

Tackling Real-World Autonomous Driving using Deep Reinforcement Learning

Paolo Maramotti¹, Alessandro Paolo Capasso², Giulio Bacchiani² and Alberto Broggi²

Abstract—In the typical autonomous driving stack, planning and control systems represent two of the most crucial components in which data retrieved by sensors and processed by perception algorithms are used to implement a safe and comfortable self-driving behavior. In particular, the planning module predicts the path the autonomous car should follow taking the correct high-level maneuver, while control systems perform a sequence of low-level actions, controlling steering angle, throttle and brake. In this work, we propose a model-free Deep Reinforcement Learning Planner training a neural network that predicts both acceleration and steering angle, thus obtaining a single module able to drive the vehicle using the data processed by localization and perception algorithms on board of the self-driving car. In particular, the system that was fully trained in simulation is able to drive smoothly and safely in obstacle-free environments both in simulation and in a real-world urban area of the city of Parma, proving that the system features good generalization capabilities also driving in those parts outside the training scenarios. Moreover, in order to deploy the system on board of the real self-driving car and to reduce the gap between simulated and real-world performances, we also develop a module represented by a tiny neural network able to reproduce the real vehicle dynamic behavior during the training in simulation.

I. INTRODUCTION

In recent years, the technological progress in mechanical, electronic and IT areas has led industrial and academic sectors to invest resources in the autonomous driving field. In the last decades, excellent progress has been conducted for increasing the level of vehicle automation, starting from simple, rule-based approaches [1] to implementing smart systems based on Artificial Intelligence ([2], [3]). In particular, these systems aim at solving the main limitations of the rule-based approaches that are the lack of negotiation and interaction with other road users and the poor understanding of the dynamics of the scenario.

In this paper, we focus on the development of a Deep Reinforcement Learning (DRL) [4] planner, able to drive the vehicle safely and comfortably in obstacle-free environments, using a neural network that predicts both the acceleration and the steering angle.

Indeed, Reinforcement Learning (RL) [5] is widely used to solve tasks using both discrete control space output, such as Go [6], Atari games [7] or chess [8], and continuous control space [9], as autonomous driving [10].

¹Paolo Maramotti is with Vislab - University of Parma, Parma, Italy paolo.maramotti@unipr.it

²Alessandro Paolo Capasso, Giulio Bacchiani, Alberto Broggi are with Vislab srl, an Ambaraella Inc. company - Parma, Italy acapasso@ambarella.com, gbacchiani@ambarella.com, broggi@vislab.it

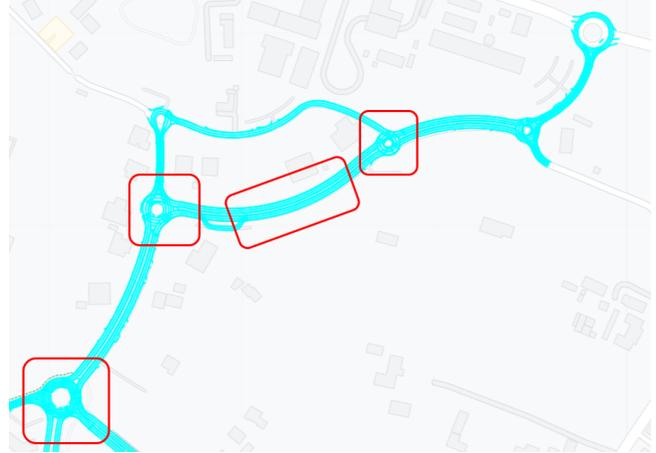


Fig. 1: Mapped area of a neighborhood of Parma. The four red rectangles contain the scenarios used for training agents. The whole light blue area was used for testing.

In particular, RL algorithms are widely used in the autonomous driving field for the development of decision-making and maneuver execution systems like lane change ([11], [12], [13]), lane keeping ([14], [15]), overtaking maneuvers [16], intersection and roundabout handling ([17], [18]) and many others.

Starting from the delayed version of Asynchronous Advantage Actor Critic (A3C) algorithm ([19], [20], [21]), we implemented a Reinforcement Learning planner training agents in a simulator based on High Definition Maps (HD Maps [22]) developed internally by the research team. In particular, we trained the model predicting continuous actions related to the acceleration and steering angle, and testing it on board of a real self-driving car on an entire urban area of the city of Parma (Fig. 1).

Moreover, we demonstrate that the system achieved good generalization performances both in simulation and in real world tests, showing that the module is able to drive smoothly and safely also in those areas not included in the training scenarios. Furthermore, in order to obtain agents that behave as similar as possible to the real self-driving car, a simple model based on a neural network (called *deep_response*) has been developed to reproduce in simulation the real self-driving vehicle dynamic behavior. In this way, we reduce the gap between the simulated agent and the real vehicle behavior, proving that the *deep_response* module is essential to obtain a comfort and safe driving style in real-world scenarios.

Finally, since RL algorithms take millions of episodes before converging to an optimal solution, especially in complex tasks as in the autonomous driving field, we decided to perform a pre-training using Imitation Learning (IL) [23] that allow us to overcome this limit. In particular, we collected a dataset using agents that drive in the scenarios following deterministic rules.

Therefore, we exploited the ability of Imitation Learning to converge more quickly to an optimal solution to obtain network parameters as initial weights for the RL training. In this way we were able to drastically reduce the time needed to obtain an optimal policy.

II. RELATED WORKS

The development of the autonomous driving systems combined with the progress of neural networks, has increasingly turned towards learn-based systems.

In [24] the authors proposed a framework for an end-to-end Deep Reinforcement Learning [4] pipeline for autonomous driving which integrates RNNs (Recurrent Neural Network) [25] to account for Partially Observable MDP (Markov Decision Process); they also integrate an attention models into the framework scenarios and test it in TORCS simulator [26]. Since 2015, when DeepMind published its work on the Deep Q-Network [7] in which, through Reinforcement Learning, they solve many Atari games, some researchers have begun to test such approach in other areas as well. In the autonomous driving field, in [27] the authors presented a RL approach that uses Deep Q-Networks to drive a vehicle in a 3D physics simulator in an end-to-end manner. Similarly, based on Deep Q-Networks and merging it with Deterministic Policy Gradient (DPG), in [28] was proposed Deep Deterministic Policy Gradient (DDPG). This is a method which can be used to solve the control problem in the high-dimensional and continuous action space and was also tested in the autonomous driving field in TORCS simulator [26].

The first work that showed that DRL is a viable approach to autonomous driving in real world is [29]. They developed a system able to perform the lane following maneuver using a single monocular camera. The system differs from the others in the literature because both training and testing are done directly on the vehicle and not only in simulation.

In this paper we used a delayed version of the original A3C [19] called Delayed-A3C (D-A3C). This algorithm was previously developed and used in [20] and [21] where it has been shown that it allows to achieve better results than A3C. In D-A3C configuration, each agent begins the episode with a local copy of the latest version of the global network, while the system collects all the contribution of the actors; the agent updates their local copy of the network at fixed time intervals but all the updates are sent to the global network only at the end of the episode, while in the classical A3C algorithm this exchange is performed at fixed time intervals.

Another approach widely used for the development of planning and control systems is Imitation Learning, which is a supervised learning technique that allows to learn a policy by demonstration, or rather by mimicking the behavior

collected in the dataset. For this purpose, NVIDIA ([30], [31]) developed an end-to-end driving system using deep convolutional neural networks. They trained the system to map images captured from a single front camera directly into steering commands. This was followed by many other works such as [32] in which they do not use as input for the system the raw data taken directly from the sensors, but process them to obtain bird’s-eye-views that contain a more concise and generalized representation of the data. In this way, unnecessary information can be discarded and only useful information for decision making and planning are kept.

Equally interesting are the works [33] and [34], where CARLA simulator was created and on which they tested their conditional imitation learning algorithms, developed to train an agent to drive following the road and avoiding obstacles, based on the behavior of a human.

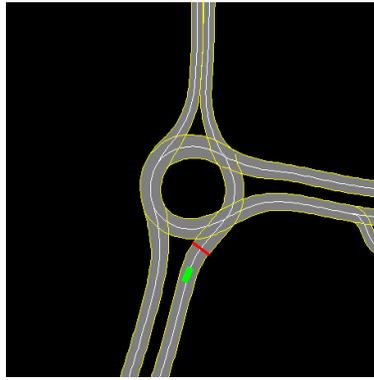
The main advantage of IL respect to RL is that is not necessary to manually design the policy model or cost function. However, IL has the main limitation of requiring a large amount of data in order to avoid overfitting and at the same time it is difficult to create datasets containing all the possible situations a vehicle could encounter.

III. NOTATION

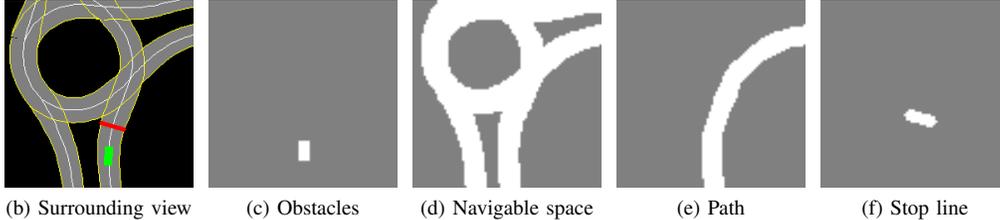
A typical Reinforcement Learning algorithm involves the interaction between at least one agent and its environment. Through the trials and errors mechanism, it learns to behave in order to achieve its goal by receiving a reward function to evaluate the goodness of the chosen action. In general, at time t the agent performs an action a_t observing a state s_t ; as a consequence of such action it will receive a reward signal r_t finding itself in a new state $s_{(t+1)}$.

The RL problem can be defined as a Markov Decision Process (MDP), $M = (S, A, P, r, \gamma)$, where S is a set of states, A is a set of discrete or continuous actions, P is the state transition probability $P(s_{t+1}|s_t, a_t)$, r represents the reward function and γ is the discount factor ($[0, 1]$) that modulates the importance of future rewards. Therefore, the goal of a Reinforcement Learning agent is to maximize the expected return $R_t = \sum_t^T r_t + \gamma r_{t+1} + \dots + \gamma^{T-t} r_T$ (where T represent the time instant of the terminal state).

In this paper we use the delayed version of Asynchronous Advantage Actor Critic (D-A3C) that belongs to the family of the so called Actor-Critic algorithms [19]. In particular, it is composed by two different entities: the *actor* and the *critic*. The purpose of the *actor* is to choose the action that the agent has to perform, while the *critic* estimates the state-value function, namely the goodness of a specific state for the agent. In other words, the *actor* is a probability distribution over actions $\pi(a|s; \theta^\pi)$ (where θ are the network parameters) and the *critic* is estimates the state-value function $v(s; \theta^v) = \mathbb{E}(R_t|s_t)$.



(a) Top-view of the training scenario



(b) Surrounding view

(c) Obstacles

(d) Navigable space

(e) Path

(f) Stop line

Fig. 2: Synthetic representation of training scenario (Fig. 2a) included in the mapped area of Fig. 1. Fig. 2b represents the 50×50 meters surrounding perceived by the agent, that is split in four channels: obstacles (Fig. 2c), navigable space (Fig. 2d), path (Fig. 2e) and stop line (Fig. 2f).

IV. TRAINING SETTINGS

A. Environment

For the development of such work, a synthetic simulator has been implemented using High Definition Maps both developed internally by the team; an example of a scenario is illustrated in Fig. 2a. This represents a portion of the mapped area (Fig. 1) in which we tested our system on board of the real self-driving car. The Fig. 2a shows a full view of such training scenario while the Fig. 2b the surrounding view perceived by the agent, corresponding to an area of 50×50 meters. This is divided into four channels: the obstacles (Fig. 2c) that at this stage of the work are only represented by the ego vehicle, the navigable space (Fig. 2d), the path (Fig. 2e) that the agent should follow and that is randomly calculated at the beginning of the episode and, finally, the stop line (Fig. 2f). The use of HD Maps in the simulator allows us to retrieve several information about the external environment like the positions or the number of lanes, the road speed limits and much more.

In this phase of the work, we focused on the achievement of a smooth and safely driving style, and for this reason we trained our agents in static scenarios, without including obstacles or other road users, learning to following the route and respecting the speed limits.

B. Neural Network

Agents are trained using the neural network illustrated in Fig. 3, that predicts steering angle and acceleration set points every 100 ms. It is divided into two main sub-modules: the first one able to define the steering angle sa and the second

one for the acceleration acc . The input of both sub-modules is represented by 4 channels (navigable space, path, obstacles and stop line) corresponding to the surrounding view of the agent (Fig. 2). Each visual input channel contains 4 images of 84×84 pixels in order to give the agent an history of the past states. Together with this visual input, the network receives 5 scalar parameters including the target speed (the road speed limit), the current speed of the agent, the ratio between the current speed and the target speed, and the last actions related to the steering angle and the acceleration.

In order to guarantee exploration, the output of the two sub-modules is sampled from two Gaussian distributions to obtain relatively acceleration ($acc = \mathcal{N}(\mu_{acc}, \sigma_{acc})$) and steering angle ($sa = \mathcal{N}(\mu_{sa}, \sigma_{sa})$). The standard deviations σ_{acc} and σ_{sa} are predicted and modulated by the neural network during the training phase such that they give an estimate of the uncertainty of the model. Furthermore, the network produces the corresponding state-value estimations (v_{acc} and v_{sa}) using two different reward functions $R_{acc,t}$ and $R_{sa,t}$ related to the acceleration and steering angle respectively.

The neural network was trained on the four scenarios contained in the red rectangles of Fig. 1. For each scenario, we created multiple instances on which the agents act independently of each other. Every agent follows the kinematic bicycle model [35] using values $[-0.2, +0.2]$ for the steering angle and $[-2.0 \frac{m}{s^2}, +2.0 \frac{m}{s^2}]$ for the acceleration. At the beginning of the episode, each agent starts driving with a random speed $([0.0, 8.0])$ and following its predetermined path and respecting the road speed limits. The road speed limit on this urban area can vary between $4 \frac{m}{s}$ to $8.3 \frac{m}{s}$.

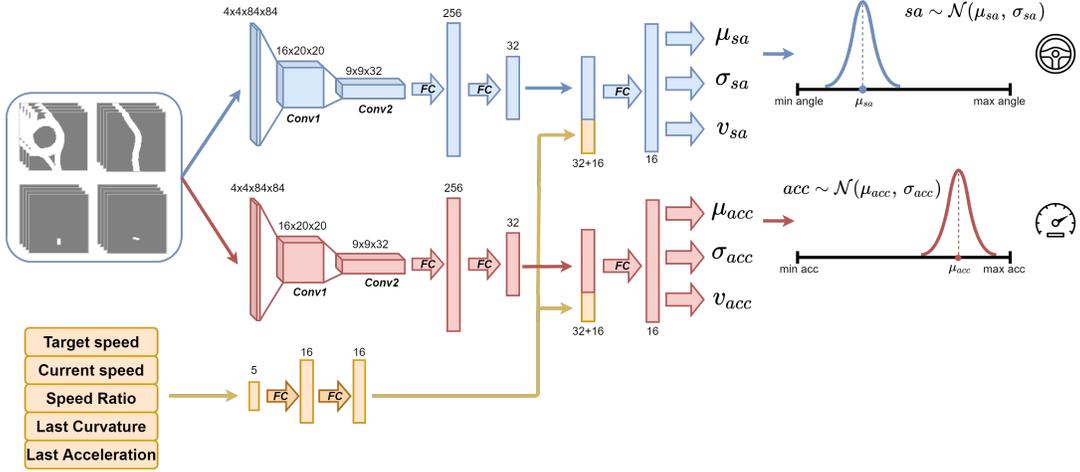


Fig. 3: The neural network architecture used for training agents to drive following their paths and observing the road speed limits. The net is composed by two sub-modules that predict acceleration and steering angle, sampled from two Gaussian distribution ($acc = \mathcal{N}(\mu_{acc}, \sigma_{acc})$, $sa = \mathcal{N}(\mu_{sa}, \sigma_{sa})$). Both the sub-modules receive the same five scalar parameters (Target speed, Current speed, Speed ratio, Last curvature, Last acceleration) and the same four visual input, each one composed by the last four sequential frame images 84×84 pixels in order to give the agent a history of the past states and, consequently understand the dynamic of the environment.

Finally, since there are no obstacles inside the training scenarios, episodes can finish in one of the following terminal states:

- *Goal reached*: the final goal position was reached by the agent.
- *Off-road*: the agent goes outside its predetermined path, predicting erroneously steering angle values.
- *Time-over*: the time to finish the episode expires; this is mainly due to cautious predictions of the acceleration output, driving the vehicle at lower speeds than the road speed limits.

C. Reward

In order to obtain a policy able to drive the car smoothly both in simulation and in the real world environment, a reward shaping is essential to achieve the desired behavior. In particular, we define two different reward functions to evaluate the two actions separately: $R_{acc,t}$ and $R_{sa,t}$ related to the acceleration and the steering angle respectively. They can be defined as follows:

$$R_{acc,t} = r_{speed} + r_{acc.indecision} + r_{terminal} \quad (1)$$

$$R_{sa,t} = r_{localization} + r_{sa.indecision} + r_{terminal} \quad (2)$$

r_{speed} is a signal given to $R_{acc,t}$ depending on the value of the ratio (sr) between the vehicle current speed and the target speed, and it encourages the agent to achieve but not to exceed the route speed limit, and it can be defined as:

$$r_{speed} = \begin{cases} sr \cdot \zeta & \text{if } sr < 1.0, \\ (sr - 1.0) \cdot \zeta & \text{otherwise,} \end{cases} \quad (3)$$

where ζ is a constant set to 0.009.

$r_{localization}$ is a penalization given to $R_{sa,t}$ when the position

or the heading of the agent differ from those of the road, and it is defined as:

$$r_{localization} = \phi \cdot (h_a - h_p) + \chi \cdot d \quad (4)$$

where ϕ and χ are constants set to 0.05, h_a and h_p are the heading of the agent and that of the road respectively, and finally d is the lateral distance between the position of the agent and the center of the lane. Both $R_{sa,t}$ and $R_{acc,t}$ have an element in the formula with the purpose of penalizing two consecutive actions that differ by a value greater than a certain threshold, δ_{acc} and δ_{sa} for acceleration and steering angle respectively. In particular, the difference between two consecutive accelerations is calculated as $\Delta_{acc} = |acc(t) - acc(t-1)|$ and $r_{acc.indecision}$ is defined as:

$$r_{acc.indecision} = \psi \cdot \min(0.0, \delta_{acc} - \Delta_{acc}) \quad (5)$$

where ψ is a constant set to 0.1 and δ_{acc} is set to $0.5 \frac{m}{s^2}$. Instead, the difference between two consecutive predictions of the steering angle is calculated as $\Delta_{sa} = |sa(t) - sa(t-1)|$, such that $r_{sa.indecision}$ is defined as:

$$r_{sa.indecision} = \lambda \cdot \min(0.0, \delta_{sa} - \Delta_{sa}) \quad (6)$$

where λ is a constant set to 0.01 and δ_{sa} is set to 0.05.

Finally, $R_{acc,t}$ and $R_{sa,t}$ depend on the terminal state achieved by the agent:

- *Goal reached*: the agents achieves the goal position, so $r_{terminal}$ is set to +1.0 for both rewards.
- *Off-road*: the agent goes off of its path and it is mainly due to an inaccurate prediction of the steering angle. For this reason we assign a negative signal of -1.0 to $R_{sa,t}$ and 0.0 to $R_{acc,t}$.
- *Time-over*: the available time to finish the episode expires, and this is mainly due to an overly cautious

acceleration predictions of the agent; for this reason $r_{terminal}$ assumes the value of -1.0 for $R_{acc,t}$ and 0.0 for $R_{sa,t}$.

D. Deep Model

One of the main problems related to the use of simulators consists in the difference between simulated and real data, caused by the difficulty of faithfully reproducing real-world situations inside the simulator. To overcome this problem we used a synthetic simulator (Fig. 2) in order to simplify the input of the neural network and to reduce the gap between simulated and real world data. Indeed, the information contained in the 4 channels (obstacles, navigable space, path and stop line) passed as input to the neural network, can be easily reproduced by perception and localization algorithms and by HD Maps embedded on the real self-driving car.

Moreover, another relevant problem in the use of a simulator is related to the difference in how simulated agents perform a target action compared to how the autonomous car would behave executing that command. Indeed, a target action computed at time t can ideally be executed with immediate effect at the same precise instant of time in simulation. Differently, this not happens on board of a real vehicle in which such target action will be performed with a certain dynamics, that results in a delay in the execution ($t + \delta$). For this reason, it is necessary to introduce such response time in simulation, in order to train agents to handle such delay on board of the real autonomous car.

At this purpose, to achieve a more realistic behavior, we initially trained agents adding a low-pass filter to the target action predicted by the neural network that the agent must perform. In Fig. 4, the blue curve represents the ideal and instantaneous response time that occurs in simulation applying target actions (steering angle in the illustrated example). Then, the green curve identifies the simulated agent response time once a low-pass filter is introduced. Instead, the orange curve shows the behavior of the autonomous vehicle while performing the same steering action. However, as we can notice from the figure, the difference in the response time between simulated and real vehicle still remains relevant.

Indeed, the acceleration and steering angle set points predicted by the neural network (Fig. 3) are not feasible commands and do not take into account some factors such as the inertia of the system, the delay of the actuators and other non-idealities. For this reason, a model consisting in a small neural network composed by 3 fully connected layers (which we will call *deep_response*) has been developed in order to reproduce as realistically as possible, the dynamics of the real vehicle. A plot of the *deep_response* behavior is illustrated by the red dotted curve in Fig. 4 and it is possible to notice that is very similar to the orange curve which represent that of the real self-driving car. Given the absence of obstacles and traffic vehicles in the training scenarios, the described problem was more evident for the actuation of the steering angle, but the same idea has been applied to the acceleration output.

We trained the *deep_response* model using a dataset col-

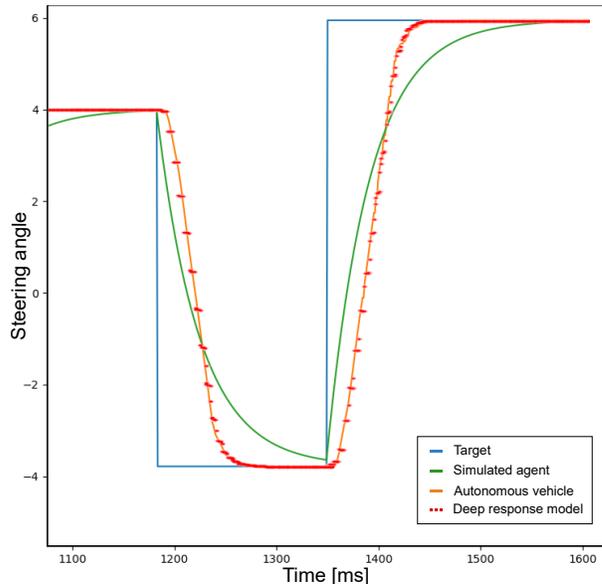


Fig. 4: Comparison between simulated agents (green curve), real vehicle (orange curve) and *deep_response* model (red dotted curve) behaviors executing target actions (blue curve).

lected on board of the autonomous car, in which the inputs correspond to the commands given to the vehicle by a human driver (accelerator pressure and steering wheel movements) and the output correspond to the throttle, brake and curvature of the vehicle that can be measured with GPS, odometry or other techniques. In this way, we embedded such model in the simulator in order to obtain a more scalable system aims at reproducing the self-driving car behavior. Therefore, the *deep_response* module was essential for the correction of the steering angle but even, if in a less evident way, it is also necessary for the acceleration and this will be clearly visible with the introduction of obstacles.

V. EXPERIMENTS

To validate the proposed system, some experiments were conducted both in simulation and in the real world. The input of the neural network used in the autonomous vehicle can be obtained with perception algorithms and through HD Maps; in this way, it is therefore possible to reconstruct the same input used in simulation.

Initially, we tested two different policies on real data to verify the impact of the *deep_response* model on the system. Subsequently, we verified that the vehicle followed the path correctly and that it respected the speed limits retrieved by the HD Maps. Finally, we demonstrate that a pre-training of the neural network with Imitation Learning drastically reduces the overall training time.

A. Test on Real Data

We analyze the behavior of the vehicle over the entire mapped area with two different policies:

- *Policy 1*: trained without using the *deep_response* model, but simulating the response of the real vehicle to a target action with a low-pass filter.

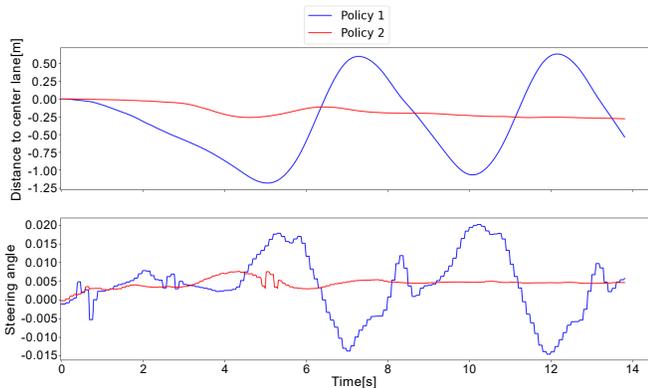


Fig. 6: The distance to the center lane and the steering angle output predicted by *Policy 1* (blue curves) and *Policy 2* (red curves) on a short time window of the real world test on board the self-driving car.

- *Policy 2*: trained by introducing the *deep_response* model to ensure a more realistic dynamics.

The tests performed in simulation led to excellent results with both policies. Indeed, the agent was able to reach the goal in 100% of cases with a smooth and safe behavior, both in training scenarios and in those parts of the mapped area not contained in the training environments.

Different results were obtained by testing the policies on real world scenarios. *Policy 1* is not able to handle the vehicle dynamics, which will execute the predicted actions differently compared to agent in simulation; in this way, *Policy 1* will observe unexpected states as a consequence of its predictions, causing noisy and uncomfortable behavior on board of the self-driving car.

This behavior also affects the reliability of the system, indeed sometimes it was necessary an human assistance to avoid the autonomous car going out of road.

On the contrary, *Policy 2* has never required human assistance during all the real world tests performed on self-driving

car since it is aware of the vehicle dynamics and so how the system will evolve to the predicted actions. The only cases in which a human interventions has been required is to avoid other road users; however, we did not consider these cases as failures since both *Policy 1* and *Policy 2* are trained in obstacle-free scenarios.

To better understand the difference between *Policy 1* and *Policy 2* we plot (Fig. 6) the steering angles predicted by the neural network and the distance to the center lane in a short time window of the real world tests. We can notice how the behavior of the two policies are completely different noticing that the *Policy 1* (blue curve) is noisy and unsafe compared to the *Policy 2* behavior (red curve), proving that the *deep_response* module is essential for the deployment of the policy on board of the real self-driving car.

An example of the real world test performed using the *Policy 2* is illustrated in video¹.

B. Speed Limits

In the second part of the work we verify that the vehicle was able to follow the path and maintain the speed limits. In Fig. 5 we illustrate the behavior of the vehicle trained with *Policy 2* in a real urban area. In particular, the Fig. 5a represents a portion of the map (Fig. 1) with the speed limits of $4 \frac{m}{s}$, $5 \frac{m}{s}$ and $8.3 \frac{m}{s}$, corresponding to orange, blue and green curve colors respectively. In order to show a more detailed example of the longitudinal behavior of the real self-driving car, we plot the accelerations predicted by the network (Fig. 5b) and the vehicle speeds (Fig. 5c) obtained driving along the route of Fig. 5a. The results shown in the plots proving that the system is able to maintain the speed limits retrieved by the HD Maps. Looking at the trend of the acceleration curve in Fig. 5b, it is also important to notice that even if the acceleration output seems noisy, the difference between two consecutive predicted values never exceeds the threshold

¹<https://drive.google.com/file/d/1ZCI2qqNPY2CsE4U1xKfEqGzNintDxf9Q/view?usp>

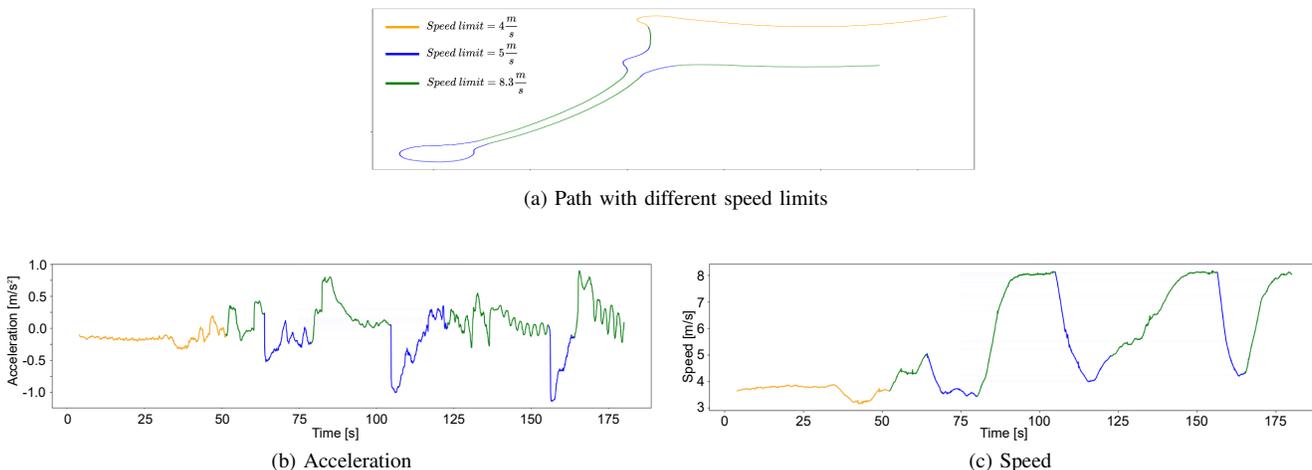


Fig. 5: Fig. 5a represents the path performed with the self-driving car in which colors identify the different speed limits faced in this area. Fig. 5b and Fig. 5c are respectively the trend of acceleration and speed during the real test.

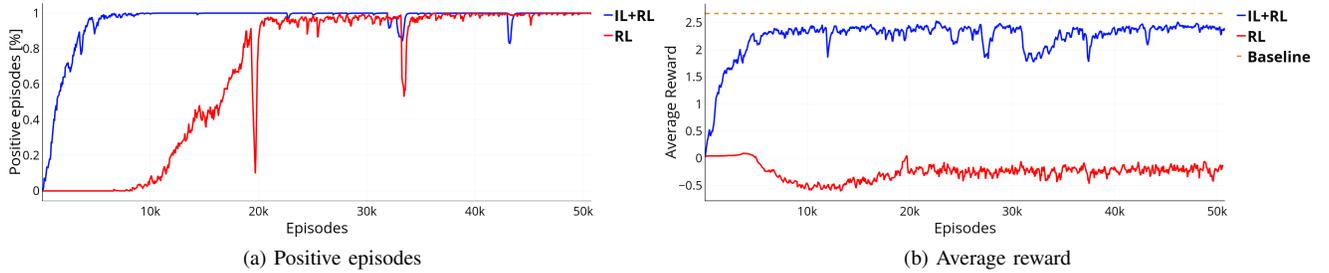


Fig. 7: Comparison between *Pure RL* training (red curve) and using a pre-training through Imitation Learning, *IL + RL* (blue curve). In Fig. 7b the orange dashed line represents the average rewards obtained in 50000 episodes using simulated agents that follow deterministic rules as those used for collecting the dataset for IL pre-training.

δ_{acc} used in the reward function in Equation 5, resulting in a smooth and comfortable longitudinal behavior perceived on-board vehicle.

C. Imitation Learning Pre-training

In order to overcome a known limit of RL that is the need for millions of episodes to reach the optimal solution, we performed a pre-training through Imitation Learning. Moreover, even if the trend in using IL is to train large models, we used the same small neural network (Fig. 3) (about a million parameters), since the idea is to continue training the system using the RL framework in order to ensure more robustness and generalization capabilities. In this way, we do not increase the hardware resources usage and it is essential considering possible future multi-agent trainings [36].

The dataset used during the IL training phase is generated by simulated agents that move following rule-based approaches. In particular, for the curvature we use the tracking algorithm called Pure Pursuit [37], where the agent's purpose is to move following specific waypoints. Instead, we use the Intelligent Driver Model (IDM) [38] to control the longitudinal acceleration of the agent.

To create the dataset, the rule-based agents move on the four training scenarios (Fig. 1) and every 100 ms they save the scalar parameters and the four visual input. Instead, the output is given by the Pure Pursuit algorithm and the Intelligent Driver Model.

The two lateral and longitudinal controls, which correspond to the output, only represent the tuple (μ_{acc}, μ_{sa}) . Consequently, during the IL training phase we do not estimate the values of the standard deviations $(\sigma_{acc}, \sigma_{sa})$ and neither the value functions (v_{acc}, v_{sa}) . These features, together with the *deep_response* module, are learned during the *IL + RL* training phase.

In Fig. 7 we show the comparison of the results obtained by training the same neural network starting with the pre-training phase (blue curve, *IL + RL*) and using only Reinforcement Learning (red curve, *Pure RL*) on all the four scenarios (Fig. 1). Both approaches achieve good success rates (Fig. 7a), even if *IL + RL* training requires less episodes than the *Pure RL* and the trend is also more stable.

Furthermore, the reward curves illustrated in Fig. 7b prove that the policy obtained with the *Pure RL* approach (red curve) did not even achieve an acceptable solution requiring more training time, while the *IL + RL* policy reaches an optimal solution in few episodes (blue curve in Fig. 7b). In this case an optimal solution is represented by the dashed orange curve. This baseline represents the average rewards obtained using simulated agents that perform 50000 episodes in the 4 scenarios of Fig. 1. The simulated agents follow deterministic rules as those used for collecting the dataset for IL pre-training, which therefore use Pure Pursuit for curvature and IDM for longitudinal acceleration. This gap between the two approaches may be more evident training the system performing more complex maneuvers in which agents interactions could be required.

VI. CONCLUSIONS

In this paper we developed a planner based on Reinforcement Learning that allows the vehicle to drive safely and comfortably in an urban environment. The agents in simulation were trained using a neural network that predict acceleration and steering angle each 100 ms. However, at this stage of work, we used obstacle-free training environments in which agents were only trained to follow their paths and to observe the road speed limits.

We also presented the *deep_response* system that allows to reproduce the real vehicle behavior in simulation. In particular, the simulated agents learn to perform target actions evolving their states with a dynamics as similar as possible to that of the real vehicle. This allowed us to achieve a more comfortable driving style testing the system on real scenarios. Currently, no studies have been conducted to verify the generalization capability of the *deep_response* model, because a fleet of different autonomous vehicles would be required to test the system.

With the current system we performed real tests in the mapped area (Fig. 1), that correspond to a low-traffic neighborhood of Parma. We observed that the system reached good performances in the entire area and was able to drive smoothly both in those parts used as training scenarios and in unknown ones, proving that did not overfit on the only training environments, but it have a good generalization

capability.

Finally, we have also shown that with a pre-training using Imitation Learning, we can drastically reduce training times compared to the pure Reinforcement Learning.

Future developments of this work are aimed at creating a multi-agent system in which the vehicle is able to perform other types of maneuvers, such as Adaptive Cruise Control, obstacle avoidance or intersection handling.

Furthermore, an interesting future development would be to solve the current limit on the size of visual inputs and consequently of their content. Currently the visual inputs (Fig. 2) contain a surrounding of 50×50 meters of the agent and this allows us to represent only an area of 40 m forward and 25 m on the sides. However, choosing a larger visual input definitely requires a more complex feature extractor than the current one. A possible solution could be the Variational Autoencoder (VAE) [39], but also in this case we will have to be careful to balance the system to avoid too long training times and computationally too expensive processes.

REFERENCES

- [1] G. Marsden, M. McDonald, and M. Brackstone, "Towards an understanding of adaptive cruise control," *Transportation Research Part C: Emerging Technologies*, vol. 9, no. 1, pp. 33–51, 2001.
- [2] C. Chen, A. Seff, A. Kornhauser, *et al.*, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE international conference on computer vision*, pp. 2722–2730, 2015.
- [3] Y. Ma, Z. Wang, H. Yang, *et al.*, "Artificial intelligence applications in the development of autonomous vehicles: a survey," *IEEE/CAA Journal of Automatica Sinica*, vol. 7, no. 2, pp. 315–329, 2020.
- [4] K. Arulkumar, M. P. Deisenroth, M. Brundage, *et al.*, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [5] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [6] D. Silver, J. Schrittwieser, K. Simonyan, *et al.*, "Mastering the game of go without human knowledge," *nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [7] V. Mnih, K. Kavukcuoglu, D. Silver, *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [8] D. Silver, T. Hubert, J. Schrittwieser, *et al.*, "A general reinforcement learning algorithm that masters chess, shogi, and go through self-play," *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.
- [9] Y. Duan, X. Chen, R. Houthoofd, *et al.*, "Benchmarking deep reinforcement learning for continuous control," in *International conference on machine learning*, pp. 1329–1338, PMLR, 2016.
- [10] B. R. Kiran, I. Sobh, V. Talpaert, *et al.*, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [11] P. Wang, C.-Y. Chan, and A. de La Fortelle, "A reinforcement learning based approach for automated lane change maneuvers," in *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1379–1384, IEEE, 2018.
- [12] C.-J. Hoel, K. Wolff, and L. Laine, "Automated speed and lane change decision making using deep reinforcement learning," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2148–2155, IEEE, 2018.
- [13] B. Mirchevska, C. Pek, *et al.*, "High-level decision making for safe and reasonable autonomous lane changing using reinforcement learning," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2156–2162, IEEE, 2018.
- [14] A. E. Sallab, M. Abdou, E. Perot, *et al.*, "End-to-end deep reinforcement learning for lane keeping assist," *arXiv preprint arXiv:1612.04340*, 2016.
- [15] A. Feher, S. Aradi, and T. Becsi, "Q-learning based reinforcement learning approach for lane keeping," in *2018 IEEE 18th International Symposium on Computational Intelligence and Informatics (CINTI)*, pp. 000031–000036, IEEE, 2018.
- [16] M. Kaushik, V. Prasad, *et al.*, "Overtaking maneuvers in simulated highway driving using deep reinforcement learning," in *2018 IEEE intelligent vehicles symposium (iv)*, pp. 1885–1890, IEEE, 2018.
- [17] L. García Cuenca, E. Puertas, J. Fernandez Andres, *et al.*, "Autonomous driving in roundabout maneuvers using reinforcement learning with q-learning," *Electronics*, vol. 8, no. 12, p. 1536, 2019.
- [18] A. Capasso, G. Bacchiani, and D. Molinari, "Intelligent roundabout insertion using deep reinforcement learning," in *Proceedings of the 12th International Conference on Agents and Artificial Intelligence - Volume 2: ICAART*, pp. 378–385, INSTICC, SciTePress, 2020.
- [19] V. Mnih, A. P. Badia, M. Mirza, *et al.*, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*, pp. 1928–1937, PMLR, 2016.
- [20] A. P. Capasso, G. Bacchiani, and A. Broggi, "From simulation to real world maneuver execution using deep reinforcement learning," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1570–1575, IEEE, 2020.
- [21] A. P. Capasso, P. Maramotti, A. Dell'Eva, *et al.*, "End-to-end intersection handling using multi-agent deep reinforcement learning," in *2021 IEEE Intelligent Vehicles Symposium (IV)*, pp. 443–450, IEEE, 2021.
- [22] R. Liu, J. Wang, and B. Zhang, "High definition map for automated driving: Overview and analysis," *The Journal of Navigation*, vol. 73, no. 2, pp. 324–341, 2020.
- [23] A. Hussein, M. M. Gaber, E. Elyan, *et al.*, "Imitation learning: A survey of learning methods," *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–35, 2017.
- [24] A. E. Sallab, M. Abdou, E. Perot, *et al.*, "Deep reinforcement learning framework for autonomous driving," *Electronic Imaging*, vol. 2017, no. 19, pp. 70–76, 2017.
- [25] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," tech. rep., California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [26] B. Wymann, E. Espié, C. Guionneau, *et al.*, "Torcs, the open racing car simulator," *Software available at http://torcs.sourceforge.net*, vol. 4, no. 6, p. 2, 2000.
- [27] P. Wolf, C. Hubschneider, M. Weber, *et al.*, "Learning how to drive in a real world simulation with deep q-networks," in *2017 IEEE Intelligent Vehicles Symposium (IV)*, pp. 244–250, IEEE, 2017.
- [28] T. P. Lillicrap, J. J. Hunt, A. Pritzel, *et al.*, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [29] A. Kendall, J. Hawke, D. Janz, *et al.*, "Learning to drive in a day," in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8248–8254, IEEE, 2019.
- [30] M. Bojarski, D. Del Testa, D. Dworakowski, *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [31] M. Bojarski, P. Yeres, A. Choromanska, *et al.*, "Explaining how a deep neural network trained with end-to-end learning steers a car," *arXiv preprint arXiv:1704.07911*, 2017.
- [32] 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)*, pp. 2884–2890, IEEE, 2019.
- [33] A. Dosovitskiy, G. Ros, F. Codevilla, *et al.*, "Carla: An open urban driving simulator," in *Conference on robot learning*, pp. 1–16, PMLR, 2017.
- [34] F. Codevilla, M. Müller, A. López, *et al.*, "End-to-end driving via conditional imitation learning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4693–4700, IEEE, 2018.
- [35] J. Kong, M. Pfeiffer, G. Schildbach, *et al.*, "Kinematic and dynamic vehicle models for autonomous driving control design," in *2015 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1094–1099, IEEE, 2015.
- [36] A. Dorri, S. S. Kanhere, and R. Jurdak, "Multi-agent systems: A survey," *Ieee Access*, vol. 6, pp. 28573–28593, 2018.
- [37] R. C. Coulter, "Implementation of the pure pursuit path tracking algorithm," tech. rep., Carnegie-Mellon UNIV Pittsburgh PA Robotics INST, 1992.
- [38] M. Treiber, A. Hennecke, and D. Helbing, "Congested traffic states in empirical observations and microscopic simulations," *Physical review E*, vol. 62, no. 2, p. 1805, 2000.
- [39] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.