

# Formulation and validation of a car-following model based on deep reinforcement learning

Fabian Hart<sup>a</sup>, Ostap Okhrin<sup>a,b</sup>, Martin Treiber<sup>a,b,\*</sup>

<sup>a</sup>*TU Dresden*

<sup>b</sup>*Possible second address*

---

## Abstract

To be written at the end

*Keywords:* reinforcement learning, car-following model, stochastic processes, string stability, validation, trajectory data

---

## 1. Introduction

Autonomous driving technologies are seen as promising solutions to improve road safety, where human errors account for 94% of the total accidents ?. Moreover congestion, energy consumption and emissions are intended to be reduced. But autonomous driving is a challenging task since the transportation traffic can be dynamic and unpredictable. On the way to fully autonomous driving, Advanced Driver Assistance Systems (ADAS) has been developed to solve tasks like lane-keeping, lane-changing, car-following, emergency braking, etc. Since Deep Learning methods has been demonstrated to surpass humans in certain domains, they are also adopted in the area of autonomous driving. Especially Deep Reinforcement Learning (DRL), which combines the power of Deep Learning in tackling large, complicated problems with Reinforcement Learning, has shown its potential in a broad variety of autonomous driving tasks. In ? and ?, DRL is used to guide an autonomous vehicle safely from on-ramp to freeway. In ?, ? and ?, DRL methods are used to manage traffic of autonomous vehicles at

---

\*Corresponding author

Email address: [Martin.treiber@tu-dresden.de](mailto:Martin.treiber@tu-dresden.de) (Martin Treiber)

URL: [www.mtreiber.de](http://www.mtreiber.de) (Martin Treiber)

intersections, optimizing safety and efficiency. In ?, DRL is used to solve lane change maneuvers.

A further task in autonomous driving is to model the vehicle behavior under car following scenarios, where suitable accelerations has to be computed in order to achieve a safe and comfortable driving. Approaches to solve this task are classical car-following models, such as the Intelligent Driver Model Treiber et al. (2000) or stochastic car-following models, such as in Treiber and Kesting (2018). Furthermore data-driven approaches use Machine Learning methods to train a car-following model based on experimental car-follower data, such as in ? or ?. Downside of this approach is that the model tries to emulate human driver behavior which still can be suboptimal.

To overcome this issue, DRL methods train non-human car-following models that can optimize metrics such as safety, efficiency and comfort. One approach is to train on scenarios, where the leading vehicles trajectory, used for training, is based on experimental data, such as in ? or ?. Similar approaches suggest a standardized driving cycle, which functions as leading vehicle trajectory, such as ? or ?, which uses the New European Driving Cycle. A disadvantage coming along with these approaches is, that for scenarios, which are not in the scope of the training data, the performance decreases, indicating inadequate machine learning generalization ?.

Another issue of car-following models is string-stability. There are several studies focusing on dampening traffic oscillations by using a sequence of trained vehicles, such as Qu et al. (2020), ? and ?.

All the mentioned DRL car-following models have three disadvantages in common: At first the acceleration range is limited in a way, that full-brakes are not considered. This results in models that are just designed for non-safety-critical scenarios. Second these models just consider car-following scenarios, while free driving or the transition between both is not reflected in the reward function. Third the trained models have just been validated on data that is similar to the training data set, so that the generalization capabilities cannot be proved.

This motivated us to design a model, which overcomes these issues. To our knowledge, no RL based car-following model has been proposed which has the following features combined:

- The proposed model considers free driving, car-following, as well as the transition between both in a way, that approaching of the leading vehicle is smooth and comfortable.
- The proposed model has a wider range of possible accelerations, which leads to a collision-free behavior also in safety-critical situations such as full-braking of the Leader Vehicle.
- The proposed model is trained on leading trajectories, based on an AR(1)-process, e. g. Honerkamp (1993), the parameters reflecting kinematics of real drivers. This leads to high generalization capabilities and a model usage in a broader variety of traffic scenarios.
- Different driver characteristics can be modeled by adjusting the parameters of the reward function.
- The proposed model shows string-stability.

Another feature of this work is a thorough validation of the above mentioned properties in scenarios based on both synthetic and naturalistic trajectory data, bringing the model to its limits. In all cases, the model proved to be accident free and string stable. In a further experiment the proposed model is compared to an Intelligent Driver Model, calibrated on the same data. The results indicate a better performance of the proposed model.

[short textual enumeration of the sections to come]

## 2. Model specification

The following vehicle [only named things are written in uppercase; special concepts/definitions such as Reinforcement Learning can be written both upper and lowercase; abbreviation always uppercase (RL); the “reward function” is

generic, hence lowercase] is controlled by a Reinforcement Learning (RL) agent. By interaction with an environment, the RL agent optimizes a sequential decision making problem. At each time step  $t$ , [always comma before the main clause starts] the agent observes an environment state  $s_t$  and, based on that state, selects an action  $a_t$ . After conducting action  $a_t$ , the RL agent receives a reward  $r(a_t, s_t)$ . The agent aims to learn an optimal state-action mapping policy  $\pi$  that maximizes the expected accumulated discounted reward  $r_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ , with  $\gamma = (0, 1]$  describing the discount factor. The crucial elements of the RL-based [use your abbreviation] control strategy are described in detail as follows:

### 2.1. Action space

In this study, the acceleration of the following vehicle is considered as the action of the RL agent. To enable comfortable driving and allow for a maximum braking deceleration in safety-critical situations, the acceleration is a continuous variable in the range between  $a_{\min} = -9 \text{ m/s}^2$  and  $a_{\max} = 2 \text{ m/s}^2$ . [In contrast to variables or symbolic constants which are set in italic, units are set in roman (upright) font. If super- or subscripts are variables, they are set italic. If they are abbreviations (e.g., min, max) they are set upright. Multi-letter standard function names such as sin, tan, cosh or exp are written upright (use  $\sin$  etc)]

### 2.2. State space

The state space defines the observations that the following vehicle can receive from the environment. To compute an optimal acceleration, the following vehicle observes its own acceleration  $a$ , its own speed [velocity=vector, speed=abs(velocity)=scalar]  $v$ , the relative speed  $\Delta v = v_l - v$  where  $v_l$  denotes the speed of the leading vehicle [crucial since both definitions  $v_l - v$  or  $v - v_l$  are common] and the (bumper-to-bumper) gap  $s$  to the leader [again, some use the space headway (including the leading vehicle length), and some the gap]. These observations are normalized to the range  $[-1, 1]$ . [How? If it is just linear translation and scaling, you need min and max values of the state variables. How to determine them?]

Table 1: RL agent parameters and default values

Parameter	Description	Value
$a_{\min}$	Minimum acceleration	$-9m/s^2$
$a_{\max}$	Maximum acceleration	$2m/s^2$
$b_{\text{comf}}$	Comfortable deceleration	$-2m/s^2$
$v_{\text{des}}$	Desired velocity	$16.6m/s$
$T_{\text{gap}}$	Desired time gap to Leader	$1.5s$
$s_{\min}$	Desired minimum space gap to Leader	$2m$
$T_{\text{var}}$	Time gap to describe the variance of the normal probability function (see Equation 1 - 4)	$0.7s$
$T_{\text{lim}}$	Upper time gap limit for zero reward (see Equation 1 - 4)	$15s$

### 2.3. Reward Function

The goal of the RL strategy is to reduce the crash risk, while maintaining comfortable driving in non-safety-critical situations and minimizing the travel time [without that, creeping which is safe and comfortable would be the optimum strategy]. The reward function includes a set of parameters that can be adjusted to realize different driving styles.  $v_{\text{des}}$  is the desired velocity that should not be exceeded.  $a_{\min}$  and  $a_{\max}$  are the minimum and maximum possible accelerations, as described in Section 2.1. All parameters are described in Table 1.

The reward function consists of four parts, described as follows: [The standard normal distribution has the symbol  $\Phi(\cdot)$  and its density  $\phi(\cdot)$  or  $\varphi(\cdot)$ ; I prefer  $\varphi$  since, then, the optical difference to the CDF  $\Phi$  is greater]

$$r_1 = \begin{cases} \frac{\varphi(s, s_{\text{opt}}, s_{\text{var}})}{\varphi(s_{\text{opt}}, s_{\text{opt}}, s_{\text{var}})}, & \text{if } s < s^* \\ \frac{\varphi(s^*, s_{\text{opt}}, s_{\text{var}})}{\varphi(s_{\text{opt}}, s_{\text{opt}}, s_{\text{var}})} \left(1 - \frac{s - s^*}{s_{\text{lim}} - s^*}\right) & \text{otherwise} \end{cases} \quad (1)$$

[Use `\left( ... \right)` to make “bigger” parentheses, brackets, braces around

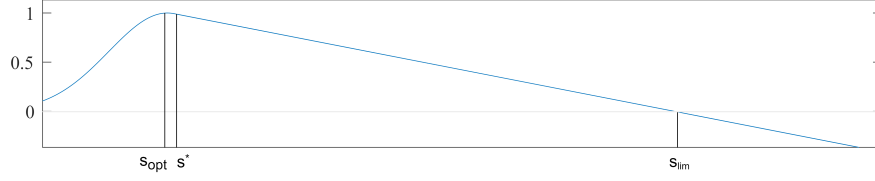


Figure 1: Reward function part 1 maximizes the reward for car following with time gap  $T_{\text{gap}}$

fractions or super- and subscript quantities] with

$$s_{\text{opt}} = vT_{\text{gap}} + s_{\text{min}}, \quad (2)$$

$$s_{\text{var}} = vT_{\text{var}} + 0.5s_{\text{min}}, \quad (3)$$

$$s_{\text{lim}} = vT_{\text{lim}} + 2s_{\text{min}}, \quad (4)$$

and  $\varphi(x, \mu, \sigma^2)$  describing a normal probability density function. [No empty line if no new paragraph; otherwise, you get an indent that looks strange] [ $s^*$  is not defined].

The first part of the reward function aims to maintain a reasonable distance to the Leader Vehicle. Figure 1 illustrates the reward function for  $r_1$ , containing the parameter  $s_{\text{opt}}$ ,  $s^*$  and  $s_{\text{lim}}$ . The reward function is designed in a way, that for high velocities  $v$  of the following vehicle the time gap between Follower and Leader Vehicle tends to  $T_{\text{gap}}$ , while for low velocities the distance between both tends to  $s_{\text{min}}$ . Different values of  $T_{\text{opt}}$  result in different driving styles in a way, that for higher values of  $T_{\text{opt}}$  the Follower drives up closer, resulting in a more aggressive driving style. The results for different values of  $T_{\text{opt}}$  can be found in Section 4.4. Different functions for  $s > s^*$  has been applied, but the best results regarding a smooth and comfortable approaching of the following vehicle has been reached with a linear function. Also, a high value of  $T_{\text{lim}}$  results in an early deceleration and comfortable approaching.

$$r_2 = \begin{cases} \tanh\left(\frac{b_{kin}-b_{comf}}{a_{min}}\right), & \text{if } b_{kin} > b_{comf} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

with

$$b_{kin} = \frac{\Delta v^2}{s} \quad (6)$$

The second part of the reward function addresses the vehicle behavior in safety-critical situations. For a deceleration with the delay  $b_{kin}$  the braking distance is equal to the current distance  $s$ . Thus, the kinematic deceleration  $b_{kin}$  represents the minimum deceleration necessary to avoid a collision. A situation is considered as "safety-critical", if the kinematic deceleration  $b_{kin}$  is greater than the comfortable deceleration  $b_{comf}$ . Thus, just in safety-critical situations the RL agent is getting penalized, illustrated in figure xy.

$$r_3 = -\left(\frac{da}{dt}\right)^2 \quad (7)$$

The third part of the reward function aims to reduce the jerk in order to achieve comfortable driving.

$$r_4 = \begin{cases} -\min\left(1, (v - v_{des})^2\right), & \text{if } v > v_{des} \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

The fourth part of the reward function penalizes the RL agent, if the current velocity  $v$  is above the desired velocity  $v_{des}$ .

$$r = 0.6r_1 + 1.1r_2 + 0.001r_3 + r_4 \quad (9)$$

The weights of each reward part has been found experimentally and can further be optimized in future studies.

#### 2.4. RL algorithm

In various similar control problems, the Deep Deterministic Policy Gradient (DDPG) Algorithm has been used and proven to perform well on tasks with

Table 2: DDPG parameter values

Parameter	Value
Learning rate	0.001
Reward discount factor	0.95
Experience buffer length	100000
Mini batch size	32
Gaussian noise mean	0.15
Gaussian noise variance	0.2
Gaussian noise variance decay	1e-5
Number of hidden layers	2
Neurons per hidden layer	32

a continuous action and state space, such as in ?, ? or ?. The original work can be found in ?. In order to reduce the variance of policy gradients and increase learning speed, DDPG is an actor-critic method. The actor determines the action, while the critic judges about the quality of the action and how the policy should be adjusted. In this study, both networks are feed-forward neural networks with two layers of hidden neurons and 32 neurons each hidden layer. All DDPG parameters are presented in Table 2.

### 2.5. Reward Function

learning input (leader speed time series)

[also relate parameters to driving style attributes such as desired speed, accelerations, decelerations, desired time gap, minimum gap]

## 3. Model training

The training of the model comprises two important components, which have to be defined in advance. There is the generation of leading trajectories and the general definition of an training episode, that will be discussed in the following.



### 3.1. Generating synthetic leading trajectories

The leading trajectory is based on an AR(1) process, whose parameters reflect the kinematics of real leaders. The AR(1) process describes the speed of the Leader Vehicle and is defined as

$$v_l(t) = c + \phi v_l(t-1) + \epsilon, \text{ with } E(\epsilon) = 0, Var(\epsilon) = \sigma \quad (10)$$

With reaching of stationarity, this process has

$$\text{an expected value of } E(v_l) = \frac{c}{1-\phi}, \quad (11)$$

$$\text{the variance } Var(v_l) = \frac{\sigma^2}{1-\phi^2}, \quad (12)$$

$$\text{the autocorrelation } ACF(dt) = \phi^{dt}, \quad (13)$$

$$\text{and the correlation time } \tau = -\frac{1}{\ln(\phi)}, \quad (14)$$

with  $d$  corresponding to the simulation step size, which is globally set to  $100ms$ .

To adjust the parameters of the AR(1) process, typical values for real leader trajectories has to be defined: With  $v_{l,des}$  as the desired velocity for the leader, the mean of the AR(1) process is set to be  $v_{l,des}/2$  and the standard deviation is set to be  $v_{l,des}$ . The acceleration  $a_{phys}$  corresponds to typical physical leader accelerations. With these values and by using Equation 11 - 14, the parameters of the AR(1) process can be calculated as:

$$\phi = e^{\left(\frac{-2da_{phys}}{v_{l,des}}\right)} \quad (15)$$

$$c = (1-\phi) \frac{v_{l,des}}{2} \quad (16)$$

$$\sigma^2 = (1-\phi^2) \frac{v_{l,des}^2}{4} \quad (17)$$

Table 3: Assumed typical values for leading trajectories and the resulting values of the AR(1) process parameters

Real trajectory		AR(1) process	
$v_{l,des}$	$15m/s$	$\phi$	0.9933
$a_{phys}$	$1m/s^2$	$c$	$0.05m/s$
		$\sigma^2$	$0.75m^2/s^2$

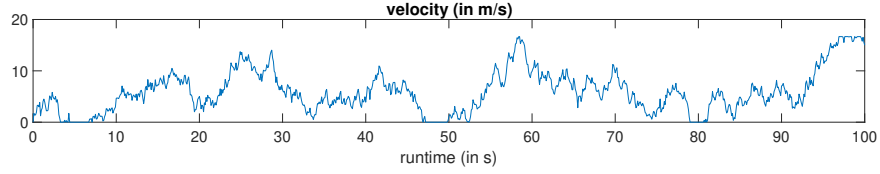


Figure 2: Example of a leading trajectory based on the parametrized AR1 process used to train the RL agent

The assumed typical values for  $v_{l,des}$  and  $a_{phys}$  as well as the resulting values of the AR(1) process parameters can be found in Table 3.

Figure 2 shows an example trajectory of the leading vehicle based on the AR(1) process using the parameters of Table 3. If the velocity exceeds the defined range of  $[0, v_{l,des}]$ , it is manually set to the range limits.

### 3.2. Training episode definition

To train the RL agent, a training episode has to be defined. One episode has a simulation time of 50s with a step size of  $d = 100ms$ , resulting in a episode length of 500 steps. The initial velocities of Follower and Leader Vehicle is a randomly set in the range  $[0, v_{des}]$  respectively  $[0, v_{l,des}]$ . The initial gap between both is set to 120m.

### 3.3. Results of the RL training process

Figure 3 shows an example of the training process. The red line shows the moving average reward of the last 30 episodes. After approximately 570

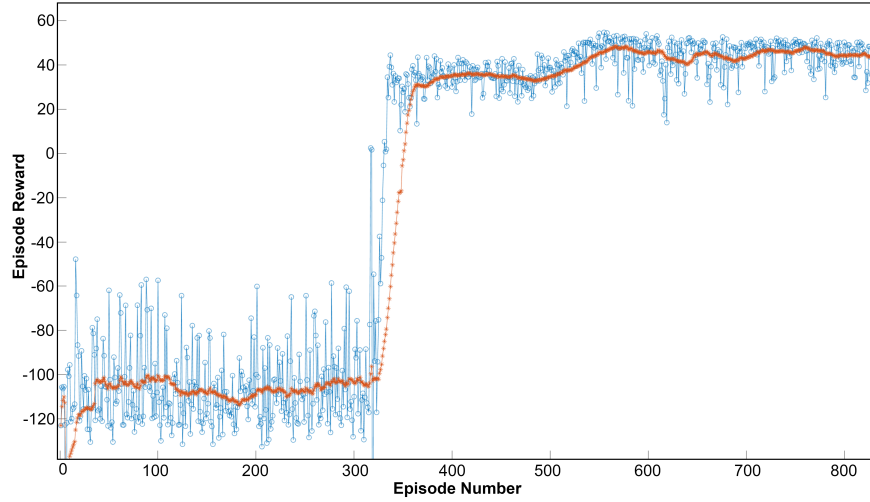


Figure 3: RL training process, one episode containing 500 steps

episodes the maximum average reward has been reached. Reaching saturation the learning process has been stopped after 850 episodes.

#### 4. Validation

The goal is not to minimize some error measure as in usual calibration/validation but to check if the driving style is safe, effective, and comfortable. The RL strategy is evaluated with respect to these metrics in different driving scenarios, described in the following.

##### 4.1. Response to an external leading vehicle speed profile

The first scenario is designed in order to evaluate the transition between free driving and car-following as well as the Follower behavior in safety-critical situations. Figure 4 shows a driving scenario with an external Leading Vehicle speed profile. The initial gap between Follower and Leader is 200 meters, which refers to a free driving scenario. The Follower accelerates with  $a_{\max} = 2m/s^2$  until the desired speed  $v_{\text{des}} = 16.6m/s$  is reached and approaches the standing Leader Vehicle. When the gap between both drops below 90 meters, the Follower starts to decelerate with a approximately  $b_{\text{comf}} = -2m/s^2$  (transition between

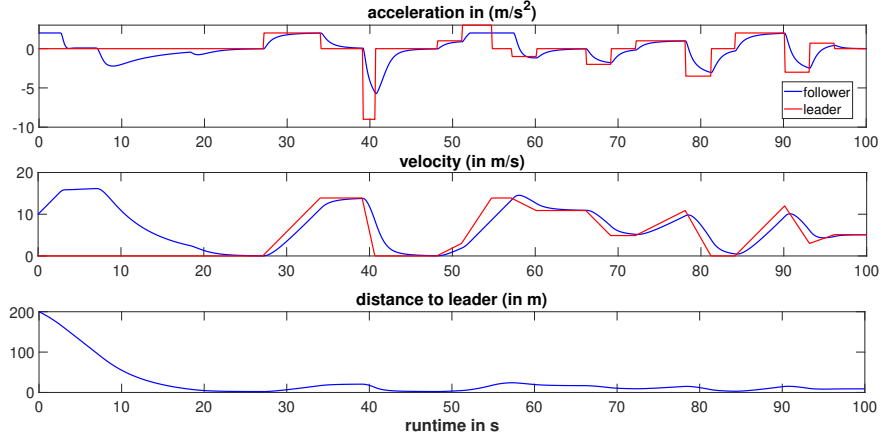


Figure 4: Response to an external leading vehicle speed profile

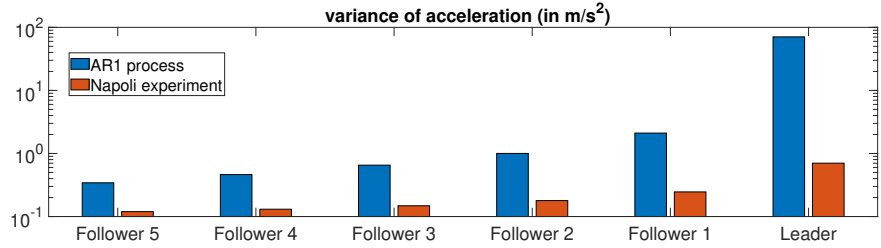


Figure 5: Comparison of the acceleration variance between Leader and Follower for AR(1) and Napoli experiment

free driving and car-following) and comes to a standstill with a final gap of approximately  $s_{\min} = 2m$ . In the following the Leader Vehicle makes some random acceleration and deceleration. At the time  $t = 40s$  the Leading Vehicle makes a full brake, resulting in a comfortable and safe braking of the following vehicle. The transition between different accelerations happens in a comfortable way, reducing the resulting jerk.

#### 4.2. String stability

The second scenario, shown in Figure 6, consists of a Leader based on the AR(1) process, followed by five vehicles, each controlled by the trained RL agent. The results show that traffic oscillations can effectively be dampened

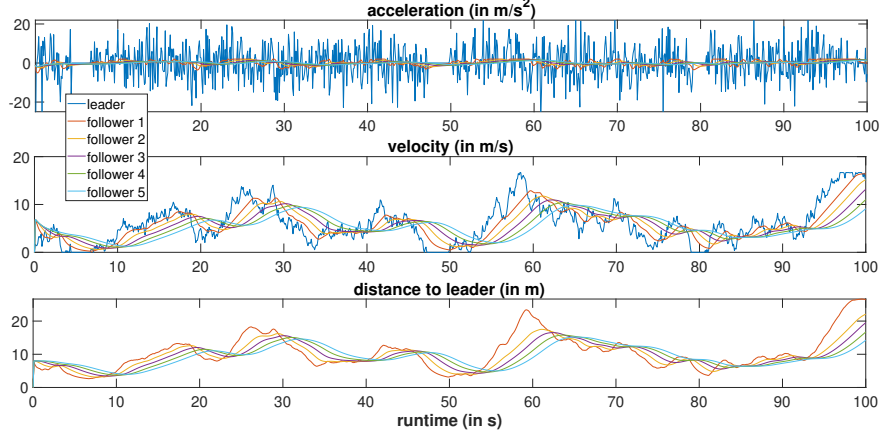


Figure 6: Response to a leader trajectory based on a AR(1) process

with a sequence of trained Followers, even if the Leader shows large outliers in acceleration. Figure 5 illustrates the difference of accelerations between Leader and the Followers (blue bars). The last Follower shows the lowest variance of acceleration, which demonstrates the ability of the RL agent to flatten the speed profile, to dampen oscillations and thus to increase comfort.

#### 4.3. Response to a real leader trajectory

In a further scenario, the abilities of the RL strategy are evaluated with a real leader trajectory. This trajectory comes from experiments in Napoli, where exact data from Leader and Follower were obtained (reference to Punzo et al.). Figure 7 shows the result of a sequence of five vehicles following a real leader trajectory. Similar to the experiment from Section 4.2 string stability and the reduction of acceleration variance, shown by the red bars in Figure 5, is demonstrated. At time  $t = 140\text{s}$  the Leader stands still, and it can be observed, that all following vehicles are keeping the minimum distance  $s_{\min}$  to the Leader.

#### 4.4. Response of different driver characteristics

As mentioned in Section 2.3, different driving styles can be achieved by adjusting the parameters of the reward function. Three RL agents has been

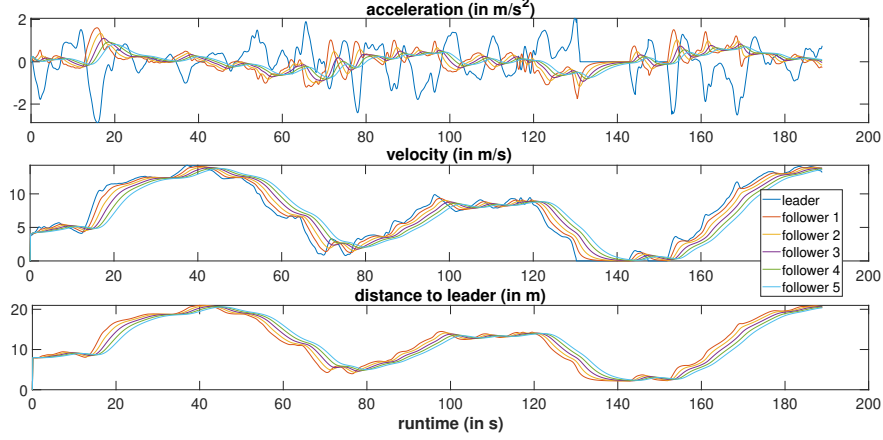


Figure 7: Response to a real leader trajectory

trained on a reward function, that differs in the desired time gap  $T_{\text{gap}}$  between Follower and Leader Vehicle ( $T_{\text{gap},1} = 1.0s$ ,  $T_{\text{gap},2} = 1.5s$ ,  $T_{\text{gap},3} = 2.0s$ ). Figure 8 shows the result of these agents, following the real leader trajectory from Napoli. It can be observed, that a lower value for  $T_{\text{gap}}$  results in closer driving to the Leader, higher accelerations and decelerations and thus in a more aggressive driving behavior.

#### 4.5. Cross validation with the Intelligent Driver Model

To compare the performance of the trained RL agent with a classical car-following model, the RL agent and an IDM are calibrated on the same set of parameters. These parameters are obtained from one and the same follower from the Napoli data set (reference) and can be found in Table 4. The IDM is formulated as

$$a = a_{\max} \left( 1 - \left( \frac{v}{v_{\text{des}}} \right)^4 - \left( \frac{s^*(v, \Delta v)}{s} \right)^2 \right), \quad (18)$$

with

$$s^*(v, \Delta v) = s_{\min} + vT_{\text{gap}} + \frac{v\Delta v}{2\sqrt{a_{\max}b_{\text{comf}}}}, \quad (19)$$

$$\Delta v = v - v_l \quad (20)$$

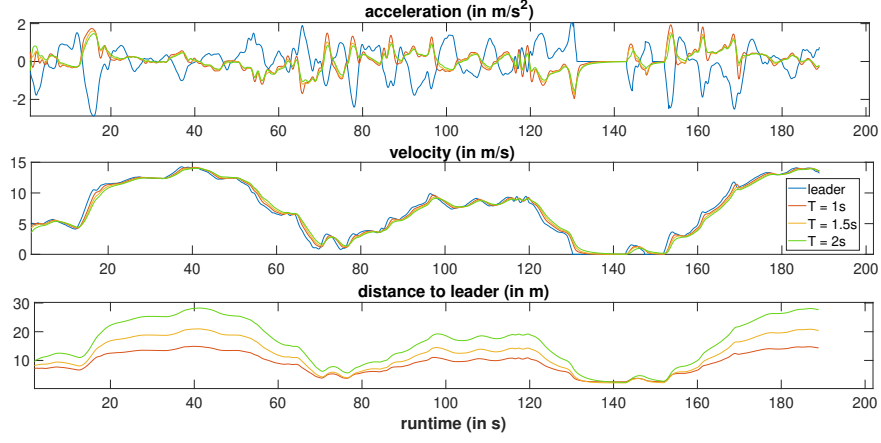


Figure 8: Impact of different parametrized RL agents on driving behavior

Table 4: Parameters obtained from experimental data and used to calibrate the RL agent and the IDM

Parameter	Value
$T_{\text{gap}}$	$0.83s$
$s_{\text{min}}$	$4.90m$
$a_{\text{max}}$	$4.32m/s^2$
$b_{\text{comf}}$	$2.34m/s^2$
$v_{\text{des}}$	$33.73m/s$

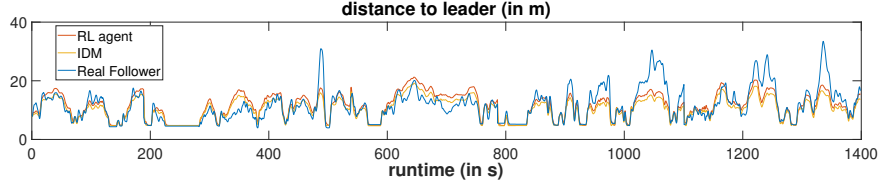


Figure 9: Comparison between IDM and RL agent, calibrated on the same set of parameters

Table 5: Comparison between calibrated RL agent and IDM for accumulated Reward and Goodness-of-Fit Function  $SSE(log(s))$

	RL agent	IDM
$SSE(log(s))$	389.10	418.05
Accumulated Reward	$8.23 \times 10^3$	$8.08 \times 10^3$

Figure 9 shows the results for: First, the RL agent, calibrated on the real follower data. Second, the IDM, calibrated on the follower data. And third, the real follower of the Napoli experiment. To evaluate the performance, for both approaches the respective objective function has been computed. The objective function of the RL agent correspondences to the reward function, while the Goodness-of-Fit Function  $SSE(s)$  defines the objective function of the IDM. A comparison between RL agent and IDM for both, the reward function and the Goodness-of-Fit Function, is shown in Table 5. For both objectives the RL agent shows a better performance.

## 5. Conclusion/Discussion

evaluation of safety and comfort, comparison to IDM

N. P. Farazi, T. Ahamed, L. Barua, B. Zou, Deep reinforcement learning and transportation research: A comprehensive review, arXiv preprint arXiv:2010.06187 (2020).



- X. Qu, Y. Yu, M. Zhou, C.-T. Lin, X. Wang, Jointly dampening traffic oscillations and improving energy consumption with electric, connected and automated vehicles: a reinforcement learning based approach, *Applied Energy* 257 (2020) 114030.
- M. Treiber, A. Hennecke, D. Helbing, Congested traffic states in empirical observations and microscopic simulations, *Physical Review E* 62 (2000) 1805–1824.
- M. Treiber, A. Kesting, *Traffic Flow Dynamics: Data, Models and Simulation*, Springer, Berlin, 2013.
- M. Treiber, A. Kesting, The intelligent driver model with stochasticity – new insights into traffic flow oscillations, *Transportation Research Part B: Methodological* 117 (2018) 613 – 623. TRB:ISTTT-22.
- J. Honerkamp, *Stochastic dynamical systems: concepts, numerical methods, data analysis*, John Wiley & Sons, 1993.