
Autonomous RL Racing with Strategic Decision Making

Souren Pashangpour
Department of Mechanical
And Industrial Engineering
University of Toronto
Toronto, ON M5S 1A1
souren.pashangpour@mail.utoronto.ca

Gibran A. Rajput
Department of Mechanical
And Industrial Engineering
University of Toronto
Toronto, ON M5S 1A1
gibran.rajput@mail.utoronto.ca

Zicheng Wang
Department of Mechanical
And Industrial Engineering
University of Toronto
Toronto, ON M5S 1A1
zchengw.wang@mail.utoronto.ca

Abstract

1 This paper presents a hierarchical reinforcement learning (HRL) framework
2 for autonomous racing that integrates macro-action planning with resource
3 management, focusing on tire wear and fuel consumption. Our model
4 utilizes a two-layer architecture: a high-layer for strategic decisions such as
5 pit stops, and a low-layer for direct vehicle control. By incorporating
6 typically omitted physical constraints, our framework improves the realism
7 and strategic depth of race simulations. The model's effectiveness is
8 demonstrated through superior performance and resource optimization in
9 simulated environments, compared to human benchmarks. Code is available
10 at: <https://github.com/ghssx19/RL-Project-Gym?tab=readme-ov-file>

11 1 Introduction

12 Autonomous racing and driving are a field where the majority of efforts are concentrated at the
13 intersection of deep learning and automobiles. Furthermore, Autonomous racing has led to
14 community projects such as F1Tenth, A2RL, AWS DeepRacer, and the Nvidia JetRacer Project.
15 The solutions proposed by its competitors in the aforementioned projects and challenges are
16 especially useful within game design, pre-race motorsport simulations, and motorsport real-time
17 decision-making support systems. However, previous work has overtly omitted factors such as
18 tyre wear, energy consumption, and physical restraints of a vehicle such as battery or breaking
19 system temperatures. The omission of these factors which we will refer to as physical constraints
20 leads to simplified decision making and omits strategic decision-making. Therefore, the developed
21 works are not well-suited for decision support systems or accurate motorsport simulations.
22 Moreover, research such as [1], which has addressed the constraint to some degree is difficult to
23 reproduce due to confidentiality limitations which make it difficult for the research community to
24 progress.

25 The choice of model to address the problem of autonomous racing has been primarily
26 reinforcement learning (RL), and Deep RL (DRL), where an agent in control of the car learns
27 directly from experiences collected in a simulation environment. The simulation environments
28 used to train and deploy these models are often proprietary such as AWS DeepRacer, or packages
29 such as Gazebo, OpenAI Gym, or Ansys. The omission of the physical constraints can be justified

30 thus far as simulation platforms lacked parallelization capabilities until the introduction of
31 IsaacLab and therefore it would lead to increased computational overhead to include the physical
32 constraints, while the problem of autonomous racing itself was relatively unexplored at the time.

33 Thus far only one paper has incorporated macro actions which are temporally extended actions
34 which are constituted of a series of primitive actions each last one single timestep[1]. Simply put
35 macro actions refer to strategic decisions such as deciding to take a pit stop, maintain current
36 position in race, or overtake in the race an agent can mitigate the effects of physical constraints to
37 grant itself an age over other competitors, during race time, where each of these actions consists of
38 low-level agent actions such as directional movement. This additional aspect of making high-level
39 strategic decisions is often referred to as macro action planners and have been applied in
40 Warehouse delivery [2].

41 In this paper, we introduce a Hierarchal RL-based multi-agent racing approach that accounts for
42 physical constraints, namely 1) tire wear, and 2) fuel consumption. Our unique contributions are
43 the development of a multi-agent Hierarchical RL project which can be expanded on by other
44 researchers with the objective of advancing macro action planning in the area of autonomous
45 racing called Hierarchical -RL Racing Model.

46 **2 Related work**

47 Existing autonomous racing and driving methods can be categorized as: 1) Hierarchical control
48 [3], [4], [5], [6], [7], [8], [9], or 2) single model control [10], [11], [12], [13], [14].

49 **2.1 Hierarchical RL for autonomous driving and racing**

50 Hierarchical RL methods for autonomous driving and racing separate the macro and micro actions
51 that an agent can take and delegate it to its respective deep learning model or algorithm. More
52 specifically hierarchical methods either use a higher level model to 1) generate way points for the
53 lower level model to drive to [4], [6], [7], 2) is used to identify the current state of the environment
54 and delegates the agent's actions to the expert model or algorithm developed for the agents current
55 state [5], [8], [9].

56 Furthermore, the method by which the high level RL model is embedded with the low level model
57 is either 1) interconnected, where the outputs from the higher level model are used as an output to
58 the lower level model which allows both components to learn simultaneously and increase the
59 adaptability of the framework [3], [5] . Or 2) where the inputs to the lower level model are strictly
60 from the environment observation space and there is no communication between the two models
61 [4], [6], [8], [9]. In the latter case a series of if/else statements are used to select the correct expert
62 model for the current state of the environment based on the high-level model's output.

63 **2.2 Single model control**

64 Single model RL methods used in autonomous racing are generally applied to the problem of
65 optimizing the racing line. In single model RL architectures, the reward functions are carefully
66 hand crafted to extract as much complex information from the environment as possible and
67 address a specific issue, for example minimizing time around the track given physical constraints
68 such as friction, maximum speed, maximum turning radius and etc.... The single model control
69 architectures are often then adopted in hierarchical RL approaches to serve as the expert models
70 for the different cases possible in the environment. Some of the common models trained and tuned
71 with hand crafted reward function for autonomous driving and racing include PPO [10], Soft
72 Actor Critic [10], twin delayed deep determinist policy gradient [11], Deep Deterministic Policy
73 Gradient [12], Deep Q-Network [13], and Probabilistic Inference for Learning Control [14]

74 **2.3 Summary of limitation**

75 In summary, there has been a significant amount of research done in applying models such as
 76 Hierarchical Program Triggered Reinforcement Learning (HPRL) [9] and Soft Twin Delayed
 77 Deep Deterministic Policy Gradient (Soft-TD3) [11] in applying single model control to different
 78 aspects of autonomous driving such as optimizing the path taken around a given track, or in road
 79 situations or maneuvers such as changing lanes or merging onto traffic. As regards to hierarchical
 80 approaches the majority of the work done is based around either 1) generating waypoints or areas
 81 which the lower level model should navigate to or 2) selecting an expert model best suited for a
 82 given environment's conditions, such as if the objective is to make a left turn the higher level
 83 model will apply the model trained for the task of making a left turn.

84 To the authors' knowledge hierarchical control for making strategic decisions has only been done
 85 once in a non-reproducible manner in [3], our aim with this project is to allow the community to
 86 have a maintained and working code base which they can use to address the problems introduced
 87 in the next section.

88 3 Hierarchical-RL racing methodology

89 In this section, we define the autonomous racing problem and present our Hierarchical DRL-based
 90 Racing model architecture. Figure 1 shows the overall model architecture.

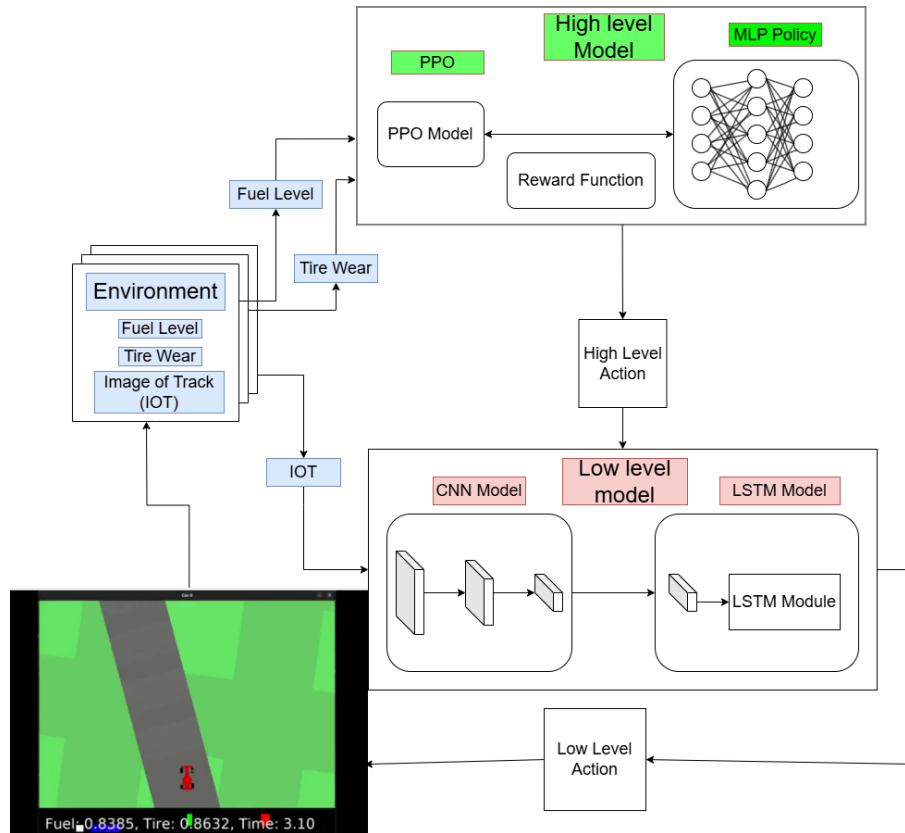


Figure 1: Model Architecture

91 3.1 Problem definition

92 The autonomous racing problem involves a vehicle V navigating a racetrack while considering
 93 key resource constraints: gas volume (v_g) and tire wear condition (t_w). The racetrack is
 94 represented as a continuous path within an environment characterized by sharp turns. Within the

95 simulation we can observe gas volume and tire wear per time step. The rate of consumption for
 96 these resources are dependent on factors such as periods of acceleration, deceleration, and the
 97 turning radius.

98 The objective is to minimize 1) the total amount of supplies (gas, tire thread) used, and 2) the lap
 99 time shown in Eq. (1). Therefore, we need to maximize the speed of the car at any given time
 100 around the track. The aforementioned recursion and the fact that a 15-second 1 pitstop can only be
 101 taken at a brief section of the track lead to the complexity of creating a model for predicting
 102 pitstops. In this paper, we aim to start preliminary work which can serve as a platform for future
 103 researchers to address the defined problem. Moreover, we will also focus on minimizing the
 104 frequency of pit stops, contributing to the overarching goal of optimizing race performance.

$$105 \quad \min T_{total} = \min \left[T_d + \sum^{n_{pit}} T_{pit} \right] \quad (1)$$

106 where T_{total} denotes total race time, T_d denotes total time spent driving on the track and T_{pit}
 107 denotes time spent in the pit stop, lastly n_{pit} denotes total number of pits stops. Specifically, this
 108 paper focuses on minimizing the frequency of pit stops, contributing to the overarching goal of
 109 optimizing race time.

110 **3.2 Autonomous racing architecture**

111 The proposed hierarchical RL racing architecture consists of two distinct layers: (1) a high-layer
 112 RL that handles strategic decision-making, and (2) a low-layer RL responsible for low-level
 113 vehicle control on the racetrack. This separation allows the architecture to efficiently divide
 114 complex tasks into manageable sub-tasks, where the high-layer RL focuses on task-specific
 115 optimization (Macro actions) [4], [15], and the low-layer RL specializes in directly interacting
 116 with the environment and producing control actions (Micro actions) [31]. The following sections
 117 describe the architecture’s key components in detail.

118 **3.2.1 Pit stop model**

119 The Pit Stop RL module forms the high-layer RL and is responsible for strategic decision-making
 120 regarding whether the vehicle should continue driving or take a pit stop for refueling or tire
 121 replacement. This decision is critical in optimizing the race performance while adhering to
 122 resource constraints, such as tire wear (t_w) and fuel volume (v_g). The module is implemented
 123 using Proximal Policy Optimization (PPO) [16], [17], [18]] with policy.

124 PPO is an advanced gradient method designed to improve the stability and efficiency of policy
 125 updates during training. It addresses the limitations of earlier methods such as Trust Region Policy
 126 Optimization (TRPO) [19] by using a clipped objective function to constrain policy changes within
 127 a trust region, thus preventing overly large updates that may destabilize training.

128 The clipped surrogate objective in PPO is defined in Eq. (2).

$$129 \quad L^{CLIP}(\theta) = E_t[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)] \quad (2)$$

130 where: $r_t(\theta)$ is the ratio of new to old policy probabilities, A_t is the estimated advantage
 131 function, ϵ is a clipping parameter (e.g., 0.2), which limits the extent of policy updates.

132 This clipped objective ensures that updates are neither too small (inefficient learning) nor too large
 133 (destabilizing), leading to a more robust and sample-efficient training process.

134 The advantage function A_t is computed in Eq. (3).

$$135 \quad A_t = \delta_t + (\gamma\lambda)\delta_{t+1} + (\gamma\lambda)^2\delta_{t+2} + \dots \quad (3)$$

136 Where:

$$137 \quad \delta_t = r_t + \gamma V(s_{t+1}) - V(s_t) \quad (4)$$

138 r_t is the reward at time t . γ is the discount factor, controlling how far into the future rewards are
139 considered. λ is the GAE (Generalized Advantage Estimation) decay parameter, balancing bias
140 and variance in the estimation. $V(s_t)$ is the value function estimate for state s_t . The advantage
141 function quantifies the relative quality of the chosen action a_t in state s_t , aiding the policy in
142 learning which actions are advantageous.

143 The MLP policy is trained to maximize this clipped objective. The multi-layer perceptron (MLP)
144 policy consists of: 1) Input Layer: Two state variables (v_g, t_w). 2) Hidden Layers: Fully connected
145 layers with non-linear activation functions (ReLU) [20], designed to capture complex state-action
146 relationships such as trade-off between tire wear and speed. 3) Output Layer: A SoftMax [21]
147 layer that outputs the probabilities of each high-level action (a_h). If $a_h = 1$, then pit stop is
148 needed. If $a_h = 0$, then continue driving. The policy is trained to maximize cumulative rewards,
149 where the reward function reflects race performance metrics, pit stop efficiency, and resource
150 utilization.

151 3.2.2 Vehicle control model

152 The low-layer RL module is responsible for low-level vehicle control, ensuring the vehicle
153 navigates the racetrack effectively and adheres to the high-layer RL’s strategic decisions. Unlike
154 the high-layer RL, which determines macro-level strategies, the low-layer RL focuses on
155 continuous control tasks such as steering, throttle, and braking to maintain the vehicle’s optimal
156 trajectory and speed. To improve efficiency and reduce training overhead, the low-layer RL
157 module leverages a pre-trained model. This model, based on PPO combined with a Long
158 Short-Term Memory (LSTM) architecture [22], is specifically designed for controlling vehicles in
159 continuous action spaces. The architecture combines the advantages of PPO’s robust policy
160 optimization with LSTM’s ability to capture temporal dependencies, making it well-suited for
161 tasks requiring sequential decision-making in dynamic environments.

162 The pre-trained model is composed of several key components. The first component is the
163 Observation Encoder, which processes the image-based observation space provided by the
164 environment. The input consists of $64 \times 64 \times 3$ RGB frames that capture a bird’s-eye view of the
165 racetrack and the vehicle’s surroundings. A convolutional neural network (CNN) [23] is used to
166 encode these images into compact, high-level feature representations, allowing the model to
167 extract meaningful spatial and contextual information essential for understanding the racetrack
168 layout and obstacles.

169 The encoded features are then passed through a Recurrent Layer implemented as a LSTM
170 network. This layer is designed to maintain a memory of past states, enabling the model to account
171 for temporal dependencies that are critical in vehicle control tasks. For instance, the LSTM allows
172 the model to anticipate upcoming turns or adjust speed in response to approaching obstacles. By
173 maintaining hidden states over time, the LSTM ensures that decisions are informed by both
174 current observations and historical context, improving the model’s ability to navigate dynamic
175 environments.

176 The Policy Network translates the output of the LSTM layer into actionable control commands for
177 the vehicle. Specifically, it maps the processed features to a continuous action space comprising
178 three control parameters: steering angle, throttle, and braking. The steering angle determines the
179 direction of the vehicle, while the throttle and braking control acceleration and deceleration,
180 respectively. These outputs are represented as probability distributions, which allow for stochastic
181 exploration during training and deterministic execution during deployment, balancing learning
182 efficiency and real-time performance.

183 Finally, the model includes a Value Function network, which estimates the expected return from
184 a given state. This component is essential for the PPO algorithm used in the pre-trained model. By
185 providing an estimate of future rewards, the value function reduces the variance in the reward
186 signal and stabilizes the policy-gradient optimization process. Together, these components enable
187 the pre-trained model to perform robust and efficient control in the continuous action space of
188 autonomous racing.

189 **3.2.3 Reward function**

190 The reward function is designed to guide the vehicle’s decision-making process, with the objective
 191 of minimizing pit stop frequency while ensuring optimal performance based on the vehicle’s
 192 resource conditions, specifically fuel volume (v_g) and tire wear (t_w). The rewards and penalties
 193 are carefully assigned to incentivize appropriate actions under different scenarios, as described in
 194 Eq. (5).

$$195 \quad r = \begin{cases} 50, & a_h = 1, v_g > 0.1, t_w > 0.1 \\ -1000, & a_h = 1, v_g \leq 0.1, t_w \leq 0.1 \\ 200, & a_h = 0, v_g \leq 0.1, t_w \leq 0.1 \\ -40, & a_h = 0, v_g > 0.5, t_w > 0.5 \end{cases} \quad (5)$$

196 $a_h = 1$ means pit stop action. $a_h = 0$ means continue driving action. For the pit stop action, the
 197 model is penalized with a reward of -40 when resources are in a healthy state ($v_g > 0.5$ and
 198 $t_w > 0.5$), discouraging unnecessary pit stops. Conversely, a reward of 200 is assigned if the pit stop
 199 action is taken under critical resource conditions ($v_g \leq 0.1$ and $t_w \leq 0.1$), encouraging the agent to
 200 take corrective measures when resources are nearly depleted.

201 For the continued driving action, the model receives a reward of 50 when resources are sufficient,
 202 promoting uninterrupted progress on the track. However, a substantial penalty of -1000 is applied
 203 if the model continues driving under critical conditions ($v_g \leq 0.1$ and $t_w \leq 0.1$), strongly
 204 discouraging risky behavior that could lead to failure.

Algorithm 1: Multi-Car Racing Simulation Algorithm

Input: *env*: MultiCarRacing environment,
low_level_model: Pre-trained low-level driving model,
pit_model: Pre-trained high-level pit decision model,
fuel_rate: Rate of fuel consumption,
tire_rate: Rate of tire wear,
max_steps: Maximum number of simulation steps,
FPS: Frames per second.
Output: *total_rewards*: Total rewards accumulated by the car

```

1 Initialization:
2 Set fuel_level  $\leftarrow$  1.0, tire_level  $\leftarrow$  1.0, total_rewards  $\leftarrow$  0.0,
   done  $\leftarrow$  False, step_counter  $\leftarrow$  0 ;
3 Reset the environment: obs  $\leftarrow$  env.reset() ;
4 while done = False and step_counter < max_steps do
5   | Predict low-level action: action_low  $\leftarrow$  low_level_model.predict(obs)
   | ;
6   | Update fuel and tire levels:
7   |   fuel_level  $\leftarrow$   $\max(\text{fuel\_level} - \frac{\text{fuel\_rate}}{\text{FPS}}, 0.0)$  ;
8   |   tire_level  $\leftarrow$   $\max(\text{tire\_level} - \frac{\text{tire\_rate}}{\text{FPS}}, 0.0)$  ;
9   | Formulate high-level observation: obs_high  $\leftarrow$  [fuel_level, tire_level]
   | ;
10  | Predict high-level action:
   |   action_high  $\leftarrow$  pit_model.predict(obs_high) ;
11  | if action_high = PIT then
   | | Refill fuel and tire levels: fuel_level, tire_level  $\leftarrow$  1.0 ;
   | | Step the environment:
   | |   obs, rewards, done  $\leftarrow$  env.step(action_low) ;
   | |   total_rewards  $\leftarrow$  total_rewards + rewards ;
14  | | else
   | | Step the environment without pitting:
   | |   obs, rewards, done  $\leftarrow$  env.step(action_low) ;
   | |   total_rewards  $\leftarrow$  total_rewards + rewards ;
17  | | end if
   | | Increment step counter: step_counter  $\leftarrow$  step_counter + 1 ;
18  | end while
19 return total_rewards

```

Figure 2: Pseudo code of the Hierarchical Model

205 **4 Training**

206 The Hierarchal-RL Racing model was trained in simulated racing environments using the PPO
207 algorithm. Each training environment represents a racetrack with dynamic elements such as
208 varying friction and obstacles. To simulate realistic race conditions, the environment incorporates
209 both sharp turns and straight sections, requiring the agent to learn diverse driving strategies. The
210 training framework uses a time-step-based episodic structure, where each episode ends when the
211 vehicle completes a lap or exhausts its resources.

212 The PPO model was configured with the following hyperparameters: a learning rate of $7 \times$
213 10^{-5} , $n_{steps} = 2048$, batch size of 64, and $n_{epochs} = 40$. A discount factor (γ) of 0.99 was used
214 to prioritize long-term rewards, while the Generalized Advantage Estimation (GAE) parameter
215 (λ) was set to 0.95 to reduce variance in the advantage function. The clipping range for PPO's
216 surrogate objective was 0.2, and the entropy coefficient was set to 0.01 to encourage policy
217 exploration. The value function coefficient was set to 0.5, and the gradient clipping threshold was
218 fixed at 0.5 to ensure stable updates during training.

219 The policy network utilized MLP architecture with customized policy keyword arguments
220 designed for efficient feature extraction. The MLP policy employed leaky ReLU as the activation
221 function for all layers except the output layer, ensuring smooth gradient propagation. The training
222 was performed over 250,000 episodes with an experience batch size of 16. The PPO agent was
223 trained using the cumulative reward as feedback, balancing strategic decision-making in the
224 high-layer RL with low-layer control optimization.

225 Training was executed over approximately 28 hours. The pseudo code for the model simulation
226 and training is in Figure 2. Figure 3,4 is the plot showing total reward per testing epoch, while
227 figure 4 is the accumulation of the rewards during training.



Figure 3: Reward per testing epoch



Figure 4: Accumulation of rewards during training

228 **5 Simulated experiments**

229 The simulated experiments include a comparison study to evaluate the performance of
230 Hierarchal-RL Racing model with human controlling.

231 **5.1 Comparison study**

232 We evaluated the performance of the Hierarchical RL racing model against human drivers based
233 on total race time (TRT). The goal of this comparison was to assess how effectively the proposed
234 model balances resource management and driving efficiency relative to human performance.

235 **5.1.1 Race environment**

236 The evaluation was conducted in a simulated racetrack environment designed to mimic real-world
237 racing conditions. The track included sharp turns, and straights, requiring adaptive strategies for
238 both speed optimization and resource management. The model and human participants evaluated
239 under identical initial resource states ($v_g=1.0$, $t_w=1.0$) and environmental settings. These
240 settings are the following: friction, driving backward, driving forward, on grass, steering, braking,
241 and speed. These settings are sent to the low-level model as an Image of track (IOT).

242 **5.1.2 Human participants**

243 One human participant controlled the vehicle using a standardized interface with keyboard inputs,
244 with feedback on fuel and tire conditions displayed in real time. Humans were instructed to
245 prioritize minimizing total race time while managing resources to complete the race successfully.

246 **5.1.3 Procedure**

247 Once the higher-level model was trained, it was imported into a separate file with the lower-level
248 model with an instance of an environment. The seed level was fixed so the generated track would
249 be the same. Another instance of the environment was created for Human with the same seed
250 value for the same track to be generated. Trials were run by their model and by the human for 10
251 laps.

252 **5.1.4 Results**

253 Table 1 presents a comparison of tabulated results between two TRT’s: the TRT of the model and
254 the TRT of the human performance.

Table 1: TRT comparison

Lap number	Model TRT (second)	Human TRT (second)
1st Lap	27.61	53.84
2nd Lap	59.36	77.18
3rd Lap	28.87	82.81
4th Lap	18,95	105.38
5th Lap	25.67	137.54
6th Lap	27.02	61.56
7th Lap	26.54	67.89

8th Lap	29.30	63.11
9th Lap	30.08	31.67
10th Lap	24.04	48.93

255 6 Discussions

256 This work introduces a novel hierarchical reinforcement learning architecture for autonomous
 257 racing, comprising a high-layer RL module for strategic decision-making and a low-layer RL
 258 module for precise vehicle control. The framework advances the state-of-the-art by incorporating
 259 macro-action planning and efficient resource management, addressing critical factors such as fuel
 260 consumption and tire wear that are often overlooked in prior research. The comparison study
 261 demonstrates the model’s ability to outperform human drivers in resource optimization and total
 262 race time, highlighting its effectiveness in minimizing unnecessary pit stops while maintaining
 263 race efficiency.

264 While the model does not yet achieve the overarching objective of fully optimizing all aspects of
 265 autonomous racing due to time constraints, the hierarchical architecture represents a significant
 266 step toward that goal. The ability to reduce pit stop frequency through strategic decision-making
 267 validates the potential of this approach in addressing real-world racing challenges.

268 Future work will concentrate on improving the adaptability of the high-layer RL module to handle
 269 a broader range of conditions and developing specialized low-layer RL modules tailored to
 270 dynamic scenarios. It is also important to extend our methodology to real-life racing scenarios to
 271 validate the model’s effectiveness under real-world conditions and further enhance its
 272 applicability. Efforts will also aim to bridge the interaction gap between the high-layer and
 273 low-layer RL modules, creating a more integrated learning framework where both layers adapt
 274 based on shared inputs and reciprocal feedback. By establishing a two-way communication
 275 between the top-level and low-level models, mutual learning can be enhanced, leading to greater
 276 control over the system’s overall performance. For example, directly linking race time to both
 277 layers allows the low-level model to refine its driving efficiency, while the high-level model
 278 concurrently learns to optimize pit stop strategies based on low-level driving behaviors. These
 279 enhancements will further evolve architecture, contributing to a more comprehensive and robust
 280 solution for autonomous racing challenges.

281 Acknowledgments

282 We extend our heartfelt gratitude to the individuals and communities whose contributions were
 283 pivotal to this project. Special thanks to [igilitschenski & tseyde](#) [24], for their foundational work
 284 on the environment, and to [sb3](#) for developing the core model that guided our research. We also
 285 acknowledge the dedicated developers and maintainers of the libraries, frameworks, and tools that
 286 empowered our project.

287 References

- 288 [1] C. Amato, G. D. Konidaris, and L. P. Kaelbling, “Planning with Macro-Actions in
 289 Decentralized POMDPs”.
- 290 [2] A. H. Tan, F. P. Bejarano, Y. Zhu, R. Ren, and G. Nejat, “Deep Reinforcement Learning for
 291 Decentralized Multi-Robot Exploration With Macro Actions,” Feb. 26, 2024, *arXiv*:
 292 arXiv:2110.02181. doi: 10.48550/arXiv.2110.02181.
- 293 [3] X. Liu, A. Fotouhi, and D. Auger, “Formula-E Multi-Car Race Strategy Development—A
 294 Novel Approach Using Reinforcement Learning,” *IEEE Trans. Intell. Transport. Syst.*, vol.
 295 25, no. 8, pp. 9524–9534, Aug. 2024, doi: 10.1109/TITS.2024.3389155.

- 296 [4] R. S. Thakkar, A. S. Samyal, D. Fridovich-Keil, Z. Xu, and U. Topcu, "Hierarchical Control
297 for Head-to-Head Autonomous Racing," *FR*, vol. 4, no. 1, pp. 46–69, Jan. 2024, doi:
298 10.55417/fr.2024002.
- 299 [5] Z. Qiao, Z. Tyree, P. Mudalige, J. Schneider, and J. M. Dolan, "Hierarchical Reinforcement
300 Learning Method for Autonomous Vehicle Behavior Planning," in *2020 IEEE/RSJ*
301 *International Conference on Intelligent Robots and Systems (IROS)*, Las Vegas, NV, USA:
302 IEEE, Oct. 2020, pp. 6084–6089. doi: 10.1109/IROS45743.2020.9341496.
- 303 [6] C. Chung, H. Seong, and D. H. Shim, "Learning from Demonstration with Hierarchical
304 Policy Abstractions Toward High-Performance and Courteous Autonomous Racing," Nov.
305 07, 2024, *arXiv*: arXiv:2411.04735. doi: 10.48550/arXiv.2411.04735.
- 306 [7] D. Kalaria, Q. Lin, and J. M. Dolan, "Towards Optimal Head-to-head Autonomous Racing
307 with Curriculum Reinforcement Learning," Aug. 25, 2023, *arXiv*: arXiv:2308.13491. doi:
308 10.48550/arXiv.2308.13491.
- 309 [8] J. Duan, S. Eben Li, Y. Guan, Q. Sun, and B. Cheng, "Hierarchical reinforcement learning
310 for self-driving decision-making without reliance on labelled driving data," *IET Intelligent*
311 *Trans Sys*, vol. 14, no. 5, pp. 297–305, May 2020, doi: 10.1049/iet-its.2019.0317.
- 312 [9] B. Gangopadhyay, H. Soora, and P. Dasgupta, "Hierarchical Program-Triggered
313 Reinforcement Learning Agents For Automated Driving," *IEEE Trans. Intell. Transport.*
314 *Syst.*, vol. 23, no. 8, pp. 10902–10911, Aug. 2022, doi: 10.1109/TITS.2021.3096998.
- 315 [10] B. Petryshyn, S. Postupaiev, S. Ben Bari, and A. Ostreika, "Deep Reinforcement Learning
316 for Autonomous Driving in Amazon Web Services DeepRacer," *Information*, vol. 15, no. 2,
317 p. 113, Feb. 2024, doi: 10.3390/info15020113.
- 318 [11] Z. Lu, C. Zhang, H. Zhang, Z. Wang, C. Huang, and Y. Ji, "Deep Reinforcement Learning
319 Based Autonomous Racing Car Control With Prior Knowledge," in *2021 China*
320 *Automation Congress (CAC)*, Beijing, China: IEEE, Oct. 2021, pp. 2241–2246. doi:
321 10.1109/CAC53003.2021.9728289.
- 322 [12] A. Remonda, S. Krebs, E. Veas, G. Luzhnica, and R. Kern, "Formula RL: Deep
323 Reinforcement Learning for Autonomous Racing using Telemetry Data," Jun. 13, 2022,
324 *arXiv*: arXiv:2104.11106. doi: 10.48550/arXiv.2104.11106.
- 325 [13] P. Aldape and S. Sowell, "Reinforcement Learning for a Simple Racing Game".
- 326 [14] Y. Zhu and D. Zhao, "Vision-based control in the open racing car simulator with deep and
327 reinforcement learning," *J Ambient Intell Human Comput*, vol. 14, no. 12, pp. 15673–
328 15685, Dec. 2023, doi: 10.1007/s12652-019-01503-y.
- 329 [15] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep Reinforcement
330 Learning: A Brief Survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38,
331 Nov. 2017, doi: 10.1109/MSP.2017.2743240.
- 332 [16] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy
333 Optimization Algorithms," Aug. 28, 2017, *arXiv*: arXiv:1707.06347. doi:
334 10.48550/arXiv.1707.06347.
- 335 [17] M. Hausknecht and P. Stone, "Deep Recurrent Q-Learning for Partially Observable MDPs".
- 336 [18] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol.
337 518, no. 7540, pp. 529–533, Feb. 2015, doi: 10.1038/nature14236.
- 338 [19] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust Region Policy
339 Optimization," Apr. 20, 2017, *arXiv*: arXiv:1502.05477. doi: 10.48550/arXiv.1502.05477.
- 340 [20] X. Glorot, A. Bordes, and Y. Bengio, "Deep Sparse Rectifier Neural Networks," in
341 *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*
342 *Statistics*, JMLR Workshop and Conference Proceedings, Jun. 2011, pp. 315–323.
343 Accessed: Nov. 24, 2024. [Online]. Available:
344 <https://proceedings.mlr.press/v15/glorot11a.html>
- 345 [21] I. Kouretas and V. Paliouras, "Simplified Hardware Implementation of the Softmax
346 Activation Function," in *2019 8th International Conference on Modern Circuits and*
347 *Systems Technologies (MOCASST)*, May 2019, pp. 1–4. doi:
348 10.1109/MOCASST.2019.8741677.
- 349 [22] S. Hochreiter and J. Schmidhuber, "Long Short-term Memory," *Neural computation*, vol. 9,
350 pp. 1735–80, Dec. 1997, doi: 10.1162/neco.1997.9.8.1735.

- 351 [23] Y. LeCun, L. Bottou, Y. Bengio, and P. Ha, "Gradient-Based Learning Applied to
352 Document Recognition," 1998.
- 353 [24] W. Swarting *et al.*, "Deep Latent Competition: Learning to Race Using Visual Control
354 Policies in Latent Space," Feb. 19, 2021, *arXiv*: arXiv:2102.09812. doi:
355 10.48550/arXiv.2102.09812.