum_k \sum_j (\phi^{k}_{i+1}{'})^j\right)$
10       *Graph (Option 2):*
11       compute $H'^{(L)}_{i+1}$ by GCN with $H'^{(0)}_{i+1} = [\Phi^1, ..., \Phi^K]^\top$ and $\Phi^k = [(\phi'_{i+1})^1, ..., (\phi'{}_{i+1})^{\text{seq len}}]$
12       get $\rho'_{i+1} = \rho'\left(\sum_k \sum_j H'^{(L)}_{i+1}\right)$
13    $y_i = r_{i+1} + \gamma \max_a Q'(\rho'_{i+1}, x^{\text{static}}_{i+1}, a)$
14    perform a gradient step on loss and update target network as in Algorithm 1.

---

graphs as input. Graph Convolutional Networks (GCNs) [14] operate on graphs defined by a set of node features $X^{\text{dyn}} = [x^1, .., x^{\text{seq len}}]^\top$ and a set of edges represented by an adjacency matrix $A$. The propagation rule of the GCN is $H^{(l)} = \sigma(D^{\frac{1}{2}}\tilde{A}D^{\frac{1}{2}}H^{(l-1)}W^{(l-1)})$ with $1 \leq l \leq L$, where we set $H^{(0)} = [\phi(x_1), ..., \phi(x_{\text{seq len}})]^\top$ using an encoder module similar as in the Deep Sets approach. $\tilde{A} \in \mathbb{R}^{N \times N}$ is an adjacency matrix with added self-connections, $D_{i,i} = \sum_j \tilde{A}_{i,j}$, $\sigma$ the activation function, $H^{(l)} \in \mathbb{R}^{N \times F}$ hidden layer activations and $W^{(l)}$ the learnable matrix of the $l$-th layer. The dynamic input representation can be computed from the last layer $L$ of the GCN: $\Psi(X^{\text{dyn}}) = \rho\left(\sum_{x \in X^{\text{dyn}}} H^{(L)}\right)$, where $\phi$ is a neural network and the output vector $\phi(\cdot) \in \mathbb{R}^F$ has length $F$. The Q-values can be computed by $Q_{\mathcal{G}} = Q(\Psi(X^{\text{dyn}})||x^{\text{static}})$. We call the corresponding Q-learning algorithm Graph-Q, see Algorithm 1.

*C. Deep Scene-Graphs*

The graph representation can be extended to deal with multiple variable-length lists of different object types $X^{\text{dyn}_1}, ..., X^{\text{dyn}_K}$ by using $K$ encoder networks. As node features, we use $H^{(0)} = [\Phi^1, ..., \Phi^K]^\top$ and $\Phi^k = [\phi^k(x_1), ..., \phi^k(x_{\text{seq len}_k})]$ for $1 \leq k \leq K$, and compute the dynamic input representation from the last layer of the GCN:

$$\Psi(X^{\text{dyn}_1}, ..., X^{\text{dyn}_K}) = \rho\left(\sum_k \sum_{x \in X^{\text{dyn}_k}} H^{(L)}\right),$$

with $1 \leq k \leq K$. Similar to the Deep Scene-Sets architecture, $\phi^k$ are neural network modules with output vector length $D$ and parameter sharing in the last layer. To create a fixed vector representation, we combine all node features by the sum into an encoded scene. The Q-values can be computed by $Q_{\mathcal{D}} = Q(\Psi(X^{\text{dyn}_1}, ..., X^{\text{dyn}_K})||x^{\text{static}})$. This module can replace the DeepScene-Sets module in DeepScene-Q as shown in Algorithm 2 (Option 2) and in Figure 1 (b).

*D. Graph Construction*

We propose two different strategies to construct bidirectional edge connections between vehicles for Graphs and

Deep Scene-Graphs representations:

1) Close agent connections: Connect agent vehicle to its direct leader and follower in its own and the left and right neighboring lanes ($6 \cdot 2$ edges).
2) All close vehicles connections: Connect all vehicles to their leader and follower in their own and the left and right lanes ($K \cdot 6 \cdot 2$ edges for $K$ surrounding vehicles).

Edge weights are computed by the inverse absolute distance between two vehicles, as shown in [23]. A fully-connected graph is avoided due to computational complexity.

*E. MDP Formulation*

The feature representations of the the surrounding cars and lanes are shown in section V-B. The action space $\mathcal{A}$ consists of a discrete set of three possible actions in lateral direction: *keep lane*, *left lane-change* and *right lane-change*. Acceleration and collision avoidance are controlled by low-level controllers, that are fixed and not updated during training. Maintaining safe distance to the preceding vehicle is handled by an integrated safety module, as proposed in [11], [5]. If the chosen lane-change action is not safe, the agent keeps the lane. The reward function $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is defined as: $r(s, a) = 1 - \frac{|v_{\text{current}}(s) - v_{\text{desired}}(s)|}{v_{\text{desired}}(s)} - p_{\text{lc}}(a)$, where $v_{\text{current}}$ and $v_{\text{desired}}$ are the actual and desired velocity of the agent, $p_{\text{lc}}$ is a penalty for choosing a lane-change action and minimizing lane-changes for additional comfort.

## V. EXPERIMENTAL SETUP

We use the open-source SUMO [29] traffic simulation to learn lane-change maneuvers.

| Driver Type | maxSpeed | lcCooperative | accel/ decel | length | lcSpeedGain |
|---|---|---|---|---|---|
| agent driver | 10 | - | 2.6/4.5 | 4.5 | - |
| passenger drivers 1 | $\mathcal{U}(8,12)$ | 0.2 | 2.6/4.5 | $\mathcal{U}(4,5)$ | $\mathcal{U}(5,10)$ |
| passenger drivers 2 | $\mathcal{U}(5,9)$ | 1.0 | 2.6/4.5 | $\mathcal{U}(4,5)$ | $\mathcal{U}(5,10)$ |
| passenger drivers 3 | $\mathcal{U}(3,7)$ | 0.8 | 2.6/4.5 | $\mathcal{U}(4,5)$ | $\mathcal{U}(5,10)$ |
| truck drivers | $\mathcal{U}(2,4)$ | 0.4 | 1.3 / 2.25 | $\mathcal{U}(9.5,14.5)$ | $\mathcal{U}(0,3)$ |
| motorcycle drivers | $\mathcal{U}(7,11)$ | 0.2 | 3.0/5.0 | $\mathcal{U}(2,3)$ | $\mathcal{U}(15,20)$ |

TABLE I

SUMO PARAMETERS FOR DIFFERENT DRIVER TYPES. IN EACH SCENARIO, TRUCKS AND MOTORCYCLES ARE SAMPLED WITH 10% AND 5%
PROBABILITY, PASSENGER CARS AND THEIR DRIVER TYPES ARE SAMPLED UNIFORMLY FOR THE REMAINING NUMBER OF VEHICLES.

### A. Scenarios

*a) Highway:* To evaluate and show the advantages of Graph-Q, we use the $1000\,\mathrm{m}$ circular highway environment shown in [12] with three continuous lanes and one object class (passenger cars). To train our agents, we used a dataset with 500.000 transitions.

*b) Fast Lanes:* To evaluate the performance of DeepScene-Q, we use a more complex scenario with a variable number of lanes, shown in Figure 2. It consists of a $1000\,\mathrm{m}$ circular highway with three continuous lanes and additional fast lanes in two $250\,\mathrm{m}$ sections. At the end of lanes, vehicles slow down and stop until they can merge into an ongoing lane. The agent receives information about additional lanes in form of traffic signs starting $200\,\mathrm{m}$ before every lane start or end. Further, different vehicle types with different behaviors are included, i.e. cars, trucks and motorcycles with different lengths and behaviors. For simplicity, we use the same feature representation for all vehicle classes. As dataset, we collected 500.000 transitions in the same manner as for the *Highway* environment.

### B. Input Features

In the *Highway* scenario, we use the same input features as proposed in [12]. For the *Fast Lanes* scenario, the input features used for vehicle $i$ are:

- *relative distance*: $dr_i = (p_i - p_{\mathrm{agent}})/d_{\mathrm{max}} \in \mathbb{R}$,
  $p_{\mathrm{agent}}$, $p_i$ are longitudinal positions in a curvilinear coordinate system of the lane.
- *relative velocity*: $dv_i = (v_i - v_{\mathrm{agent}})/v_{\mathrm{allowed}}$
- *relative lane index*: $dl_i = l_i - l_{\mathrm{agent}} \in \mathbb{N}$,
  where $l_i$, $l_{\mathrm{agent}}$ are lane indices.
- *vehicle length*: $\mathrm{len}_i/10.0$

The state representation for lane $j$ is:

- *lane start and end*: distances (km) to lane start and end
- *lane valid*: lane currently passable
- *relative lane index*: $dl_j = l_j - l_{\mathrm{agent}} \in \mathbb{N}$,
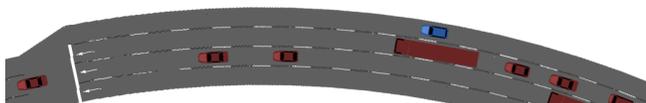  where $l_j$, $l_{\mathrm{agent}}$ are lane indices.



Fig. 2. *Fast Lanes* scenario in SUMO. The agent (blue) is overtaking other vehicles (red) on the fast lane and has to merge before the lane ends.

For the agent, the normalized velocity $v_{\mathrm{current}}/v_{\mathrm{desired}}$ is included, where $v_{\mathrm{current}}$ and $v_{\mathrm{desired}}$ are the current and desired velocity of the agent. Passenger cars, trucks and motorcycles use the same feature representation. When the agent reaches a traffic sign indicating a starting (ending) lane, the lane features get updated until the start (end) of the lane.

### C. Training & Evaluation Setup

All agents are trained off-policy on datasets collected by a rule-based agent with enabled SUMO safety module integrated, performing random lane changes to the left or right whenever possible. For training, traffic scenarios with a random number of $n \in (30,60)$ vehicles for *Highway* and with $n \in (30,90)$ vehicles for *Fast Lanes* are used. Evaluation scenarios vary in the number of vehicles $n \in (30,35,...,90)$. For each fixed $n$, we evaluate 20 scenarios with different *a priori* randomly sampled positions and driver types for each vehicle, to smooth the high variance.

In SUMO, we set the time step length to $0.5\,\mathrm{s}$. The action step length of the reinforcement learning agents is $2\,\mathrm{s}$ and the lane change duration is $2\,s$. Desired time headway $\tau$ and minimum gap are $0.5\,\mathrm{s}$ and $2\,\mathrm{m}$. All vehicles have no desire to keep right (lcKeepRight = 0.0). The sensor range of the agent is $d_{\mathrm{max}} = 80\,\mathrm{m}$. *LC2013* is used as lane-change controller for all other vehicles. To simulate traffic conditions as realistic as possible, different driver types are used with parameters shown in Table I.

### D. Comparative Analysis

Each network is trained with a batch size of $64$ and optimized by Adam [30] with a learning rate of $10^{-4}$. As activation function, we use Rectified Linear Units (ReLu) in all hidden layers of all architectures. The target networks are updated with a step-size of $\tau = 10^{-4}$. All network architectures, including the baselines, were optimized using Random Search with the same budget of 20 training runs. We preferred Random Search over Grid Search, since it has been shown to result in better performance using budgets in this range [31]. The Deep Sets architecture and hyperparameter-optimized settings for all encoder networks are used from [12]. The network architectures are shown in Table II. Graph-Q is compared to two other interaction-aware Q-learning algorithms, that use input modules originally proposed for supervised vehicle trajectory prediction. To support our architecture choices for the Deep Scene-Sets, we compare to a

| Social CNN | VBIN | GCN |
|---|---|---|
| Input($B \times 80 \times 5$) | Input($B \times 15$) | Input($B \times$ seq $\times 3$) |
| $\phi$: FC(20), FC(80) $16 \times$ Conv2D($3 \times 1$) $32 \times$ Conv2D($3 \times 1$) | $\phi$: FC(20), FC(80) concat($\cdot$) $\rho$: FC(80), FC(20) | $\phi$: FC(20), FC(80) $1 \times$ GCN(80) sum($\cdot$) |
| concat($\cdot$, Input($B \times 3$)) FC(100)*, FC(100), Linear(3) | | |

| Deep Scene-Sets | Deep Scene-Graphs |
|---|---|
| Input($B \times$ seq$_0 \times 4$) and Input($B \times$ seq$_1 \times 4$) | |
| $\phi_0$: FC(20), FC(80),FC(80)** $\phi_1$: FC(20), FC(80), FC(80)** sum($\cdot$) $\rho$: FC(80), FC(80) | $\phi_0$: FC(20), FC(80),FC(80)** $\phi_1$: FC(20), FC(80),FC(80)** $1 \times$ GCN(80) sum($\cdot$) |
| concat($\cdot$, Input($B \times 3$)) FC(100), FC(100), Linear(3) | |

TABLE II

NETWORK ARCHITECTURES. FC($\cdot$) ARE FULLY-CONNECTED LAYERS.
THE CNN USES STRIDES OF ($2 \times 1$). (*) FOR VBIN FC(200). (**)
PARAMETERS OF THE LAST LAYERS ARE SHARED.

| Architecture | Parameter | Configuration Space |
|---|---|---|
| Encoders | $\phi$: num layers | $1, 2, 3$ |
| | $\phi$: hidden/ output dims | $5, 20, 80, 100$ |
| Deep Sets | $\rho$: num layers | $1, 2, 3$ |
| | $\rho$: hidden/ output dims | $5, 20, 100$ |
| GCN | num GCN layers | $1,2,3$ |
| | hidden and output dim | $20, 80$ |
| | use edge weights | True, False |
| SocialCNN | CONV: num layers | $2, 3$ |
| | kernel sizes | $([7, 3, 2], [2, 1])$ |
| | strides | $([2, 1], [2, 1])$ |
| | filters | $8, 16, 32$ |
| VBIN | $\phi$ : output dim | $20, 80$ |
| | $\rho$ : hidden dim | $20, 80, 160, 200$ |
| | $Q$ : hidden dim | $100, 200$ |
| Deep Scene-Sets | $\rho$ : output dim | $20, 80$ |
| | shared parameters | True, False |
| Deep Scene-Graphs | use $\rho$ network | True, False |
| | $\rho$ : output dim | $20, 80$ |
| | shared parameters | True, False |

TABLE III

RANDOM SEARCH CONFIGURATION SPACE. FOR EVERY ARCHITECTURE,
WE SAMPLED 20 CONFIGURATIONS TO FIND THE BEST SETTING.

modification with separate $\rho$ networks. We use the following baselines[1]:

*a) Rule-Based Controller:* Naive, rule-based agent controller, that uses the SUMO lane change model *LC2013*.

*b) Convolutional Social Pooling (SocialCNN):* In [25], a social tensor is created by learning latent vectors of all cars by an encoder network and projecting them to a grid map in order to learn spatial dependencies.

*c) Vehicle Behaviour Interaction Networks (VBIN):* In [24], instead of summarizing the output vectors as in the Deep Sets approach, the vectors are concatenated, which results in a limitation to a fixed number of cars. We consider the 6 vehicles surrounding the agent (leader and follower on own, left and right lane).

*d) Multiple $\rho$-networks:* Deep Scene architecture where all object types are processed separately by using $K$ different $\rho$-network modules. The $K$ resulting output vectors are concatenated as $\left[ \rho^1 \left( \sum_{x \in X^{\text{dyn}_1}} \phi^1(x) \right), ..., \rho^K \left( \sum_{x \in X^{\text{dyn}_K}} \phi^K(x) \right) \right]$ and fed into the Q-network module.

*E. Implementation Details & Hyperparameter Optimization*

All networks were trained for $1.25 \cdot 10^6$ optimization steps. The Random Search configuration space is shown in Table III. For all approaches except VBIN, we used the same $\phi$ and $Q$ architectures. Due to stability issues, adapted these parameters for VBIN. For SocialCNN, we used the optimized grid from [12] with a size of $80 \times 5$. The GCN architectures were implemented using the pytorch gemoetric library [32].

## VI. RESULTS

The results for the *Highway* scenario are shown in Figure 3. Graph-Q using the GCN input representation (with

---

[1]Since we do not focus on including temporal context, we adapt recurrent layers to fully-connected layers in all baselines.

all close vehicle connections) is outperforming VBIN and Social CNN. Further, the GCN input module yields a better performance compared to Deep Sets in all scenarios besides in very light traffic with rare interactions between vehicles. While the Social CNN architecture has a high variance, VBIN shows a better and more robust performance and is also outperforming the Deep Sets architecture in high traffic scenarios. This underlines the importance of interaction-aware network modules for autonomous driving, especially in urban scenarios. However, VBIN are still limited to fixed-sized input and additional gains can be achieved by combining both variable input and interaction-aware methods as in Graph Networks. To verify that the shown performance increases are significant, we performed a T-Test exemplarily for 90 car scenarios:

- Independence of the mean performances of DeepSet-Q and Graph-Q is highly significant ($< 0.001$) with a p-value of 0.0011.
- Independence of the mean performances between Graph-Q and VBIN is significant ($< 0.1$) with a p-value of 0.0848. Graph-Q is additionally more flexible and can consider a variable number of surrounding vehicles.

Figure 3 (right) shows the performance of the two graph construction strategies. A graph built with connections for all close vehicles outperforms a graph built with close agent connections only. However, the performance increase is only slight, which indicates that interactions with the direct neighbors of the agent are most important.

The evaluation results for *Fast Lanes* are shown in Figure 4 (left). The vehicles controlled by the rule-based controller rarely use the fast lane. In contrast, our agent learns to drive on the fast lane as much as possible (39.0% of the driving time). We assume, that the Deep Scene-Sets are outperforming Deep Scene-Graphs slightly, because the agent has to
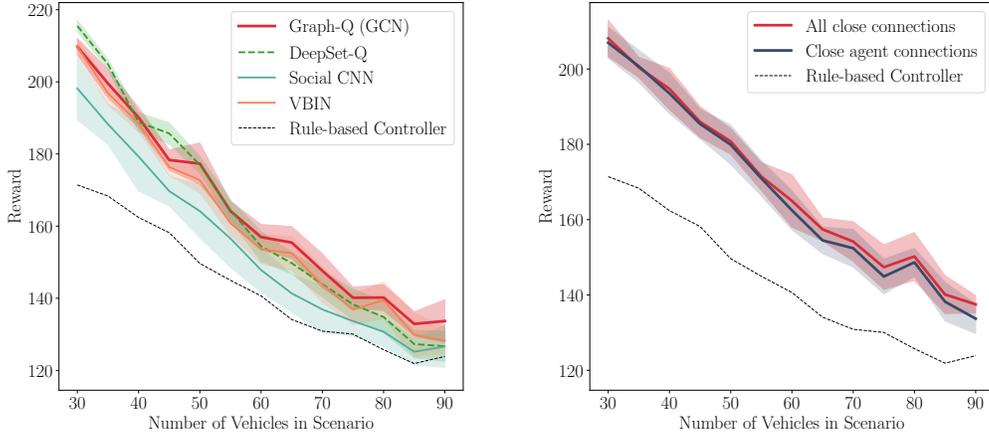
Fig. 3. Mean performance and standard deviation in the *Highway* scenario over 10 training runs for Graph-Q with all close vehicle connections, the Deep Sets [12] and two other interaction-aware Q-function input modules (left), and Graph-Q using the two proposed graph construction strategies (right). The number of vehicles indicates the traffic intensity, from light to dense traffic.
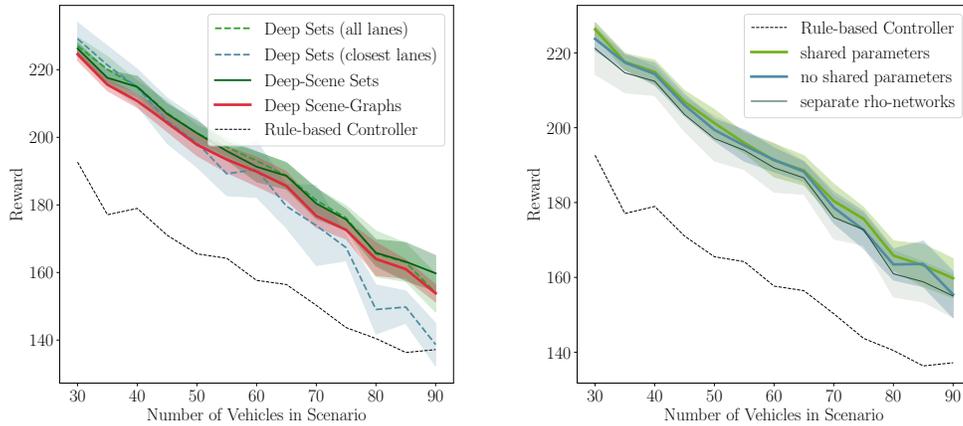


Fig. 4. Mean performance and standard deviation in the *Fast Lanes* scenario over 10 training runs for Deep Scene-Sets, Deep Scene-Graphs and the rule-based controller from SUMO (left), and different architecture choices of the Deep Scenes (right). The number of vehicles indicates the traffic intensity.

deal with less interactions than in the *Highway* scenario. Finally, we compare Deep Scene-Sets to a basic Deep Sets architecture with a fixed feature representation. Using the exact same lane features (if necessary filled with dummy values), both architectures show similar performance. However the performance collapse for the Deep Sets agent considering only its own, left and right lane shows, that the ability to deal with an arbitrary number of lanes (or other object types) can be very important in certain situations. Due to its limited lane representation, the Deep Sets (closest lanes) agent is not able to see the fast lane and thus significantly slower. Figure 4 (right) shows an ablation study, comparing the performance of the Deep-Scene Sets with and without shared parameters in the last layer of the encoder networks. Using shared parameters in the last layer leads to a slight increase in robustness and performance, and outperforms the architecture with separate $\rho$ networks.

## VII. CONCLUSION

In this paper, we propose Graph-Q and DeepScene-Q, interaction-aware reinforcement learning algorithms that can deal with variable input sizes and multiple object types in the problem of high-level decision making for autonomous driving. We showed, that interaction-aware neural networks, and among them especially GCNs, can boost the performance in dense traffic situations. The Deep Scene architecture overcomes the limitation of fixed-sized inputs and can deal with multiple object types by projecting them into the same encoded object space. The ability of dealing with objects of different types is necessary especially in urban environments. In the future, this approach could be extended by devising algorithms that adapt the graph structure of GCNs dynamically to adapt to the current traffic conditions. Based on our results, it would be promising to omit graph edges in light traffic, essentially falling back to the Deep Sets approach, while it is beneficial to model more interactions with increasing traffic density.

REFERENCES

[1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[2] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.

[3] M. Watter, J. T. Springenberg, J. Boedecker, and M. A. Riedmiller, "Embed to control: A locally linear latent dynamics model for control from raw images," in *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, 2015, pp. 2746–2754.

[4] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *Journal of Machine Learning Research*, vol. 17, pp. 39:1–39:40, 2016.

[5] B. Mirchevska, C. Pek, M. Werling, M. Althoff, and J. Boedecker, "High-level decision making for safe and reasonable autonomous lane changing using reinforcement learning," *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2156–2162, 2018.

[6] P. Wolf, K. Kurzer, T. Wingert, F. Kuhnt, and J. M. Zöllner, "Adaptive behavior generation for autonomous driving using deep reinforcement learning with compact semantic states," *CoRR*, vol. abs/1809.03214, 2018. [Online]. Available: http://arxiv.org/abs/1809.03214

[7] B. Mirchevska, M. Blum, L. Louis, J. Boedecker, and M. Werling, "Reinforcement learning for autonomous maneuvering in highway scenarios." *11. Workshop Fahrerassistenzsysteme und automatisiertes Fahren.*

[8] M. Nosrati, E. A. Abolfathi, M. Elmahgiubi, P. Yadmellat, J. Luo, Y. Zhang, H. Yao, H. Zhang, and A. Jamil, "Towards practical hierarchical reinforcement learning for multi-lane autonomous driving," *2018 NIPS MLITS Workshop*, 2018.

[9] M. Kaushik, V. Prasad, M. Krishna, and B. Ravindran, "Overtaking maneuvers in simulated highway driving using deep reinforcement learning," 06 2018, pp. 1885–1890.

[10] M. Mukadam, A. Cosgun, and K. Fujimura, "Tactical decision making for lane changing with deep reinforcement learning," *NIPS Workshop on Machine Learning for Intelligent Transportation Systems*, 2017.

[11] L. Fridman, B. Jenik, and J. Terwilliger, "DeepTraffic: Driving Fast through Dense Traffic with Deep Reinforcement Learning," *arXiv e-prints*, p. arXiv:1801.02805, Jan. 2018.

[12] M. Huegle, G. Kalweit, B. Mirchevska, M. Werling, and J. Boedecker, "Dynamic input for deep reinforcement learning in autonomous driving," *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2019.

[13] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola, "Deep sets," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 3391–3401. [Online]. Available: http://papers.nips.cc/paper/6931-deep-sets.pdf

[14] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *CoRR*, vol. abs/1609.02907, 2016. [Online]. Available: http://arxiv.org/abs/1609.02907

[15] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *Trans. Neur. Netw.*, vol. 20, no. 1, pp. 61–80, Jan. 2009. [Online]. Available: http://dx.doi.org/10.1109/TNN.2008.2005605

[16] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, and M. Sun, "Graph neural networks: A review of methods and applications," *CoRR*, vol. abs/1812.08434, 2018. [Online]. Available: http://arxiv.org/abs/1812.08434

[17] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. A. Riedmiller, R. Hadsell, and P. Battaglia, "Graph networks as learnable physics engines for inference and control," *CoRR*, vol. abs/1806.01242, 2018. [Online]. Available: http://arxiv.org/abs/1806.01242

[18] P. W. Battaglia, R. Pascanu, M. Lai, D. J. Rezende, and K. Kavukcuoglu, "Interaction networks for learning about objects, relations and physics," *CoRR*, vol. abs/1612.00222, 2016. [Online]. Available: http://arxiv.org/abs/1612.00222

[19] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. F. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, Ç. Gülçehre, F. Song, A. J. Ballard, J. Gilmer, G. E. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu, "Relational inductive biases, deep learning, and graph networks," *CoRR*, vol. abs/1806.01261, 2018. [Online]. Available: http://arxiv.org/abs/1806.01261

[20] H. Dai, E. B. Khalil, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," *CoRR*, vol. abs/1704.01665, 2017. [Online]. Available: http://arxiv.org/abs/1704.01665

[21] J. B. Hamrick, K. R. Allen, V. Bapst, T. Zhu, K. R. McKee, J. B. Tenenbaum, and P. W. Battaglia, "Relational inductive bias for physical construction in humans and machines," *CoRR*, vol. abs/1806.01203, 2018. [Online]. Available: http://arxiv.org/abs/1806.01203

[22] J. Jiang, C. Dun, and Z. Lu, "Graph convolutional reinforcement learning for multi-agent cooperation," *CoRR*, vol. abs/1810.09202, 2018. [Online]. Available: http://arxiv.org/abs/1810.09202

[23] F. Diehl, T. Brunner, M. Truong-Le, and A. Knoll, "Graph neural networks for modelling traffic participant interaction," *CoRR*, vol. abs/1903.01254, 2019. [Online]. Available: http://arxiv.org/abs/1903.01254

[24] W. Ding, J. Chen, and S. Shen, "Predicting vehicle behaviors over an extended horizon using behavior interaction network," *CoRR*, vol. abs/1903.00848, 2019. [Online]. Available: http://arxiv.org/abs/1903.00848

[25] N. Deo and M. M. Trivedi, "Convolutional social pooling for vehicle trajectory prediction," *CoRR*, vol. abs/1805.06771, 2018. [Online]. Available: http://arxiv.org/abs/1805.06771

[26] C. J. C. H. Watkins and P. Dayan, "Q-learning," in *Machine Learning*, 1992, pp. 279–292.

[27] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," *CoRR*, vol. abs/1509.06461, 2015. [Online]. Available: http://arxiv.org/abs/1509.06461

[28] 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, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, 2018, pp. 1582–1591. [Online]. Available: http://proceedings.mlr.press/v80/fujimoto18a.html

[29] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker-Walz, "Recent development and applications of sumo - simulation of urban mobility," *International Journal On Advances in Systems and Measurements*, vol. 3&4, 12 2012.

[30] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: http://arxiv.org/abs/1412.6980

[31] J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *J. Mach. Learn. Res.*, vol. 13, pp. 281–305, Feb. 2012. [Online]. Available: http://dl.acm.org/citation.cfm?id=2188385.2188395

[32] M. Fey and J. E. Lenssen, "Fast graph representation learning with pytorch geometric," *CoRR*, vol. abs/1903.02428, 2019. [Online]. Available: http://arxiv.org/abs/1903.02428