

# Deep Reinforcement Learning with Continuous Control in CARLA

Moritz Zanger<sup>1</sup> Florian Rottach<sup>1</sup> Izel Kilinc<sup>1</sup>

Abstract—Style:

- Written last
- 100-250 words
- Past tense

Content:

- Condensed version of entire article

## I. Introduction

The task of teaching vehicles how to drive autonomously in urban scenarios is a challenging and complex one to solve. Not only is there the problem of finding the adequate response to a given situation but also the challenge of taking into account the surrounding factors that have an influence on the state that a vehicle is in and its possible actions. To date, most approaches focus on the manual design of behavioral policies, such as defining a driving policy through the use of annotated maps (7). While these solutions might work in situations which are documented by the provided mapping infrastructure, they are often difficult to generalize or scale, as they do not necessarily enable the comprehension of any given local scene. In order to make autonomous driving truly feasible in a real-world scenario it would be better to develop systems which are able to find their way without having to rely on an explicit set of rules. One possible solution to this task is provided by reinforcement learning methods. Here, the agent, i.e. the vehicle, actively searches for the optimal driving policy whilst trying to maximize a numerical reward signal. As opposed to imitation learning techniques, which have been popular in finding driving policies (1), reinforcement learning algorithms enable a car to exceed human abilities, if applied correctly. In recent years, deep reinforcement learning methods have proven to be successful in solving complex tasks such as playing GO (2) or Atari (3) and there have been efforts in tackling various problems in the field of autonomous driving, including continuous control tasks (4).

However, two major drawbacks of reinforcement learning methods are their heavy dependency on adequate input state representations (5) and, as with other machine learning techniques, their need of a sufficient amount of accurate sample data to train on. In order to be able to train safely on an adequate amount of data, one approach is the use of data from other domains. For this purpose, the urban driving simulator CARLA has been

developed, which is used as a simulation environment for this project.

In previous projects, CARLA has been used to learn policies that allow agents to navigate through complex urban environments via imitation learning methods (15)(12). In this paper, it serves as an environment to solve a continuous control task via deep reinforcement learning. Several state-of-the-art reinforcement learning algorithms are implemented and compared, with regard to their performance considering driving tasks of increasing complexity. Additionally, reward functions for the respective problems are tested and different input representations are designed and evaluated.

## II. Related Work

Various reinforcement learning tasks have been tackled successfully in the past, leading to notable advances in learning policies for simulated and real-world robotic control problems (Levine and Koltun, de Bruin et al.) or solving complex game structures (9). The resulting deep reinforcement learning methods have recently been developed so far as to achieving similar results or beating human experts in some of these tasks (2)(8). In their paper, Silver et al. (2) introduce a new approach of teaching an agent the game of GO by combining Monte Carlo simulation with value and policy networks to evaluate board positions and choose new moves. In autonomous driving, Wolf et al. have achieved notable results using a Deep Q Network to evaluate possible steering actions from an action space of five predetermined steering options (10).

However, many tasks that could be solved with reinforcement learning approaches have continuous action spaces, to which approaches such as the DQN cannot be directly applied in a feasible manner (4). Therefore, algorithms that can handle these high-dimensional action spaces have been developed in recent years. Mnih et al. present asynchronous variations of standard reinforcement learning algorithms and report the Asynchronous Advantage Actor-Critic (A3C) to have achieved the best results in their evaluations, including a physics control task with continuous action space in the MuJoCo Physics Simulator (6). Lillicrap et al. introduce a different algorithm called Deep Deterministic Policy Gradient (DDPG), which they also evaluate in simulated physical environments, receiving near perfect results (3).

DDPG is also used by Kendall et al. to teach a full-sized autonomous vehicle a lane-following policy in an on-board manner. While already achieving good

<sup>1</sup>The authors are with FZI Research Center for Information Technology, Haid-und-Neu-Str. 10-14, 76131 Karlsruhe, Germany {uydrn, uzdnz}@student.kit.edu

results with DDPG alone, they show that the additional use of a Variational Autoencoder greatly improves the overall performance (1), suggesting an improved state representation as an area of possible further development in reinforcement learning for continuous control tasks. Another proposition to reduce the complexity of visual features that has investigated in several projects is the representation of the input images as a top-down view on the vehicle (13)(12)(11).

We extend this research by applying the aforementioned algorithms to autonomous driving tasks and evaluating their performance in OpenAI’s CarRacing-v0 environment and in the CARLA driving simulator (14) while investigating input states of differing complexity, such as bird’s-eye view or latent space representations.

### III. Background

We regard the typical reinforcement learning setting where an agent interacts with an environment  $\mathcal{E}$ . At each one of a number of discrete timesteps  $t$  the agent decides on taking an action  $a_t$  from a given set of actions  $\mathcal{A}$ . This is done based on the state  $s_t$  that the agent is currently in and following a policy  $\pi$ , which is a mapping of the possible states to the action space  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ . In this case,  $\pi$  is stochastic, as it returns the probability distribution over the possible actions, and the action space is continuous.

As a result, the agent receives a reward  $r_t$  and the subsequent state  $s_{t+1}$  of his environment. The setup is assumed to follow the properties of a Markov decision process, where besides state space  $\mathcal{S}$ , action space  $\mathcal{A}$  and reward function  $r(s_t, a_t)$ , we also include the transition function to the future states  $p(s_{t+1}|s_t, a_t)$ . The objective of the agent is to maximize his expected return  $R_t$ , which is the cumulative reward of the current state  $s_t$  and the rewards of future states, discounted with a factor  $\gamma \in [0, 1]$ .

In actor-critic methods the critic calculates the action value  $Q^\pi(s, a) = E[R_t|s_t = s, a]$  describing the expected cumulative return for taking action  $a$  in state  $s$  under the given policy  $\pi$ , thus evaluating the selected action. The actor on the other hand estimates an optimal policy, following  $\pi(s) = \operatorname{argmax}_a Q(s, a)$ . Similar to the action value,  $V^\pi(s) = E[R_t|s_t = s]$  defines the value of state  $s$  under the policy  $\pi$  and simply describes the expected return for pursuing a policy  $\pi$  from state  $s$ .

In our project, we consider a continuous control task in the context of autonomous driving where three different actions can be selected: steering in a range of  $[-1, 1]$ , throttle and acceleration in a range of  $[0, 1]$  respectively.

### IV. Concept/Methods and Models

In the course of this project, three different algorithms are used. Following the work of Lillicrap et al. (4) and Mnih et al. (6), the (DDPG)- and (A3C) algorithms are implemented and compared according to their performance in the Gym environment CarRacing-v0. Later

on, the PPO algorithm is applied to solve a continuous control task in the CARLA simulator. TODO: Hier unser Konzept analog zur Präsentation vorstellen (Slide mit den 4 Kästen)

#### A. DDPG

One very notable advance in reinforcement learning has been made by the development of the so-called ”Deep Q Network” (Mnih et al., 2015). The DQN is able to solve tasks with high-dimensional observation spaces. However, it is only efficiently capable of working with discrete and low-dimensional action spaces. In order to adapt a (DQN) for the successful use with continuous control problems, as given in CarRacing-v0, a discretization of the action space has to be carried out, which can lead to two main difficulties: an explosion in the number of possible actions and the loss of important information (4).

To evade these obstacles (Lillicrap et al.) propose a new approach, called the Deep Deterministic Policy Gradient (DDPG), which is a model-free, off-policy actor-critic algorithm. They adopt the advantages of (DQN) and combine them with the actor-critic framework, resulting in the stabilization of Q-learning by using a replay buffer and soft updates on the target networks  $\theta'$  of both actor and critic, through

$$\tau < 1 : \theta' \leftarrow \tau\theta + (1 - \tau)\theta', \quad (1)$$

as well as finding a deterministic policy.

In order to enable better exploration, we add noise in the form of the Ornstein-Uhlenbeck process to our policy. For our continuous control problem, we chose the Ornstein-Uhlenbeck process due to the fact that it produces temporally correlated noise, which enables smoother transitions, compared to e.g. Gaussian noise.

#### B. A3C

The Asynchronous Advantage Actor-Critic (A3C), introduced by Mnih et al. [?], has become a go-to algorithm in Deep Reinforcement Learning due to its performance, robustness, and ability to perform well on high-dimensional action- and state-spaces. A key characteristic of A3C is the utilization of multiple agents, each equipped with its own environment instance and its own set of network parameters. One of the advantages of this approach is the diversification of the collected experience but yields the challenge of handling gradient update mismatches between the asynchronously collected network parameter updates. In our implementation of A3C we adapted several alterations as opposed to the original version by Mnih et al.. Despite being an on-policy method, we decided to include a numerically efficient n-step return as proposed by jaromiru [?].

Furthermore, experiences are buffered in a global update queue and updates are only performed by a master network. This feature helps decorrelating experiences, emphasizes exploration and allows for more gpu-friendly

batch-learning. However, one has to keep in mind that this contradicts the original update rule by mnih et al. and might lead to policy lag. A proper update frequency in the master network is therefore of major importance.

### C. PPO

The Proximal Policy Optimization (PPO) algorithm was developed to improve training stability by avoiding parameter updates that excessively change a given policy at each step. The algorithm that PPO is based on is called Trust Region Policy Optimization (TRPO) and realizes training stability through a KL divergence constraint on the policy update range at each iteration (16). PPO simplifies this concept by incorporating the hard constraint into the objective function through Lagrangian duality and, thus, softening the constraint. Additionally, PPO clips the objective function if the ratio between old and new policy lies outside the range of  $[1-\epsilon, 1+\epsilon]$ , which discourages large policy changes. Tests on benchmark tasks have shown PPO to achieve outstanding results while reducing TRPOs complexity significantly (17).

### D. Reward function

The design of the reward function turned out to be a much more difficult procedure as expected initially. The already provided rewards in gym's CarRacing environment were rather simple and still led to a good result in the end. For CARLA however, it was necessary to invest more time into the engineering and designing of a suitable inducement system, because the driving behavior depended on more factors and the performance metric was not just driving as fast as possible. In the following, the process of arriving at our final reward function will be described.

In the first step, we investigated possible sensor and measurement inputs that might have an impact on the driving behavior and discussed on their impact. As summarized in Table 1, the outcome consists of eight possible features, that were tested in the following.

We started with incrementally adding these attributes to our reward calculation and quickly realized, that the main challenge is adjusting the weights and harmonizing the contrary effects of the terms. An example would be, that giving the velocity a relatively high weight, such as 0.8, while giving the distance to center line a weight of 0.2 results in an agent that speeds over the map and pays only few attention on lane invasions. On the opposite side, the agent will drive very slowly or even not at all, if the rewards for oscillations are too high compared to the velocity. Considering, that we have not only two, but several possible components, it results in a complex combinatorial problem that can either be solved by trail and error or applying permutational optimization techniques. To achieve initial results, we started to discover a proper reward function "by hand". We pruned the above list by applying the following

TABLE I: Identified reward function components

Component	Description	Intention
Per frame penalty	The penalty amount subtracted from the reward for each frame	Forces agent to move in order to compensate the penalties
Lane invasion increment	Is either 0 or 1 and reflects if a lane invasion took place for the current frame.	The agent should perform as few lane changes as possible
Steering angle	The absolute angle of the steering wheel	Avoids oscillations in the driving behavior
Delta heading	The angle between current street direction and car position	Agent should drive as straight as possible relatively to road direction
Position change	The angle between current street direction and car position	Maximize travelled distance
Collision binary per frame	Penalty for crashing into other vehicles, objects or road infrastructure	Avoid crashes and improve security of driving
Velocity	The agent's absolute speed retrieved from the CARLA engine	Forces the agent to move forward
Distance to middle lane	The absolute distance to center in meters - can also be squared to improve	Drive as centered as possible

considerations that resulted from testing different approaches. Firstly, some components show a redundant behavior and hence one of them can be removed. An example would be the velocity and the position change - both contain the same information. Further, the attributes lane invasion and steering angle didn't affect the driving behavior in a positive way. The two most important parameters turned out be the velocity and the delta heading, which expressed the relative angle to the current street angle. Our final reward function and the aggregation weights can be found in table 2.

We want to note here, that a performance-wise comparison between the different parameter combination is hard to measure with a metric and is not the aim of this project. However, it would make sense to develop a suitable approach to categorize and assess the performance of different reward functions towards the desired outcome.

The result (in combination with tuning the reinforcement learning models) is very promising and can be summarized as follows: The agent is capable of driving smoothly within the right lane and can perform turn

TABLE II: Final reward function terms and weights

Component	Effect	Weight
Per frame penalty	Lead to a strong initial learning behavior. Already after few steps the agent kept accelerating to compensate this penalty.	-0.01
Velocity	When selected too small, slow learning and when too high, strong lane oscillations / non-smooth driving.	+0.05
Delta heading	Introducing this term improved the driving stability enormously.	-0.005
Squared distance to middle lane	Squaring had very beneficial effects, other powers were to restrictive.	-0.01
Collision binary	Might be set to an even higher value but lead to attempts to avoid other objects	-100

maneuvers on most intersections. It attempts to drive around other vehicles, but is not capable of breaking. We trained different models on this reward function and all of them had a good performance compared to other incentive attempts. We assume, that even better results can be achieved when applying an improved optimization method to find out the best parameter-combination. However, this is a combinatorial problem and requires several simulation runs.

#### E. Input representation

Dealing with the high-dimensional environment in CARLA is one of the central challenges in implementing a well functioning RL algorithm. In consideration of an abundance of possible sensor types available in CARLA we decided to pursue increasing levels of realism which generally also correspond to increasing levels of difficulty. These aforementioned levels of realism involve the following:

- Ground truth segmented bird’s eye view
- Ground truth segmented front view
- Latent space generated from ground truth segmented images
- Latent space generated from rgb images

Notably, these input representations differ in dimensions and thus lead to different network architectures in the appended network architectures of the RL agent.

a) Ground truth segmented bird’s eye view: The Ground truth segmented bird’s eye view is a rather

unrealistic scenario, in which we assume availability of a camera positioned 20m above the performing agent. It is merely imaginable in a fully observed city where autonomous driving is part of a high-level traffic control system. Albeit rather unrealistic, we decided to implement a 13-class ground truth segmented bird’s eye view due to its similarity to the CarRacing environment and its obvious advantages in containing information central to navigational tasks. The full list of the included classes are described by dosovitskiy et al.[?]. We deemed a 1-channel grayscale image with size 64x64 large enough to contain key information for the algorithm (fig. x l.)

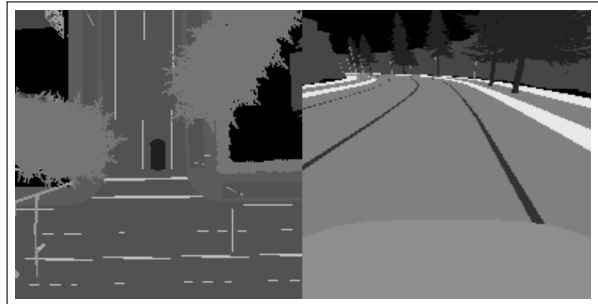


Fig. 1: left: Ground truth segmented bird’s eye view image, 64x64 in grayscale, 9 classes right: Ground truth segmented front view image, 80x80 in grayscale, 9 classes

b) Ground truth segmented frontview: Much like the Ground truth segmented bird’s eye view, the corresponding front view input representation assumes the availability of a perfect segmentation camera. Nonetheless, we increase realism with this representation as the camera is placed in front of the car. Again, we chose a 13-class, 1-channel grayscale image with a slightly increased size of 80x80 to account for a larger view angle of the camera.

c) Latent space generated from ground truth segmented images: In this section we will further discuss a modified version of the integrated encoder-decoder network that is based on dziubinski’s [?] implementation. The underlying idea of this model is to guide feature extraction towards more useful, problem-specific features by exposing the model to the additional target of reconstructing a bird’s eye view from the vehicle-based camera images. In its original architecture, the network comprises 5 input tensors and 7 output tensors with the inputs representing images of a front-, rear-, left-, right-, and top-view camera. The model is built with the Keras functional API and generally serves two purposes. On the one hand, each input tensor is encoded into a 64-sized feature vector and decoded into it’s original shape with a cross entropy loss with regard to the original input. This is achieved with separate branches of autoencoders with 3 convolutions respectively. On the other hand, a separate generative branch creates a bird’s eye view reconstruction based on the concatenated feature vectors of the vehicle-based cameras. In addition to the cross entropy loss at the reconstructed bird’s eye view output,

a mean square error loss is calculated against the output of a subtract layer between the autoencoder’s feature vector and the reconstructed feature vector to support convergence. Notably, the autoencoder branch of the top-down camera view is solely purposed for improving the latent space of the generative branch during training and is thus not required for inference.

In comparison to dziubinski’s [?] vanilla version, we adjusted the architecture to fit the single camera bottlenecks into a length 64 vector and the reconstructed bird’s eye view bottleneck into a length 128 vector. To reduce complexity, we condensed both input and target segmentation images to 3 classes and implemented a generalized weighted cross entropy[?] loss function to account for the unbalanced distribution of vehicles and obstacles.

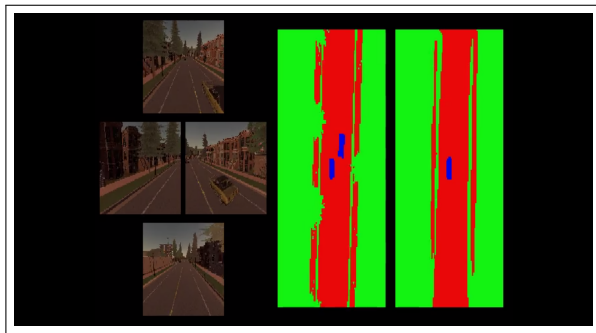


Fig. 2: left: rgb images from 4 vehicle-based cameras  
2nd from right: ground truth segmented bird’s eye view  
right: reconstructed bird’s eye view

d) Latent space generated from rgb images: To deal with the higher input complexity in RGB-images, we extended the encoder and decoder models to 5 convolutions.

#### F. Autoencoder/Latent space

TODO: Encoder-decoder architecture bzw. Generator nennen

#### G. Training

In order to align the algorithms with the varying input dimensions we had to implement a wrapper around the actual models that adjusts the network architecture accordingly. Especially the fact that the latent space is flattened and not 2-dimensional such as the ground truth pictures made it necessary to implement convolutional as well as fully connected networks. Luckily, for the stable-baselines PPO model this is implemented very quickly. The architecture of our Keras-DDPG model however, had to be adjusted completely. We performed the training in four process steps, aligned with the input representation types mentioned previously. Therefore, the initial training was carried out on the ground-truth front view, afterwards on the ground-truth bird’s eye view and finally we trained on the latent space of rgb and ground-truth segmented images.

1) DDPG: The Deep Deterministic Policy Gradient performed well on our initial attempts within gym CarRacing and we identified reasonable hyperparameters for proper learning. We applied the same model on the new input features of CARLA and came to the sobering result that the input seemed to be too complex for the model. We assume, that the learning algorithm was distracted by the advanced simulation environment. The following illustration emphasizes the poor behavior of the model (trained on circumstances where other models performed well). Hence, we decided to continue with different learning approaches and agreed due to recent developments on the PPO.

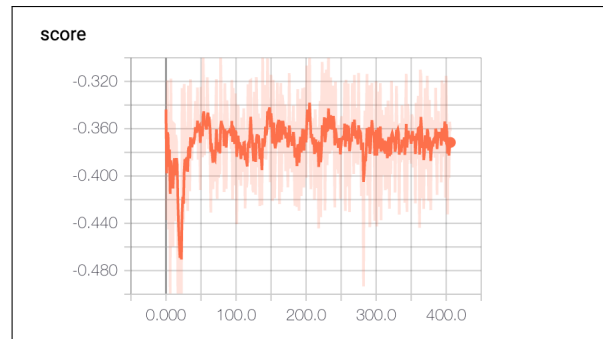


Fig. 3: Keras DDPG learning behavior

2) PPO: The stable-baselines PPO2 model is easy to train and only few hyperparameters need to be adjusted. The drawback of this approach is that customization of these pre-defined models is rather laborious. Especially a combination of diverse input types such as scalars and images requires arduous architecture adjustments. We changed the following parameters and arrived at a very fast and crisp learning behavior that enables the agent to start driving after already very few episodes.

TABLE III: Final reward function terms and weights

Parameter	Default Value	Selected Value
Learning rate	0.00025	0.0004
Clip range	0.2	0.1
Gamma	0.99	0.97
N_steps	128	1024
Environment steps	25k	200k

The learning behavior usually increased until 200k steps and started to stagnate afterwards. The achieved episode rewards and entropy loss over several runs was very stable and performance drops occurred rarely. The results on the proved input representations were very different: We trained the PPO with the same reward function, the same amount of steps at the same environment spawn point and hence assumed equal conditions. The front-cam-view resulted in a very oscillating driving

behavior and the rewards were damped by the penalties for inaccurate delta heading. The birds-eye-view on the other hand had a very smooth driving behavior and delivered the best performance in our test. This comes from the improved view on the street trajectory from above. Considering the easier and smaller latent space representation we realized that the algorithm trains very fast but can't reproduce the performance of the ground-truth. We assume that the architecture didn't contain the complete information in the layers that were extracted in our experiment and hence the latent space should usually arrive at a similar outcome as the ground-truth birds-eye-view.

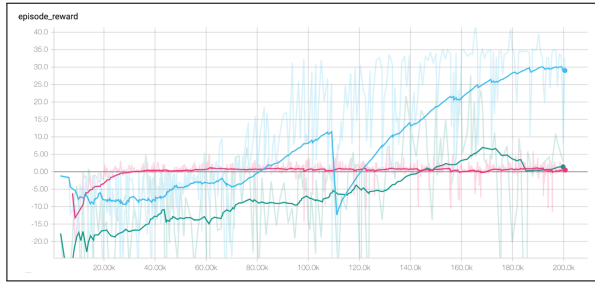


Fig. 4: Green line: Performance on front-cam-view, Blue line: Performance on birds-eye-view, Red line: Performance on latent space

## V. Evaluation

### A. CarRacing-v0

Due to the simple input structure and the straightforward reward system we achieved stable results, even with using not the most-advanced learning algorithms (We consider the PPO as an improvement over the DDPG). Also the A3C could generate promising results (average rewards around 700) with only few training time. In total, these results can be explained by the simplicity of the environment. However, here we invested more time on the finetuning of the hyperparameters compared with CARLA.

### B. CARLA

The CARLA simulation environment is far more complex than the previous environment we worked with. Additionally, it was necessary to implement the observations, terminal conditions as well as the reward function, where a lot of time was invested. The application of the stable-baselines algorithms is rather straightforward. Hence, the challenging part here was merely the environment itself.

### C. Results

TODO: Mit Concept vergleichen und formulieren welche Ziele erreicht wurden / wo es probleme gab (analog zur präsi)

— mo Despite the weighted categorical cross entropy loss, the network struggled with reconstructing pixels of the vehicle and obstacles class. Aside from improving

training data and network parameters, we consider a more sophisticated encoding model as crucial for better results. — mo

### D. Comparison algorithms/reward functions

## VI. Conclusions

As result of this practical seminar we have successfully applied different deep reinforcement approaches to the CARLA environment and finally achieved a stable and promising lane-keeping behavior for the agent. This was achieved by a combination of finding the most suitable input representation for smooth driving, finetuning of the rewardfunction and identification of the best performing hyperparameters. The outcome proves that reinforcement learning can be applied to complex environments as well as long as the desired behavior is incooperated in the reward system. Furthermore, we could prove that the birds-eye-view leads to a better performance as a single front camera and is also feasible in a real-world setup by generating the representation from four cameras surrounding the car. In the training process we recommend to not use early-termination for the episodes, especially if the reward function contains negative components, since this increases the probability that the agent starts to drive off the track as soon as possible. Finally, we emphasize that the driving behavior has the potential to be even further improved in the future by selecting better data, learning algorithms and an optimal reward function. We can also imagine, that the involvement of further attributes in the model such as the current velocity can be used to describe even more complex driving situations.

For the future we additionally suggest a combination of the birds-eye-view with a front camera in order to enable the identification of traffic signs and traffic lights, which wouldn't be possible from above. Also, we assume that the improvement of the encoder-decoder architecture could boost the learning speed and overall performance.

- Similar to the abstract but more detail
- Conclusion of the key points of each section
- Summary of main findings
- Important conclusions that can be drawn
- Discuss benefits and shortcomings of our approach
- Suggest future areas of research