

# Residual Policy Learning for Vehicle Control of Autonomous Racing Cars

Raphael Trumpp, Denis Hoornaert, and Marco Caccamo

**Abstract**—The development of vehicle controllers for autonomous racing is challenging because racing cars operate at their physical driving limit. Prompted by the demand for improved performance, autonomous racing research has seen the proliferation of machine learning-based controllers. While these approaches show competitive performance, their practical applicability is often limited. Residual policy learning promises to mitigate this drawback by combining classical controllers with learned residual controllers. The critical advantage of residual controllers is their high adaptability parallel to the classical controller’s stable behavior. We propose a residual vehicle controller for autonomous racing cars that learns to amend a classical controller for the path-following of racing lines. In an extensive study, performance gains of our approach are evaluated for a simulated car of the F1TENTH autonomous racing series. The evaluation for twelve replicated real-world racetracks shows that the residual controller reduces lap times by an average of 4.55 % compared to a classical controller and even enables lap time gains on unknown racetracks.

**Index Terms**—Residual Policy Learning, Autonomous Racing, Vehicle Control, F1TENTH,

## I. INTRODUCTION

Autonomous racing series, such as the F1TENTH [1] and Indy Autonomous Challenge, are at the frontier of autonomous racing development. In against-the-clock racing, the computed racing lines push the car’s physics to its limits, requiring highly fine-tuned controllers to deliver the expected performance. On the one hand, classical controllers, e.g., PID control, only achieve adequate performance with tedious racetrack-dependent parameter tuning. On the other hand, modern model-based controllers like model-predictive control require accurate vehicle models and have high computational requirements. The advancing field of deep reinforcement learning (DRL) promises to overcome these limitations by learning model-free vehicle controllers through trial-and-error interaction and accumulating a reward signal. It has been shown that end-to-end approaches using DRL [2] on top of a perception module learn competitive racing behavior. However, the practical application of these methods is challenging due to the sim2real gap and ample training time.

The field of residual policy learning (RPL) [3] overcomes these hurdles by combining classical controllers with a DRL-based approach to learn residual actions that correct and improve the performance of the base controller. In the context of autonomous racing, RPL promises to combine the *reliability* of

classical control approaches with the *adaptability* of learning-based controllers. The residual controller can learn the control action under various optimization targets, e.g., energy consumption or vehicle stability, alongside the ability to cope with variations in the environment in real-time by observing the current vehicle state and its environment. Opposed to classical DRL, RPL has proven to bootstrap the learning process leading to reduced training time with high real-world performance [4].

Originally motivated by robotic applications [3], [4], RPL has recently gained attention in the automotive community. In [5], the authors leverage RPL to improve traditional powertrain control policies for better fuel consumption. RPL has also been applied for suspension control to benefit passenger comfort [6]. Closest to our work is the use of RPL for autonomous racing in [7]. The authors propose the ResCar model that combines a modified artificial potential field controller with RPL. Evaluated for simulated F1TENTH cars, the performance is demonstrated in five replicated real-world racetracks showing reduced demand for training data compared to the DreamerV2-based approach in [8]. However, ResCar is trained per racetrack, leading to tedious retraining as the generalization capability of the approach is only discussed for limited scenarios. Similarly, end-to-end DRL approaches [2] offer low generalization capability while requiring higher computational power. In contrast, we show that our method generalizes to unknown racetracks directly. Augmenting model-predictive control for trajectory optimization with DRL [9] is a promising direction with high generalization potential; our work proves that high performance can also be achieved with a fully model-free pipeline and only local observation of a planned trajectory.

We propose the use of RPL to synthesize a vehicle controller for autonomous racing cars consisting of two modules: The first module acts as a *base* controller that has local access to planned waypoints to follow an optimized racing line; this racing line may be obtained by map data in advance or through online observation of the racing car’s environment. The second module, the *residual* controller, uses *local* observations to learn a residual action to amend the base controller’s action – the residual controller accounts for local optimization potential due to non-optimal tracking performance of the base controller, or a non-optimal racing line, to optimize for an overall improved performance. This adaptability can be crucial for achieving the best lap times in competitive racing, as our simulated F1TENTH racing car shows.

Our main contributions can be summarized as follows:

- Design of a novel residual controller combining RPL with path-following of a planned racing line in the F1TENTH-

R. Trumpp, D. Hoornaert, and M. Caccamo are with the TUM School of Engineering and Design, Technical University of Munich, Germany {raphael.trumpp, denis.hoornaert, mcaccamo}@tum.de.

gym simulator [1] for twelve real-world racetracks.

- The influence of RPL for improving vehicle handling is analyzed in a study on vehicle dynamics and lap time.
- We demonstrate the generalization capability of the proposed approach by evaluating its behavior in unknown racetracks and randomized starting positions. This study proves that the training environments are not just memorized; instead, a deep scene understanding is learned.
- We reason about how the proposed residual controller improves lap times by an average of 4.55 % compared to the base controller.
- The code is publicly available<sup>1</sup>.

## II. BACKGROUND

In the following, the used vehicle dynamics model and vehicle controller for path-following are introduced along with the theoretical description of RPL using DRL algorithms.

### A. Vehicle Dynamics

The FITENTH-gym simulator’s vehicle dynamics are derived from the CommonRoad framework [10] using a kinematic single-track model with tire slip; the vertical load of all four wheels, their individual spin and slipe, and non-linear tire dynamics are not represented in this model. The vehicle model’s control input is defined as  $a_t = [\dot{\delta}_t, \dot{v}_t]$  with  $\dot{\delta}_t$  as the steering velocity of the steering angle  $\delta_t$ , and  $\dot{v}_t$  as the acceleration of the vehicle at time  $t$ .

### B. Trajectory Generation

When global map knowledge of the racetrack to drive is available, trajectory optimization techniques can be used to derive a (near) optimal racing line consisting of waypoints  $w = \{x_d, y_d, \phi_d, v_d\}$  that define a global desired position  $[x_d, y_d]$  and orientation  $\phi_d$  together with a planned velocity profile  $v_d$ . The set of waypoints  $W = \{w_1, w_2, \dots\}$  is obtained by minimization of a loss function, e.g., the minimum curvature of the racing line as discussed in [11].

### C. Vehicle Controller for Path-Following

The task in path-following is to follow a planned trajectory with minimal error, i.e., the positional and velocity tracking error must be minimized. While a PID controller is commonly used to track the velocity profile, i.e., a vehicle acceleration  $\dot{v}_t$  is calculated w.r.t. the velocity error, the derivation of a well-performing control law for the steering angle  $\delta_t$  is more involved. The pure pursuit controller [12] calculates  $\delta_t$  based on a geometric definition of the reference trajectory. Given a set of planned waypoints  $W$ , the current lookahead point  $w_t \in W$  at time  $t$  is selected in a fixed lookahead distance  $l_d$  relative to the vehicle position. Using the geometric center of the circle defined by the vehicle’s rear axis, its length  $l$ , and the waypoint  $w_t$ , the steering angle is calculated by

$$\delta_t = \arctan\left(\frac{2l \sin \alpha}{l_d}\right), \quad (1)$$

with  $\alpha$  as the angle between the line connecting the vehicle’s rear with the lookahead point and the vehicle’s orientation.

### D. Residual Policy Learning

In model-free DRL, the behavior of a stochastic agent is defined by a policy  $\pi(a_t|s_t)$  that maps the observed state  $s_t \in \mathcal{S}$  to a probabilistic action space  $\mathcal{P}(\mathcal{A})$  of (continuous) actions  $a_t \in \mathcal{A}$ . The state transition  $\mathcal{T}$  is a mapping  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  defining the transition probability from states to new states when an action  $a_t$  is taken. Such transitions are evaluated by the reward function  $\mathcal{R}$  with a scalar value  $r_t$ . Combined with the discount factor  $\gamma$ , a Markov decision process (MDP) with tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$  is described.

For RPL, the policy  $\pi_{\text{RB}} : \mathcal{S} \rightarrow \mathcal{A}$  is defined with

$$\pi_{\text{RB}}(s_t) = a_{\text{R},t}|_{a_{\text{R},t} \sim \pi_{\text{R}}(\cdot|s_t)} + \mu_{\text{B}}(s_t), \quad (2)$$

combining the deterministic *base* policy  $\mu_{\text{B}}(s_t) : \mathcal{S} \rightarrow \mathcal{A}$  and corresponding action  $a_{\text{B},t}$  with a *learned* stochastic<sup>2</sup> residual policy  $\pi_{\text{R}} : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ . For  $\mu_{\text{B}}(s_t)$ , any deterministic control algorithm can be used while  $\pi_{\text{R}}(a_t|s_t)$  is typically learned by DRL. Because  $\pi_{\text{R}}(a_t|s_t)$  is defined as a probability distribution, actions must be sampled with  $a_{\text{R},t} \sim \pi_{\text{R}}(\cdot|s_t)$ . Note that for the environment to be Markovian, the base action  $a_{\text{B},t}$  must be part of or be observable from  $s_t$ .

### E. Deep Reinforcement Learning Algorithms

DRL combine the use of neural networks (NNs) with an iterative learning approach. Typically, a value  $V(s_t)$  or Q function  $Q(a_t, s_t)$  is learned through interaction with an environment and observable transition tuples  $e_t = \{s_t, a_t, r_t, s_{t+1}\}$  along a trajectory  $\tau = \{e_0, e_1, \dots, e_T\}$  of  $T$  fixed-sized steps.

The *on-policy* algorithm proximal policy optimization (PPO) [13] learns a stochastic policy  $\pi_{\theta_k}$ , parameterized with weights  $\theta_k$  at iteration  $k$ , by updating the policy according to

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{s_t, a_t \sim \pi_{\theta_k}} [L(s_t, a_t, \theta_k, \theta)] \quad (3)$$

with the objective function

$$L(s_t, a_t, \theta_k, \theta) = \min \left( \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_k}(a_t | s_t)} A^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right) \quad (4)$$

where  $g(\epsilon, A)$  is defined as

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0. \end{cases} \quad (5)$$

with the advantage function  $A = A^{\pi_{\theta_k}}(s_t, a_t)$ . In (5),  $\epsilon$  is a hyperparameter that implies that a gradient update step has limited incentive for the new policy to be far away in probability space from the previous one through clipping of the function, i.e., policy changes are regularized. This version of PPO, often called PPO-Clip, requires calculating the advantage

<sup>2</sup>For deterministic policies, (2) can be rewritten to  $\pi_{\text{RB}}(s_t) = \pi_{\text{R}}(s_t) + \mu_{\text{B}}(s_t)$ .

<sup>1</sup><https://github.com/raphajaner/rpl4f110>.

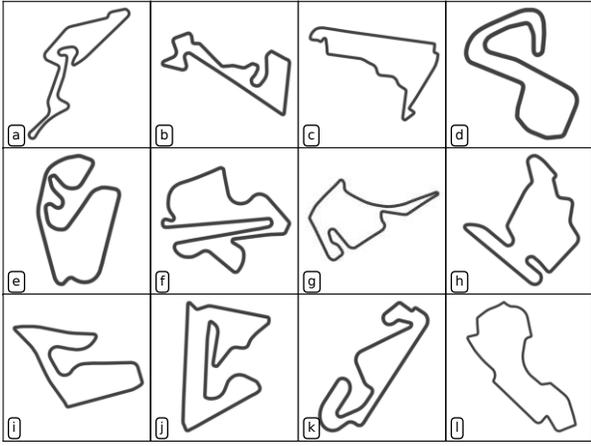


Fig. 1: Racetracks for training and testing: (a) Nürburgring, (b) Moscow Raceway, (c) Mexico City, (d) Brands Hatch, (e) Sao Paulo, (f) Sepang, (g) Hockenheim, (h) Budapest, (i) Spielberg, (j) Sakhir, (k) Catalunya, and (l) Melbourne.

function  $A$ . Typically, the generalized advantage estimate (GAE) [14] approach is used which requires the training of a value function  $V_{\phi_k}(s_t)$  with weights  $\phi_k$  in parallel to the optimization of the policy  $\pi_{\theta_k}$ . In practice, the expectation in (3) is estimated by batches of  $B$  recorded  $T$ -step long trajectories  $\mathcal{D}_k = \{\tau_0, \tau_2, \dots, \tau_{B-1}\}_k$  that must have been realized by using the current policy  $\pi_{\theta_k}$ . The value function  $V_{\phi_k}$  is updated by regression of

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T (V_{\phi_k}(s_t) - \hat{R}_t)^2, \quad (6)$$

with the reward-to-go  $\hat{R}_t = \sum_{t'=t}^T r_{t'}$ .

### III. SYSTEM DESIGN

This work proposes an RPL-based controller for vehicle control of autonomous racing cars. The approach is developed in a simulator for FITENTH cars. In the following, the simulation environment and the components of the developed vehicle controller are presented.

#### A. Learning Environment

The FITENTH-gym environment [1] simulates vehicle dynamics using the single-track model (including slip) as presented in Section II-A. The simulator allows the use of arbitrary map data; we use the replication of twelve real-world racetracks from the FITENTH-racetracks [15] repository<sup>3</sup> adjusted to the FITENTH scale. As it can be seen in their visualization in Fig. 1, these racetracks offer various track designs with challenging high- and low-speed sections.

These real-world FITENTH-racetracks offer an associated racing line for each racetrack in the form of a set of waypoints that include their global positions, an optimized velocity profile, and curvature information. According to the authors

TABLE I: Vehicle model parameters for a FITENTH RC-car.

Parameter	Value	Unit	Description
$h$	0.074	m	Height
$w$	0.27	m	Width
$l$	0.51	m	Length
$l_f$	0.15875	m	Length from COG to front
$l_r$	0.17145	m	Length from COG to rear
$m$	3.47	kg	Total mass
$I$	0.04712	kg·m <sup>2</sup>	Vehicle inertia
$\delta^{\max}$	0.4189	rad	Steering angle bounds
$\dot{\delta}^{\max}$	3.2	rad/s	Steering velocity bounds
$v^{\text{switch}}$	7.319	m/s	Dynamic bound for full acceleration
$v^{\max}$	8.0	m/s	Maximal velocity
$\dot{v}^{\max}$	7.51	m/s <sup>2</sup>	Maximal acceleration
$C_{S,f}$	4.718	N/m	Tire stiffness front
$C_{S,r}$	5.4562	N/m	Tire stiffness rear
$\mu$	0.8	-	Friction coefficient (estimated)

of the repository, the racing lines have been obtained via minimum curvature trajectory planning such as described in [11]. We use realistic vehicle parameters as measured from our real FITENTH car; see Table I for a full list of the used parameters. The simulator runs with a synchronous simulation and control frequency of 100 Hz. As implemented in the simulator, a low-level controller is used that sets the model input  $[\dot{\delta}_t, \dot{v}_t]$  given a high-level control action  $a_t = [\delta_t, v_t]$ .

#### B. Residual Controller

We proposed a residual controller architecture that consists of two sequential control modules:

- 1) The *base* controller consisting of a classical controller outputting a deterministic control action  $a_{B,t}$ .
- 2) The *residual* controller that learns a control action  $a_{R,t}$  by RPL and interaction with the environment.

As shown in Fig. 2, our presented controller obtains the control action  $a_{RB,t}$  by combining both control modules

$$a_{RB,t} = a_{R,t} + a_{B,t}, \quad (7)$$

with  $a_{RB,t} = [v_t, \delta_t]$  as the high-level control input to the simulator's low-level controller.

The base module's action  $a_{B,t}$  is derived from a pure pursuit controller for lateral control; see Section II-C. The steering angle  $\delta_t$  is calculated by minimizing the tracking error of the nearest waypoint of the planned racing line along with the velocity taken from the planned velocity profile of this waypoint; a lookahead distance of 0.82 m is used.

The residual controller is implemented by a PPO agent which observes only local information

$$s_t = [L_t, w_{\text{rel},t}, v_{x,t}, v_{y,t}, \dot{v}_{x,t}, \dot{v}_{y,t}, \psi_t, \dot{\psi}_t, \beta_t, a_{B,t}, a_{RB,t-1}], \quad (8)$$

with  $L_t \in \mathcal{R}$  as a 1D-lidar representation with 1080 points spread along a 270° field-of-view measuring the vehicle's distance to other objects, e.g., the wall defining the racetrack. A limited set of future waypoints  $w_{\text{rel},t}$  are given with relative coordinates to the vehicle's current position. This information is limited to local observation; only the waypoints ahead

<sup>3</sup>[https://github.com/fl1tenth/fl1tenth\\_racetracks](https://github.com/fl1tenth/fl1tenth_racetracks), commit: b95c4ef.

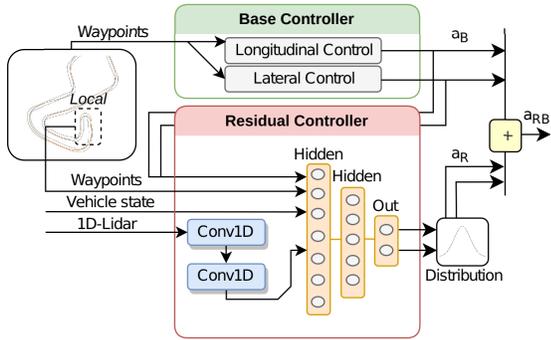


Fig. 2: Proposed control architecture consisting of a base controller that outputs control action  $a_B$  which is added to the residual action  $a_R$  to form the combined action  $a_{RB}$ . The action  $a_R$  of the residual controller is sampled from a distribution whose parameters are defined by a NN.

in 30 m distance are considered. The vehicle state with longitudinal and lateral velocity and acceleration, respectively, is concatenated to the yaw angle  $\psi$ , yaw rate  $\dot{\psi}$ , and the slip angle  $\beta$ . Additionally, the agent observes the action planned by the base controller  $a_{B,t}$  and the previous action that was applied  $a_{RB,t-1}$ . We stack the last three frames to introduce a limited memory to the agent; only the recent observation of  $L_t$  and  $w_{rel,t}$ : at  $t$  are used to reduce the dimension of the state. The reward function used for learning is defined as

$$r_{t+1} = \tau_1 \cdot v_{x,t} + \tau_2 \cdot v_{y,t}^2 + \rho \cdot \begin{cases} 1, & \text{if collision}_t = \text{True} \\ 0, & \text{otherwise} \end{cases}, \quad (9)$$

with  $\tau_1 = 0.003$  to encourage the agent to optimize for higher velocity while not increasing the lateral velocity, and therefore instability, too much due to  $\tau_2 = -0.003$ . Moreover,  $\rho = -50$  is set to prevent collisions.

### C. Network Design

In PPO, the policy  $\pi_\theta$  is defined as a NN with weights  $\theta$  along with a value network  $V_\phi$  and its parameters  $\phi$ . As shown in Fig. 2, our proposed architecture uses a perception module in form of two consecutive layers of 1D-convolution with 16 and 32 filters, respectively, with ReLu activation and average pooling. The perception module learns an embedding of the 1D-lidar data  $L_t$ ; its weights are shared between the policy and value network. In the policy and value network, respectively, the multilayer perceptrons (MLPs) consists of two hidden linear layers with  $\{400, 300\}$  neurons and ReLu activation and an output layer of fitting output dimension. The policy network learns the parameters of a TanhNormal distribution, i.e., a Gaussian distribution which is tanh transformed for projection to a bounded space  $[-1, 1]$ , with state-dependent mean  $\mu(s_t)$  but state-independent variance  $\sigma^2$ .

### D. Training Procedure

We use a customized version of PPO as implemented in [16]. The residual controller is trained for  $1e7$  episodes which correspond to approximately 2000 full race laps; the high-control frequency of the simulator requires training with long

trajectories of 2048 steps. Training trajectories are collected from 36 environments running in parallel. Observations and rewards are normalized by a running statistics calculation. While we set  $\gamma = 0.998$ , a batch size of 128, and stop epoch updates early for steps larger than a KL-divergence of 0.01, other used hyperparameters are unchanged from [16]. A scaling factor of  $[0.05, 1]^T$  is multiplied with the policy output to bound the residual  $a_{R,t}$  for the steering angle and the vehicle velocity by  $\pm 0.05$  and  $\pm 1.0$ , respectively. We use this setting for safety reasons to prevent an overly large influence of the residual controller on the base controller's behavior.

To evaluate the generalization capability of the proposed residual controller, training is only done on *nine* training racetracks (a-i) while the *three* test racetracks (j-l) are used for validation and are not seen during training. The environment resets after two full rounds on a racetrack are completed. Additional variability is introduced by randomly selecting a waypoint on the racing line every time the environment resets with a new starting position. While this makes the driving task more challenging, it is crucial for a meaningful evaluation of the controller's generalization capabilities.

## IV. EVALUATION

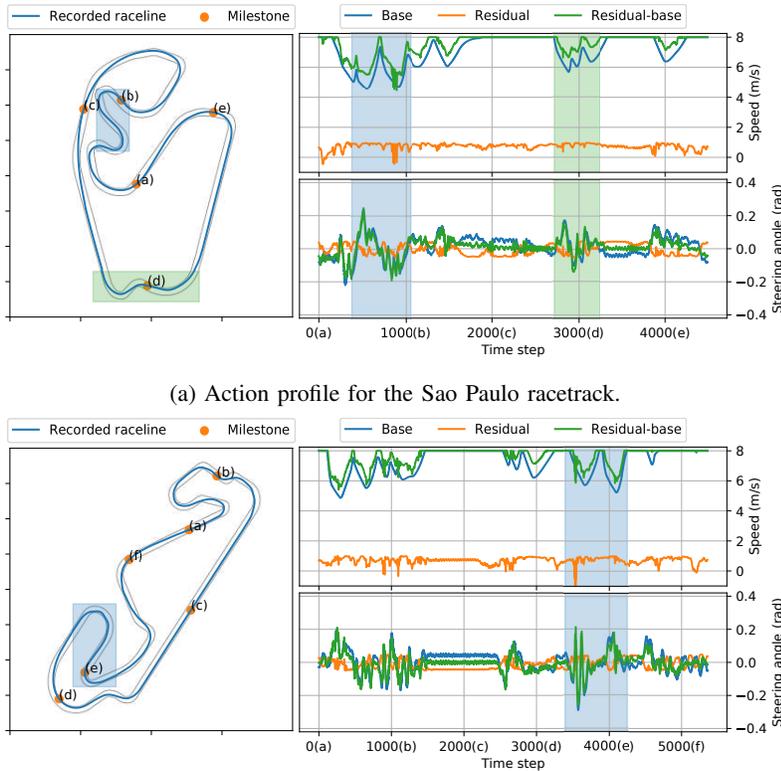
This section assesses the improvements linked to the use of our residual controller. For evaluation, the mean of the residual controller's stochastic policy  $\pi_\theta$  is taken instead of sampling from the learned distribution. We also introduce the *standard* controller as a baseline that uses the base controller's action  $a_{B,t}$  but without the added residual controller for comparison.

First, a high-level comparison of the lap time for all available racetracks is provided in Section IV-A. Reported lab times are the median results of two running start laps with random starting points over three independent full training runs, i.e., six lap time estimates overall, with changing seeds to estimate training stability. Secondly, Section IV-B presents an analysis of trajectories for a sub-set of racetracks for the best of the three trained controllers. Finally, Section IV-C compares the stability of this controller with the baseline.

### A. Lap Times

In this experiment, the residual controller's performance gains are quantified against the standard controller w.r.t. recorded lap times. Table II shows that lap times for the training set are directly improved by adding the residual component. The relative improvements range from 2.07 % to 7.06 % with an overall average improvement of 4.55 %. While improvements are expected for the training racetracks, similar gains can also be observed when the residual controller is deployed on the test set of racetracks. Lap times on the Sakhir, Catalunya, and Melbourne racetracks are improved by 4.34 %, 5.24 %, and 2.44 %, respectively. This evaluation demonstrates the capability of our proposed approach to generalize to unknown racetracks with comparable performance to the results observed on the training racetracks.

Looking closer at the characteristics of the racetracks and the achieved relative improvements, a correlation between



(a) Action profile for the Sao Paulo racetrack.

(b) Action profile for the Catalunya racetrack.

Fig. 3: Overview of the residual controller’s action profile split into base  $a_{B,t}$  (blue line), residual  $a_{R,t}$  (orange line), and residual-base  $a_{RB,t} = a_{R,t} + a_{B,t}$  (green line). Milestones (a) to (f) on the left inset correspond to the profiles (right insets) time steps (x-axis ticks). Green and blue areas provide a more precise overview of the vehicle behavior in the corresponding highlighted regions of the racetrack.

the racetrack curvature and the relative improvements can be observed. Racetracks with high curvatures, such as Moscow, Sao Paulo, and Hockenheim, form the top-3 improvements, whereas low curvature racetracks like Brand-Hatch tail the ranking. A similar observation can be made for the test set, with the Catalunya racetrack presenting a higher curvature than Melbourne and, thus, a higher lap time improvement.

### B. Residual Controller Behaviour

Given the hypothesis that the residual approach improves lap times when handling curves, we analyze the residual controller’s behavior profile in action-state space for one racetrack each from the training and test set.

Fig. 3 displays the vehicle’s recorded trajectory for the Sao Paulo and Catalunya racetracks w.r.t. the action profile for the residual controller, split into the action  $a_{B,t}$  of the base controller, the pure residual action  $a_{R,t}$  itself, and the combination of both as residual-base  $a_{RB,t}$ . Overall, it can be seen that the residual controller (i) has a strong tendency to increase the velocity; the action oscillates close to  $1 \frac{m}{s}$ , and (ii) only decreases the velocity in rare occurrences. Note that velocities above  $v^{\max} = 8 \frac{m}{s}$  are clipped due to the racing car’s

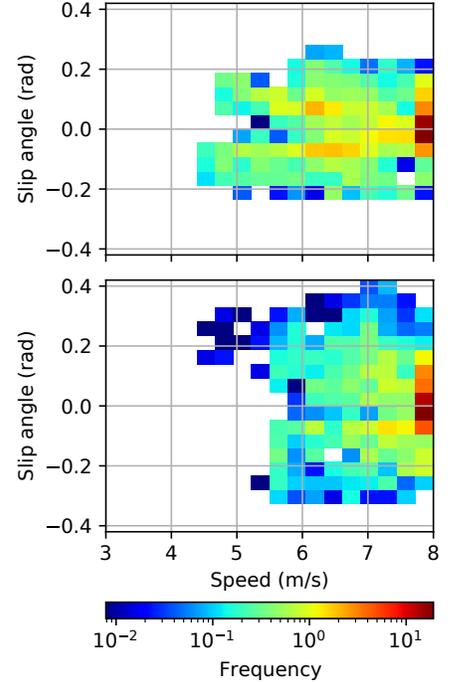


Fig. 4: Frequency of vehicle slip  $\beta$  aggregated over single laps on the Sao Paulo and Catalunya racetracks. The base controller (top inset) has a maximum  $|\beta|$  of 0.27 rad while the residual controller’s higher velocity (bottom) comes at the cost of a higher slip angle up to 0.43 rad.

TABLE II: Lap time results of the *standard* and the *residual* controllers for all racetracks as medians of three training runs.

	#	Racetrack Name	Lap time		Improvement	
			Standard	Residual	Diff.	Rel.
Training	a	Nürburgring	60.84 s	58.07 s	2.77 s	4.55 %
	b	Moscow	46.75 s	43.45 s	3.30 s	7.06 %
	c	Mexico City	49.12 s	46.76 s	2.36 s	4.80 %
	d	Brands Hatch	45.92 s	44.97 s	0.95 s	2.07 %
	e	Sao Paulo	47.92 s	44.92 s	3.00 s	6.26 %
	f	Sepang	66.24 s	63.18 s	3.06 s	4.62 %
	g	Hockenheim	49.96 s	47.35 s	2.61 s	5.22 %
	h	Budapest	54.33 s	51.67 s	2.66 s	4.90 %
	i	Spielberg	45.33 s	43.93 s	1.40 s	3.09 %
Test	j	Sakhir	60.34 s	57.72 s	2.62 s	4.34 %
	k	Catalunya	56.50 s	53.54 s	1.49 s	5.24 %
	l	Melbourne	61.03 s	59.54 s	2.96 s	2.44 %
Overall averages			53.69 s	50.85 s	2.43 s	4.55 %

physical limits, which occurs when the racing line’s planned velocity profile is already close to this limit.

For the Sao Paulo racetrack in Fig. 3a, it can be seen that even in the curvy section in the blue area of interest, the residual controller increases the velocity set by the base controllers

most times. Interestingly, the residual controller has learned to decelerate the car when approaching the curves' apex. This behavior yields more efficient steering behavior and a more aggressive racing line. Similar behavior is observed in the s-shaped section highlighted by the green area of interest, where the joint optimization of velocity and steering allows for taking dynamics effects into account. Furthermore, as can be seen from before milestone (c) and up to milestone (d), the residual controller has learned to comprehend the vehicle dynamics as it astutely counterbalances the base controller's steering to alter the trajectory enabling stable and high velocities.

Evaluation of the Catalunya racetrack in Fig. 3b supports these findings. As highlighted in the blue area of interest, the residual controller reduces the velocity before the curve's apex. Analysis of the straight section between milestones (b) and (d) validates that the residual controller has learned to adequately counter-steer at high velocities. Observing the same characteristics on a racetrack issued from the test set further confirms that the residual agent is capable of generalization.

### C. Vehicle Stability

Based on the analysis in Section IV-B, the observed improvements in lap times can be credited to higher velocity when taking turns. Doing so incurs more constraints on the racing car, potentially leading to reduced control quality and a detrimental deterioration of the car's stability. To evaluate the latter, we monitor the slip angle  $\beta$  of the kinematic single-track model with a linear relation to the tire forces. While a small slip angle is desirable w.r.t. stability, overall racing performance can only be improved by exploiting the car's dynamics adequately. Fig. 4 reports an aggregation of the experienced slip angles when using the residual controller. Compared to the standard controller as a baseline, it is clear that the residual controller experiences larger slip angles in the range of  $[-0.31, 0.43]$  rad compared to the standard controller's range of  $[-0.23, 0.27]$  rad. This observation confirms the intuition described above that higher speeds and short lap times are achieved at the expense of the car's stability. However, the fact that the residual controller can exploit the stability limit to its advantage without causing collisions demonstrates the added value of the residual component and its capability to adapt from known to unknown racetracks.

## V. CONCLUSION

We presented a new approach for vehicle control of autonomous racing cars using RPL. The proposed architecture consists of a *base* module for path-following of a racing line and a *residual* controller that learns a residual action from local observation to amend the base controller. Our study on vehicle behavior shows that the residual controller improves lap times by 4.55 % in average over a set of twelve replicated real-world racetracks. These improvements can be credited to adequate vehicle control in high curvature sections enabled by the residual controller. Moreover, the zero-shot generalization capability of our approach is demonstrated on unknown racetracks. Additionally to an ablation study of

the residual design, the envisioned future work will aim at evaluating the improvements of the residual controller for an extended set of scenarios. This includes the simulation of non-linear effects, e.g., drift, as well as a robustness analysis for noisy observations. Finally, we plan to explore the applicability of RPL to bridge the sim2real gap on real F1TENTH cars.

## ACKNOWLEDGEMENT

Marco Caccamo was supported by an Alexander von Humboldt Professorship endowed by the German Federal Ministry of Education and Research.

## REFERENCES

- [1] M. O'Kelly, H. Zheng, D. Karthik, and R. Mangharam, "F1tenth: An open-source evaluation environment for continuous control and reinforcement learning," in *NeurIPS 2019 Competition and Demonstration Track*. PMLR, 2020, pp. 77–89.
- [2] F. Fuchs, Y. Song, E. Kaufmann, D. Scaramuzza, and P. Dürri, "Super-human performance in gran turismo sport using deep reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4257–4264, 2021.
- [3] T. Johannink, S. Bahl, A. Nair, J. Luo, A. Kumar, M. Loskyll, J. A. Ojea, E. Solowjow, and S. Levine, "Residual reinforcement learning for robot control," in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 6023–6029.
- [4] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser, "Tossing-bot: Learning to throw arbitrary objects with residual physics," *IEEE Transactions on Robotics*, vol. 36, no. 4, pp. 1307–1319, 2020.
- [5] L. Kerbel, B. Ayalew, A. Ivanco, and K. Loisel, "Residual policy learning for powertrain control," *IFAC-PapersOnLine*, vol. 55, no. 24, pp. 111–116, 2022, 10th IFAC Symposium on Advances in Automotive Control AAC 2022.
- [6] A. Hynes, E. P. Sapozhnikova, and I. Dusparic, "Optimising PID control with residual policy reinforcement learning," in *Proceedings of The 28th Irish Conference on Artificial Intelligence and Cognitive Science, Dublin, Republic of Ireland, December 7-8, 2020*, ser. CEUR Workshop Proceedings, L. Longo, L. Rizzo, E. Hunter, and A. Pakrashi, Eds., vol. 2771. CEUR-WS.org, 2020, pp. 277–288.
- [7] R. Zhang, J. Hou, G. Chen, Z. Li, J. Chen, and A. Knoll, "Residual policy learning facilitates efficient model-free autonomous racing," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11 625–11 632, 2022.
- [8] A. Brunnbauer, L. Berducci, A. Brandstätter, M. Lechner, R. Hasani, D. Rus, and R. Grosu, "Latent imagination facilitates zero-shot transfer in autonomous racing," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 7513–7520.
- [9] G. Bellegarda and K. Byl, "An online training method for augmenting mpc with deep reinforcement learning," in *2020 IEEE/RSJ Intelligent Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020.
- [10] M. Althoff, M. Koschi, and S. Manzi, "Commonroad: Composable benchmarks for motion planning on roads," in *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017, pp. 719–726.
- [11] A. Heilmeyer, A. Wischniewski, L. Hermansdorfer, J. Betz, M. Lienkamp, and B. Lohmann, "Minimum curvature trajectory planning and control for an autonomous race car," *Vehicle System Dynamics*, 2019.
- [12] R. C. Coulter, "Implementation of the pure pursuit path tracking algorithm," Carnegie-Mellon UNIV Pittsburgh PA Robotics INST, Tech. Rep., 1992.
- [13] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint*, 2017.
- [14] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *arXiv preprint*, 2015.
- [15] J. Betz, H. Zheng, A. Liniger, U. Rosolia, P. Karle, M. Behl, V. Krovi, and R. Mangharam, "Autonomous vehicles on the edge: A survey on autonomous vehicle racing," *IEEE Open J. Intell. Transp. Syst.*, vol. 3, pp. 458–488, 2022.
- [16] S. Huang, R. F. J. Dossa, C. Ye, J. Braga, D. Chakraborty, K. Mehta, and J. G. Araújo, "Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms," *Journal of Machine Learning Research*, vol. 23, no. 274, pp. 1–18, 2022.