

Outracing champion Gran Turismo drivers with deep reinforcement learning

<https://doi.org/10.1038/s41586-021-04357-7>

Received: 9 August 2021

Accepted: 15 December 2021

Published online: 9 February 2022

 Check for updates

Peter R. Wurman¹✉, Samuel Barrett¹, Kenta Kawamoto², James MacGlashan¹, Kaushik Subramanian³, Thomas J. Walsh¹, Roberto Capobianco³, Alisa Devlic³, Franziska Eckert³, Florian Fuchs³, Leilani Gilpin¹, Piyush Khandelwal¹, Varun Kompella¹, HaoChih Lin³, Patrick MacAlpine¹, Declan Oller¹, Takuma Seno², Craig Sherstan¹, Michael D. Thomure¹, Houmehar Aghabozorgi¹, Leon Barrett¹, Rory Douglas¹, Dion Whitehead¹, Peter Dürr³, Peter Stone¹, Michael Spranger² & Hiroaki Kitano²

Many potential applications of artificial intelligence involve making real-time decisions in physical systems while interacting with humans. Automobile racing represents an extreme example of these conditions; drivers must execute complex tactical manoeuvres to pass or block opponents while operating their vehicles at their traction limits¹. Racing simulations, such as the PlayStation game Gran Turismo, faithfully reproduce the non-linear control challenges of real race cars while also encapsulating the complex multi-agent interactions. Here we describe how we trained agents for Gran Turismo that can compete with the world's best e-sports drivers. We combine state-of-the-art, model-free, deep reinforcement learning algorithms with mixed-scenario training to learn an integrated control policy that combines exceptional speed with impressive tactics. In addition, we construct a reward function that enables the agent to be competitive while adhering to racing's important, but under-specified, sportsmanship rules. We demonstrate the capabilities of our agent, Gran Turismo Sophy, by winning a head-to-head competition against four of the world's best Gran Turismo drivers. By describing how we trained championship-level racers, we demonstrate the possibilities and challenges of using these techniques to control complex dynamical systems in domains where agents must respect imprecisely defined human norms.

Deep reinforcement learning (deep RL) has been a key component of impressive recent artificial intelligence (AI) milestones in domains such as Atari², Go^{3,4}, StarCraft⁵ and Dota⁶. For deep RL to have an influence on robotics and automation, researchers must demonstrate success in controlling complex physical systems. In addition, many potential applications of robotics require interacting in close proximity to humans while respecting imprecisely specified human norms. Automobile racing is a domain that poses exactly these challenges; it requires real-time control of vehicles with complex, non-linear dynamics while operating within inches of opponents. Fortunately, it is also a domain for which highly realistic simulations exist, making it amenable to experimentation with machine-learning approaches.

Research on autonomous racing has accelerated in recent years, leveraging full-sized^{7–10}, scale^{11–15} and simulated^{16–25} vehicles. A common approach pre-computes trajectories^{26,27} and uses model predictive control to execute those trajectories^{7,28}. However, when driving at the absolute limits of friction, small modelling errors can be catastrophic. Racing against other drivers puts even greater demands on modelling accuracy, introduces complex aerodynamic interactions and further requires engineers to design control schemes that continuously predict and adapt to the trajectories of other

cars. Racing with real driverless vehicles still seems to be several years away, as the recent Indy Autonomous Challenge curtailed its planned head-to-head competition to time trials and simple obstacle avoidance²⁹.

Researchers have explored various ways to use machine learning to avoid this modelling complexity, including using supervised learning to model vehicle dynamics^{8,12,30} and using imitation learning³¹, evolutionary approaches³² or reinforcement learning^{16,21} to learn driving policies. Although some studies achieved super-human performance in solo driving²⁴ or progressed to simple passing scenarios^{16,20,25,33}, none have tackled racing at the highest levels.

To be successful, racers must become highly skilled in four areas: (1) race-car control, (2) racing tactics, (3) racing etiquette and (4) racing strategy. To control the car, drivers develop a detailed understanding of the dynamics of their vehicle and the idiosyncrasies of the track on which they are racing. On this foundation, drivers build the tactical skills needed to pass and defend against opponents, executing precise manoeuvres at high speed with little margin for error. At the same time, drivers must conform to highly refined, but imprecisely specified, sportsmanship rules. Finally, drivers use strategic thinking when modelling opponents and deciding when and how to attempt a pass.

¹Sony AI, New York, NY, USA. ²Sony AI, Tokyo, Japan. ³Sony AI, Zürich, Switzerland. ✉e-mail: peter.wurman@sony.com

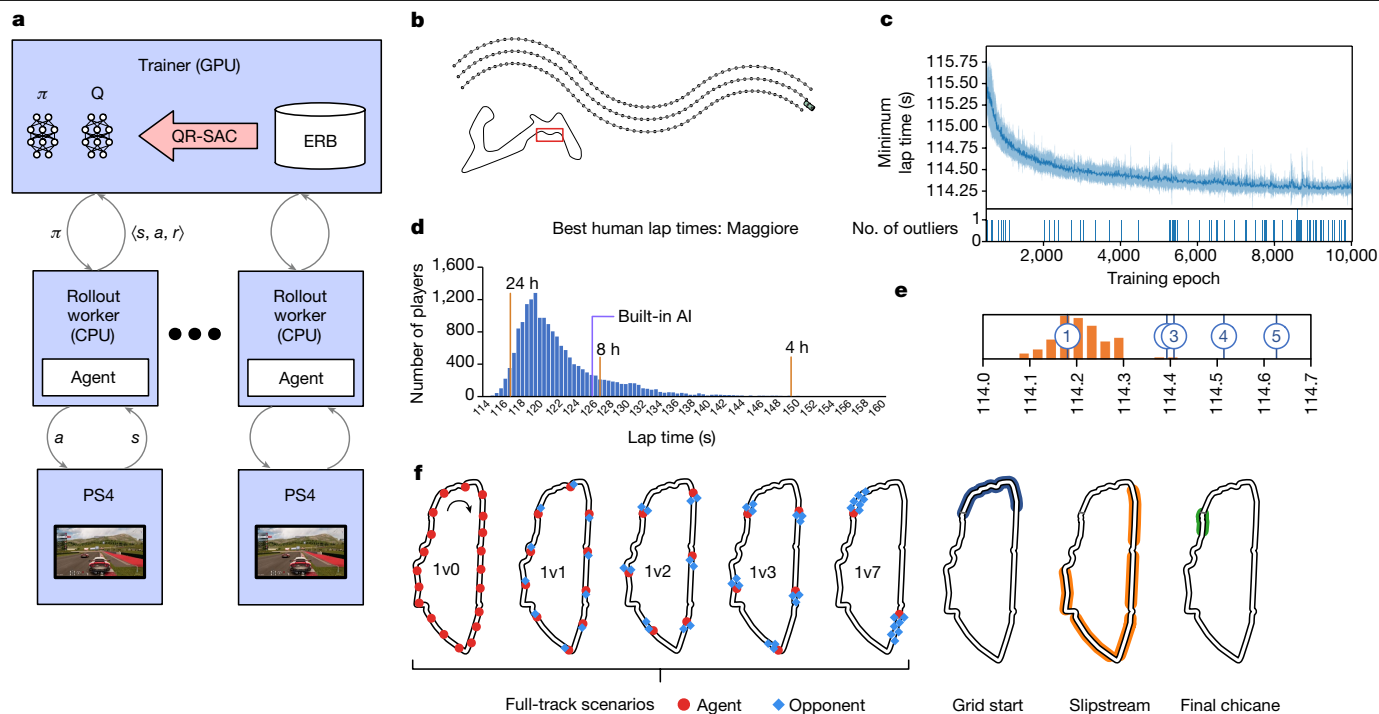


Fig. 1 | Training. **a**, An example training configuration. The trainer distributes training scenarios to rollout workers, each of which controls one PS4 running an instance of GT. The agent in the worker runs one copy of the most recent policy, π , to control up to 20 cars. The agent sends an action, a , for each car it controls to the game. Asynchronously, the game computes the next frames and sends each new state, s , to the agent. When the game reports that the action has been registered, the agent reports the state, action, reward tuple $\langle s, a, r \rangle$ to the trainer, which stores it in the ERB. The trainer samples the ERB to update the policy, π , and Q-function networks. **b**, The course representation ahead of the car on a sequence of curves on Maggiore if the car was travelling at 200 km h⁻¹. **c**, The distribution of the learning curves on Maggiore from 15 different random seeds. All of the seeds reached superhuman performance. Most reached it in 10 days of training, whereas the longest took 25 days. **d**, The distribution of individual players' best lap times on Maggiore as recorded on

Kudos Prime (<https://www.kudosprime.com/gts/rankings.php?sec=daily>). Superimposed on **d** is the number of hours that GT Sophy, using ten PlayStation 4s with 20 cars each, needed to achieve similar performance. **e**, A histogram (in orange) of 100 laps from the time-trial policy GT Sophy used on 2 July 2021 compared with the five best human drivers' best lap times (circles 1–5) in the Kudos Prime data. Similar graphs for the other two tracks are in the Supplementary Information; Maggiore is the only one of the three tracks on which the best human performance was close to GT Sophy. **f**, The training scenarios on Sarthe, including five full-track configurations in which the agent starts with zero, one, two, three or seven nearby opponents and three specialized scenarios that are limited to the shaded regions. The actual track positions, opponents and relative car arrangements are varied to ensure that the learned skills are robust.

In this article, we describe how we used model-free, off-policy deep RL to build a champion-level racing agent, which we call Gran Turismo Sophy (GT Sophy). GT Sophy was developed to compete with the world's best players of the highly realistic PlayStation 4 (PS4) game Gran Turismo (GT) Sport (<https://www.gran-turismo.com/us/>), developed by Polyphony Digital, Inc. We demonstrate GT Sophy by competing against top human drivers on three car and track combinations that posed different racing challenges. The car used on the first track, Dragon Trail Seaside (Seaside), was a high-performance road vehicle. On the second track, Lago Maggiore GP (Maggiore), the vehicle was equivalent to the Federation Internationale de l'Automobile (FIA) GT3 class of race cars. The third and final race took place on the Circuit de la Sarthe (Sarthe), famous as the home of the 24 Hours of Le Mans. This race featured the Red Bull X2019 Competition race car, which can reach speeds in excess of 300 km h⁻¹. Although lacking strategic savvy, in the process of winning the races against humans, GT Sophy demonstrated notable advances in the first three of the four skill areas mentioned above.

Approach

The training configuration is illustrated in Fig. 1a. GT runs only on PlayStation 4s, which necessitated that the agent runs on a separate computing device and communicates asynchronously with the game

by means of TCP. Although GT ran only in real time, each GT Sophy instance controlled up to 20 cars on its PlayStation 4, which accelerated data collection. We typically trained GT Sophy from scratch using 10–20 PlayStation 4s, an equal number of compute instances and a GPU machine that asynchronously updates the neural networks.

The core actions of the agent were mapped to two continuous-valued dimensions: changing velocity (accelerating or braking) and steering (left or right). The effect of the actions was enforced by the game to be consistent with the physics of the environment; GT Sophy cannot brake harder than humans but it can learn more precisely when to brake. GT Sophy interacted with the game at 10 Hz, which we claim does not give GT Sophy a particular advantage over professional gamers³⁴ or athletes³⁵.

As is common^{26,27}, the agent was given a static map defining the left and right edges and the centre line of the track. We encoded the approaching course segment as 60 equally spaced 3D points along each edge of the track and the centre line (Fig. 1b). The span of the points in any given observation was a function of the current velocity, so as to always represent approximately the next 6 s of travel. The points were computed from the track map and presented to the neural network in the egocentric frame of reference of the agent.

Through an API, GT Sophy observed the positions, velocities, accelerations and other relevant state information about itself and all opponents. To make opponent information amenable to deep learning, GT

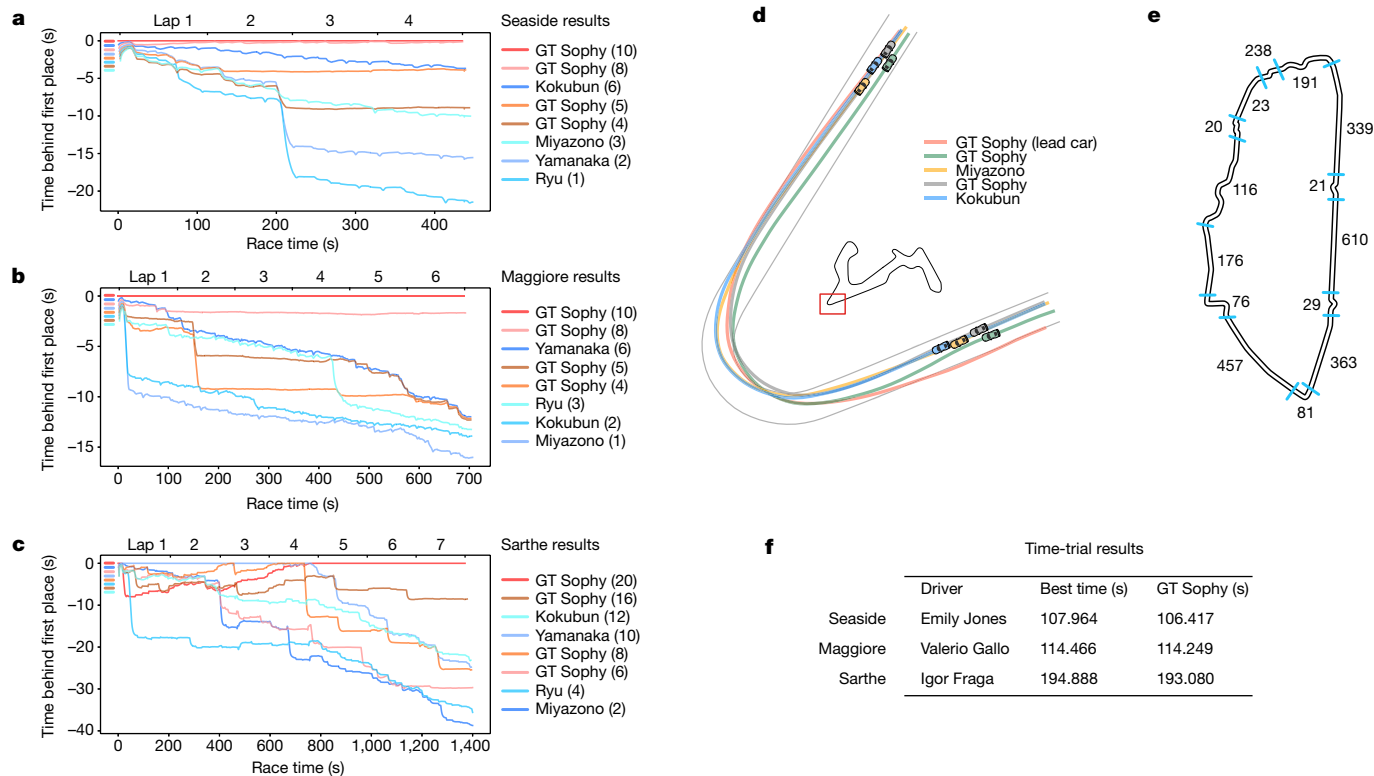


Fig. 3 | Results. a–c, How each race unfolded on Seaside (a), Maggiore (b) and Sarthe (c). The distance from the leader is computed as the time since the lead car passed the same position on the track. The legend for each race shows the final places and, in parentheses, the points for each driver. These charts clearly show how, once GT Sophy obtained a small lead, the human drivers could not catch it. The sharp decreases represent either a driver losing control or paying a penalty for either exceeding the course bounds or colliding with another driver. Sarthe (c) had the most incidents, with GT Sophy receiving two penalties for excessive contact and the humans receiving one penalty and two warnings. Both the humans and GT Sophy also had several smaller penalties for exceeding the course boundaries, particularly in the final chicane sequence. **d, An**

example from the 2 July 2021 race in which two instances of GT Sophy (grey and green) passed two humans (yellow and blue) on a corner at Maggiore. As a reference, the trajectory of the lead GT Sophy car when taking the corner alone is shown in red. The example clearly illustrates that the trajectory of GT Sophy through the corner is contextual; even though the human drivers tried to protect the inside going into the corner, GT Sophy was able to find two different, faster trajectories. **e,** The number of passes that occurred on different parts of Sarthe in 100 4v4 races between two GT Sophy policies, demonstrating that the agent has learned to pass on many parts of the track. **f,** Results from the time-trial competition in July 2021.

assignment, the resulting policies were judged much too aggressive by stewards and test drivers. For the final races, we opted for a conservative approach that penalized the agent for any collision in which it was involved (regardless of fault), with some extra penalties if the collision was likely considered unacceptable. Figure 2a–h isolates the effects of collision penalties and other key design choices made during this project.

Although many applications of RL to games use self-play to improve performance^{3,42}, the straightforward application of self-play was inadequate in this setting. For example, as a human enters a difficult corner, they may brake a fraction of a second earlier than the agent would. Even a small bump at the wrong moment can cause an opponent to lose control of their car. By racing against only copies of itself, the agent was ill-prepared for the imprecision it would see with human opponents. If the agent following does not anticipate the possibility of the opponent braking early, it will not be able to avoid rear-ending the human driver and will be assessed a penalty. This feature of racing—that one player’s suboptimal choice causes the other player to be penalized—is not a feature of zero-sum games such as Go and chess. To alleviate this issue, we used a mixed population of opponents, including agents curated from previous experiments and the game’s (relatively slower) built-in AI. Figure 2e shows the importance of these choices.

In addition, the opportunities to learn certain skills are rare. We call this the exposure problem; certain states of the world are not accessible to the agent without the ‘cooperation’ of its opponents. For example, to

execute a slingshot pass, a car must be in the slipstream of an opponent on a long straightaway, a condition that may occur naturally a few times or not at all in an entire race. If that opponent always drives only on the right, the agent will learn to pass only on the left and would be easily foiled by a human who chose to drive on the left. To address this issue, we developed a process that we called mixed-scenario training. We worked with a retired competitive GT driver to identify a small number of race situations that were probably pivotal on each track. We then configured scenarios that presented the agent with noisy variations of those critical situations. In slipstream passing scenarios, we used simple PID controllers to ensure that the opponents followed certain trajectories, such as driving on the left, that we wanted our agent to be prepared for. Figure 1f shows the full-track and specialized scenarios for Sarthe. Notably, all scenarios were present throughout the training regime; no sequential curriculum was needed. We used a form of stratified sampling⁴³ to ensure that situational diversity was present throughout training. Figure 2h shows that this technique resulted in more robust skills being learned.

Evaluation

To evaluate GT Sophy, we raced the agent in two events against top GT drivers. The first event was on 2 July 2021 and involved both time-trial and head-to-head races. In the time-trial race, three of the world’s top

drivers were asked to try to beat the lap times of GT Sophy. Although the human drivers were allowed to see a ‘ghost’ of GT Sophy as they drove around the track, GT Sophy won all three matches. The results are shown in Fig. 3f.

The head-to-head race was held at the headquarters of Polyphony Digital and, although limited to top Japanese players due to pandemic travel restrictions, included four of the world’s best GT drivers. These drivers formed a team to compete against four instances of GT Sophy. Points were awarded to the team on the basis of the final positions (10, 8, 6, 5, 4, 3, 2 and 1, from first to last), with Sarthe, the final and most challenging race, counting double. Each team started in either the odd or the even positions on the basis of their best qualifying time. The human drivers won the team event on 2 July 2021 by a score of 86–70.

After examining GT Sophy’s 2 July 2021 performance, we improved the training regime, increased the network size, made small modifications to some features and rewards and improved the population of opponents. GT Sophy handily won the rematch held on 21 October 2021 by an overall team score of 104–52. Starting in the odd positions, team GT Sophy improved four spots on Seaside and Maggiore and two on Sarthe. Figure 3a–c shows the relative positions of the cars through each race and the points earned by each individual.

One of the advantages of using deep RL to develop a racing agent is that it eliminates the need for engineers to program how and when to execute the skills needed to win the race—as long as it is exposed to the right conditions, the agent learns to do the right thing by trial and error. We observed that GT Sophy was able to perform several types of corner passing, use the slipstream effectively, disrupt the draft of a following car, block, and execute emergency manoeuvres. Figure 3d shows particularly compelling evidence of GT Sophy’s generalized tactical competence. The diagram illustrates a situation from the 2 July 2021 event in which two GT Sophy cars both pass two human cars on a single corner on Maggiore. This kind of tactical competence was not limited to any particular part of the course. Figure 3e shows the number of passes that occurred on different sections of Sarthe from 100 4v4 races between two different GT Sophy policies. Although slipstream passing on the straightaways was most common, the results show that GT Sophy was able to take advantage of passing opportunities on many different sections of Sarthe.

Although GT Sophy demonstrated enough tactical skill to beat expert humans in head-to-head racing, there are many areas for improvement, particularly in strategic decision-making. For example, GT Sophy takes the first opportunity to pass on a straightaway, sometimes leaving enough room on the same stretch of track for the opponent to use the slipstream to pass back. GT Sophy also aggressively tries to pass an opponent with a looming penalty, whereas a strategic human driver may wait and make the easy pass when the opponent is forced to slow down.

Conclusions

Simulated automobile racing is a domain that requires real-time, continuous control in an environment with highly realistic, complex physics. The success of GT Sophy in this environment shows, for the first time, that it is possible to train AI agents that are better than the top human racers across a range of car and track types. This result can be seen as another important step in the continuing progression of competitive tasks that computers can beat the very best people at, such as chess, Go, Jeopardy, poker and StarCraft. In the context of previous landmarks of this kind, GT Sophy is the first that deals with head-to-head, competitive, high-speed racing, which requires advanced tactics and subtle sportsmanship considerations. Agents such as GT Sophy have the potential to make racing games more enjoyable, provide realistic, high-level competition for training professional drivers and discover new racing techniques. The success of deep RL in this environment suggests that

these techniques may soon have an effect on real-world systems such as collaborative robotics, aerial drones or autonomous vehicles.

All references to Gran Turismo, PlayStation and other Sony properties are made with permission of the respective rights owners. Gran Turismo Sport: © 2019 Sony Interactive Entertainment Inc. Developed by Polyphony Digital Inc. Manufacturers, cars, names, brands and associated imagery featured in this game in some cases include trademarks and/or copyrighted materials of their respective owners. All rights reserved. Any depiction or recreation of real world locations, entities, businesses, or organizations is not intended to be or imply any sponsorship or endorsement of this game by such party or parties. “Gran Turismo” and “Gran Turismo Sophy” logos are registered trademarks or trademarks of Sony Interactive Entertainment Inc.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-021-04357-7>.

1. Milliken, W. F. et al. *Race Car Vehicle Dynamics* Vol. 400 (Society of Automotive Engineers, 1995).
2. Mnih, V. et al. Playing Atari with deep reinforcement learning. Preprint at <https://arxiv.org/abs/1312.5602> (2013).
3. Silver, D. et al. Mastering the game of Go with deep neural networks and tree search. *Nature* **529**, 484–489 (2016).
4. Silver, D. et al. Mastering the game of Go without human knowledge. *Nature* **550**, 354–359 (2017).
5. Vinyals, O. et al. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* **575**, 350–354 (2019).
6. Berner, C. et al. Dota 2 with large scale deep reinforcement learning. Preprint at <https://arxiv.org/abs/1912.06680> (2019).
7. Laurence, V. A., Goh, J. Y. & Gerdes, J. C. In *2017 American Control Conference (ACC)* 5586–5591 (IEEE, 2017).
8. Spielberg, N. A., Brown, M., Kapania, N. R., Kegelman, J. C. & Gerdes, J. C. Neural network vehicle models for high-performance automated driving. *Sci. Robot.* **4**, eaaw1975 (2019).
9. Burke, K. Data makes it beta: Roborace returns for second season with updateable self-driving vehicles powered by NVIDIA DRIVE. *The Official NVIDIA Blog* <https://blogs.nvidia.com/blog/2020/10/29/roborace-second-season-nvidia-drive/> (2020).
10. Leporati, G. No driver? no problem—this is the Indy Autonomous Challenge. *Ars Technica* <https://arstechnica.com/cars/2021/07/a-science-fair-or-the-future-of-racing-the-indy-autonomous-challenge/> (2021).
11. Williams, G., Drews, P., Goldfain, B., Reh, J. M. & Theodorou, E. A. In *2016 IEEE International Conference on Robotics and Automation (ICRA)* 1433–1440 (IEEE, 2016).
12. Williams, G., Drews, P., Goldfain, B., Reh, J. M. & Theodorou, E. A. Information-theoretic model predictive control: theory and applications to autonomous driving. *IEEE Trans. Robot.* **34**, 1603–1622 (2018).
13. Pan, Y. et al. In *Proc. Robotics: Science and Systems XIV* (eds Kress-Gazit, H., Srinivasa, S., Howard, T. & Atanasov, N.) <https://doi.org/10.15607/RSS.2018.XIV.056> (Carnegie Mellon Univ., 2018).
14. Pan, Y. et al. Imitation learning for agile autonomous driving. *Int. J. Robot. Res.* **39**, 286–302 (2020).
15. Amazon Web Services. AWS DeepRacer League. <https://aws.amazon.com/deepracer/league/> (2019).
16. Pyeatt, L. D. & Howe, A. E. Learning to race: experiments with a simulated race car. In *Proc. Eleventh International FLAIRS Conference* 357–361 (AAA, 1998).
17. Chaperot, B. & Fyfe, C. In *2006 IEEE Symposium on Computational Intelligence and Games* 181–186 (IEEE, 2006).
18. Cardamone, L., Loiacono, D. & Lanzi, P. L. In *Proc. 11th Annual Conference on Genetic and Evolutionary Computation* 1179–1186 (ACM, 2009).
19. Cardamone, L., Loiacono, D. & Lanzi, P. L. In *2009 IEEE Congress on Evolutionary Computation* 2622–2629 (IEEE, 2009).
20. Loiacono, D., Prete, A., Lanzi, L. & Cardamone, L. In *IEEE Congress on Evolutionary Computation* 1–8 (IEEE, 2010).
21. Jaritz, M., de Charette, R., Toromanoff, M., Perot, E. & Nashashibi, F. In *2018 IEEE International Conference on Robotics and Automation (ICRA)* 2070–2075 (IEEE, 2018).
22. Weiss, T. & Behl, M. In *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)* 1163–1168 (IEEE, 2020).
23. Weiss, T., Babu, V. S. & Behl, M. In *NeurIPS 2020 Workshop on Machine Learning for Autonomous Driving* (NeurIPS, 2020).
24. Fuchs, F., Song, Y., Kaufmann, E., Scaramuzza, D. & Dürr, P. Super-human performance in Gran Turismo Sport using deep reinforcement learning. *IEEE Robot. Autom. Lett.* **6**, 4257–4264 (2021).
25. Song, Y., Lin, H., Kaufmann, E., Dürr, P. & Scaramuzza, D. In *Proc. IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2021).
26. Theodosios, P. A. & Gerdes, J. C. In *Dynamic Systems and Control Conference* Vol. 45295, 235–241 (American Society of Mechanical Engineers, 2012).
27. Funke, J. et al. In *2012 IEEE Intelligent Vehicles Symposium* 541–547 (IEEE, 2012).

28. Kritayakirana, K. & Gerdes, J. C. Autonomous vehicle control at the limits of handling. *Int. J. Veh. Auton. Syst.* **10**, 271–296 (2012).
29. Bonkowski, J. Here's what you missed from the Indy Autonomous Challenge main event. *Autoweek* <https://www.autoweek.com/racing/more-racing/a38069263/what-missed-indy-autonomous-challenge-main-event/> (2021).
30. Rutherford, S. J. & Cole, D. J. Modelling nonlinear vehicle dynamics with neural networks. *Int. J. Veh. Des.* **53**, 260–287 (2010).
31. Pomerleau, D. A. In *Robot Learning* (eds Connell, J. H. & Mahadevan, S.) 19–43 (Springer, 1993).
32. Togelius, J. & Lucas, S. M. In *2006 IEEE International Conference on Evolutionary Computation* 1187–1194 (IEEE, 2006).
33. Schwarting, W. et al. Deep latent competition: learning to race using visual control policies in latent space. Preprint at <https://arxiv.org/abs/2102.09812> (2021).
34. Gozli, D. G., Bavelier, D. & Pratt, J. The effect of action video game playing on sensorimotor learning: evidence from a movement tracking task. *Hum. Mov. Sci.* **38**, 152–162 (2014).
35. Davids, K., Williams, A. M. & Williams, J. G. *Visual Perception and Action in Sport* (Routledge, 2005).
36. Haarnoja, T., Zhou, A., Abbeel, P. & Levine, S. In *Proc. 35th International Conference on Machine Learning* 1856–1865 (PMLR, 2018).
37. Haarnoja, T. et al. Soft actor-critic algorithms and applications. Preprint at <https://arxiv.org/abs/1812.05905> (2018).
38. Mnih, V. et al. In *Proc. 33rd International Conference on Machine Learning* 1928–1937 (PMLR, 2016).
39. Dabney, W., Rowland, M., Bellemare, M. G. & Munos, R. In *32nd AAAI Conference on Artificial Intelligence* (AAAI, 2018).
40. Lin, L.-J. *Reinforcement Learning for Robots Using Neural Networks*. Dissertation, Carnegie Mellon Univ. (1993).
41. Siu, H. C. et al. Evaluation of human-AI teams for learned and rule-based agents in Hanabi. Preprint at <https://arxiv.org/abs/2107.07630> (2021).
42. Tesauro, G. TD-Gammon, a self-teaching backgammon program, achieves master-level play. *Neural Comput.* **6**, 215–219 (1994).
43. Devore, J. L. *Probability and Statistics for Engineering and the Sciences* 6th edn (Brooks/Cole, 2004).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature Limited 2022

Methods

Game environment

Since its debut in 1997, the GT franchise has sold more than 80 million units. The most recent release, Gran Turismo Sport, is known for precise vehicle dynamics simulation and racing realism, earning it the distinction of being sanctioned by the FIA and selected as a platform for the first Virtual Olympics (<https://olympics.com/en/sport-events/olympic-virtual-motorsport-event/>). GT Sport runs only on PS4s and at a 60-Hz-dynamics simulation cycle. A maximum of 20 cars can be in a race.

Our agent ran asynchronously on a separate computer and communicated with the game by means of HTTP over wired Ethernet. The agent requested the latest observation and made decisions at 10 Hz (every 100 ms). We tested action frequencies from 5 Hz to 60 Hz and found no substantial performance gains from acting more frequently than 10 Hz. The agent had to be robust to the infrequent, but real, networking delays. The agent's action was treated the same as a human's game controller input, but only a subset of action capabilities were supported in the GT API. For example, the API did not allow the agent to control gear shifting, the traction control system or the brake balance, all of which can be adjusted in-game by human players.

Computing environment

Each experiment used a single trainer on a compute node with either one NVIDIA V100 or half of an NVIDIA A100 coupled with around eight vCPUs and 55 GiB of memory. Some of these trainers were run in PlayStation Now data centres and others in AWS EC2 using p3.2xlarge instances.

Each experiment also used a number of rollout workers, where each rollout worker consisted of a compute node controlling a PS4. In this set-up, the PS4 ran the game and the compute node managed the rollouts by performing tasks such as computing actions, sending them to the game, sending experience streams to the trainer and receiving revised policies from the trainer (see Fig. 1a). The compute node used around two vCPUs and 3.3 GB of memory. In the time-trial experiments, ten rollout workers (and therefore ten PS4s) were used for about 8 days. To train policies that could drive in traffic, 21 rollout workers were used for between 7 and 12 days. In both cases, one worker was primarily evaluating intermediate policies, rather than generating new training data.

Actions

The GT API enabled control of three independent continuous actions: throttle, brake and steering. Because the throttle and brake are rarely engaged at the same time, the agent was presented control over the throttle and brake as one continuous action dimension. The combined dimension was scaled to $[-1, 1]$. Positive values engaged the throttle (with maximum throttle at +1), whereas negative values engaged the brake (with maximum braking at -1); the value zero engaged neither the throttle nor the brake. The steering dimension was also scaled to $[-1, 1]$, where the extreme values corresponded to the maximum steering angle possible in either direction for the vehicle being controlled.

The policy network selected actions by outputting a squashed normal distribution with a learned mean and diagonal covariance matrix over these two dimensions. The squashed normal distribution enforced sampled actions to always be within the $[-1, 1]$ action bounds³⁶. The diagonal covariance matrix values were constrained to lie in the range (e^{-40}, e^4) , allowing for nearly deterministic or nearly uniform random action selection policies to be learned.

Features

A variety of state features were input to the neural networks. These features were either directly available from the game state or processed into more convenient forms and concatenated before being input to the models.

Time-trial features. To learn competent time-trial performance, the agent needed features that allowed it to learn how the vehicle behaved and what the upcoming course looked like. The list of vehicle features included the car's 3D velocity, 3D angular velocity, 3D acceleration, load on each tyre and tyre-slip angles. Information about the environment was converted into features including the scalar progress of the car along the track represented as sine and cosine components, the local course surface inclination, the car's orientation with respect to the course centre line and the (left, centre and right) course points describing the course ahead on the basis of the car's velocity. The agent also received indicators if it contacted a fixed barrier or was considered off course by the game and it received real-valued feedback for the game's view of the car's most recent steering angle, throttle intensity and brake intensity. We relied on the game engine to determine whether the agent was off course (defined as when three or more tyres are out of bounds) because the out-of-bounds regions are not exactly defined by the course edges; kerbs and other tarmac areas outside the track edges are often considered in bounds.

Racing features. When training the agent to race against other cars, the list of features also included a car contact flag to detect collisions and a slipstream scalar that indicates if the agent was experiencing the slipstream effect from the cars in front of it. To represent the nearby cars, the agent used a fixed forward and rear distance bound to determine which cars to encode. The cars were ordered by their relative distance to the agent and were represented using their relative centre-of-mass position, velocity and acceleration. The combination of features provided the information required for the agent to drive fast and learn to overtake cars while avoiding collisions.

To keep the features described here in a reasonable numerical range when training neural networks, we standardized the inputs on the basis of the knowledge of the range of each feature scalar. We assumed that the samples were drawn from a uniform distribution given the range and computed the expected mean and standard deviation. These were used to compute the z-score for each scalar before being input to the models.

Rewards

The reward function was a hand-tuned linear combination of reward components computed on the transition between the previous state s and current state s' . The reward components were: course progress (R_{cp}), off-course penalty (R_{soc} or R_{loc}), wall penalty (R_w), tyre-slip penalty (R_{ts}), passing bonus (R_{ps}), any-collision penalty (R_c), rear-end penalty (R_r) and unsporting-collision penalty (R_{uc}). The reward weightings for the three tracks are shown in Extended Data Table 1.

Owing to the high speeds on Sarthe, training for that track used a slightly different off-course penalty, included the unsporting-collision penalty and excluded the tyre-slip penalty. Note that, to reduce variance in time-sensitive rewards, such as course progress and off-course penalty, we filtered transitions when network delays were encountered. The components are described in detail below.

Course progress (R_{cp}). Following previous work²⁴, the primary reward component rewarded the amount of progress made along the track since the last observation. To measure progress, we made use of the state variable l that measured the length (in metres) along the centreline from the start of the track. The agent's centreline distance l was estimated by first projecting its current position to the closest point on the centreline. The progress reward was the difference in l between the previous and the current state: $R_{cp}(s, s') \triangleq s'_l - s_l$. To reduce the incentive to cut corners, this reward was masked when the agent was driving off course.

Off-course penalty (R_{soc} or R_{loc}). The off-course reward penalty was proportional to the squared speed the agent was travelling at to further discourage corner cutting that may result in a very large gain in

Article

position: $R_{\text{soc}}(s, s') \triangleq - (s'_o - s_o)(s'_{\text{kph}})^2$, where s_o is the cumulative time off course and s_{kph} is the speed in kilometres per hour. To avoid an explosion in values at Sarthe where driving speeds were markedly faster and corners particularly easy to cut, we used a penalty that was proportional to the speed (not squared): $R_{\text{loc}}(s, s') \triangleq - (s'_o - s_o)s'_{\text{kph}}$, and the penalty was doubled for the difficult first and final chicanes.

Wall penalty (R_w). To assist the agent in learning to avoid walls, a wall-contact penalty was included. This penalty was proportional to the squared speed of the car and the amount of time in contact with the wall since the last observation: $R_w(s, s') \triangleq - (s'_w - s_w)(s'_{\text{kph}})^2$, where s_w is the cumulative time that the agent was in contact with a wall.

Tyre-slip penalty (R_{ts}). Tyre slip makes it more difficult to control the car. To assist learning, we included a penalty when the tyres were slipping in a different direction from where they were pointing: $R_{\text{ts}}(s, s') \triangleq - \sum_i^4 \min(|s'_{\text{tsr},i}|, 1.0) |s'_{\text{ts}\theta,i}|$, where $s_{\text{tsr},i}$ is the tyre-slip ratio for the i th tyre and $s_{\text{ts}\theta,i}$ is the angle of the slip from the forward direction of the i th tyre.

Passing bonus (R_{ps}). As in previous work²⁵, to incentivize passing opponents, we included a term that positively rewarded gaining ground and overtaking opponents, and negatively rewarded losing ground to an opponent. The negative reward ensured that there were no positive-cycle reward loops to exploit and encouraged defensive play when an opponent was trying to overtake the agent. This reward was defined as $R_{\text{ps}}(s, s') \triangleq \sum_i (s_{L_i} - s'_{L_i}) \max(\mathbf{1}_{(b,f)}(s_{L_i}), \mathbf{1}_{(b,f)}(s'_{L_i}))$, where s_{L_i} is the projected centreline signed distance (in metres) from the agent to opponent i and $\mathbf{1}_{(b,f)}(x)$ is an indicator function for when an opponent is no more than b metres behind nor f metres in front of the agent. We used $b = 20$ and $f = 40$ to train GT Sophy. The max operator ensures that the reward is provided when the agent was within bounds in the previous state or in the current state. In the particularly complex first and final chicanes of Sarthe, we masked this passing bonus to strongly discourage the agent from cutting corners to gain a passing reward.

Any-collision penalty (R_c). To discourage collisions and pushing cars off the road, we included a reward penalty whenever the agent was involved in any collision. This was defined as a negative indicator whenever the agent collided with another car: $R_c(s, s') \triangleq - \max_{i \in N} s'_{c,i}$, where $s_{c,i}$ is 1 when the agent collided with opponent i and 0 otherwise, and N is the number of opponents.

Rear-end penalty (R_r). Rear-ending an opponent was one of the more common ways to cause an opponent to lose control and for the agent to be penalized by stewards. To discourage bumping from behind, we included the penalty $R_r(s, s') \triangleq - \sum_i s'_{c,i} \cdot \mathbf{1}_{>0}(s'_{L_i} - s'_l) \cdot \|s'_v - s'_{v,i}\|_2^2$, where $s_{c,i}$ is a binary indicator for whether the agent was in a collision with opponent i , $\mathbf{1}_{>0}(s_{L_i} - s_l)$ is an indicator for whether opponent i was in front of the agent, s_v is the velocity vector of the agent and $s_{v,i}$ is the velocity vector of opponent i . The penalty was dependent on speed to more strongly discourage higher speed collisions.

Unsporting-collision penalty (R_{uc}). Owing to the high speed of cars and the technical difficulty of Sarthe, training the agent to avoid collisions was particularly challenging. Merely increasing the any-collision penalty resulted in very timid agent behaviour. To discourage being involved in collisions without causing the agent to be too timid, we included an extra collision penalty for Sarthe. Like the any-collision penalty, this penalty was a negative Boolean indicator. Unlike the any-collision penalty, it only fired when the agent rear-ended or sideswiped an opponent on a straightaway or was in a collision in a curve that was not caused by an opponent rear-ending them: $R_{\text{uc}}(s, s') \triangleq - \max_{i \in N} u(s', i)$, where $u(s', i)$ indicates an unsporting collision as defined above.

Training algorithm

To train our agent, we used an extension of the soft actor-critic algorithm³⁶ that we refer to as QR-SAC. To give the agent more capacity to predict the variation in the environment during a race, we make use of a QR Q-function³⁹ modified to accept continuous actions as inputs. QR-SAC is similar to distributional SAC⁴⁴ but uses a different formulation of the value backup and target functions. We used $M = 32$ quantiles and modified the loss function of the QR Q-function with an N -step temporal difference backup. The target function, y_i , for the i th quantile, \hat{t}_i , consists of terms for the immediate reward, $R_t = \sum_{i=1}^N \gamma^{i-1} r_{t+i}$, the estimated quantile value at the N th future state, $Z_{\hat{t}_i}$, and the SAC entropy term. Like existing work using N -step backups³⁸, we do not correct for the off-policy nature of N -step returns stored in the replay buffer. To avoid the computational cost of forwarding the policy for intermediate steps of the N -step backup, we only include the entropy reward bonus that SAC adds for encouraging exploration in the final step of the N -step backup. Despite this lack of off-policy correction and limited use of the entropy reward bonus, we found that using N -step backups greatly improved performance compared with a standard one-step backup, as shown in Fig. 2d. To avoid overestimation bias, the N th state quantiles are taken from the Q-function with the smallest N th state mean value⁴⁵, indexed by k

$$k = \arg \min_{m=1,2} Q(s_{t+N}, a' | \theta'_m) \quad (1)$$

$$y_i = R_t + Z_{\hat{t}_i}(s_{t+N}, a' | \theta_k) - \alpha \log \pi(a' | s_{t+N}, \phi)$$

where θ and ϕ are parameters of the Q-functions and the policy, respectively. Using this target value, y_i , the loss function of the Q-function is defined as follows

$$\delta_{i,j} = y_i - Z_{\hat{t}_i}(s_t, a_t | \theta)$$

$$L(\theta) = \frac{1}{M^2} \sum_i \sum_j \mathbb{E}_{s_t, a_t, R_t, s_{t+1} \sim D, a' \sim \pi} \rho(\delta_{i,j}) \quad (2)$$

where D represents data from the ERB and ρ is a quantile Huber loss function³⁹. Finally, the objective function for the policy is as follows:

$$J(\phi) = \mathbb{E}_{s \sim D, a \sim \pi(a|s, \pi)} [\alpha \log \pi(a|s, \phi) - \min_{i=1,2} Q(s, a | \theta_i)] \quad (3)$$

The Q-functions and policy models used in the October race consist of four hidden layers with 2,048 units each and a rectified linear unit activation function. To achieve robust control, dropout⁴⁶ with a 0.1 drop probability is applied to the policy function⁴⁷. The parameters are optimized using an Adam optimizer⁴⁸ with learning rates of 5.0×10^{-5} and 2.5×10^{-5} for the Q-function and policy, respectively. The discount factor γ was 0.9896 and the SAC entropy temperature value α was set to 0.01. The mixing parameter when updating the target model parameters after every algorithm step was set to 0.005. The off-course penalty and rear-end-speed penalty can produce large penalty values due to the squared speed term, which makes the Q-function training unstable due to large loss values. To mitigate this issue, the gradients of the Q-function are clipped by the global norm of the of 10.

The rollout workers send state-transition tuples $\langle s, a, r \rangle$ collected in an episode (of length 150 s) to the trainer to store the data in an ERB implemented using the Reverb Python library⁴⁹. The buffer had capacity of 10^7 N -step transitions. The trainer began the training loop once 40,000 transitions had been collected and uses a mini-batch of size 1,024 to update the Q-function and policy. A training epoch consists of 6,000 gradient steps. After each epoch, the trainer sent the latest model parameters to the rollout workers.

Training scenarios

Learning to race requires mastering a gamut of skills: surviving a crowded start, making tactical open-road passes and precisely running the track alone. To encourage basic racing skills, we placed the agent in scenarios with zero, one, two, three or seven opponents launched nearby (1v0, 1v1, 1v2, 1v3 and 1v7, respectively). To create variety, we randomized track positions, start speeds, spacing between cars and opponent policies. We leveraged the fact that the game supports 20 cars at a time to maximize PlayStation usage by launching more than one group on the track. All base scenarios ran for 150 s. In addition, to ensure that the agent was exposed to situations that would allow it to learn the skills highlighted by our expert advisor, we used time-limited or distance-limited scenarios on specific course sections. Figure 1f illustrates the skill scenarios used at Sarthe: eight-car grid starts, 1v1 slipstream passing and mastering the final chicane in light traffic. Extended Data Figure 1 shows the specialized scenarios used to prepare the agent to race on Seaside (f) and Maggiore (g). To learn how to avoid catastrophic outcomes at the high-speed Sarthe track, we also incorporated mistake learning⁵⁰. During policy evaluations, if an agent lost control of the car, the state shortly before the event was recorded and used as a launch point for more training scenarios.

Unlike curriculum training where early skills are supplanted by later ones or in which skills build on top of one another in a hierarchical fashion, our training scenarios are complementary and were trained into a single control policy for racing. During training, the trainer assigned new scenarios to each rollout worker by selecting from the set configured for that track on the basis of hand-tuned ratios designed to provide sufficient skill coverage. See Extended Data Fig. 1e for an example ERB at Sarthe. However, even with this relative execution balance, random sampling fluctuations from the buffer often led to skills being unlearned between successive training epochs, as shown in Fig. 2h. Therefore, we implemented multi-table stratified sampling to explicitly enforce proportions of each scenario in each training mini-batch, notably stabilizing skill retention (Fig. 2g).

Policy selection

In machine learning, convergence means that further training will not improve performance. In RL, due to the continuing exploration and random sampling of experiences, the performance of the policy will often continue to vary after convergence (Fig. 2h). Therefore, even with the stabilizing techniques described above, continuing training after convergence produced policies that differed in small ways in their ability to execute the desired racing skills. A subsequent policy, for instance, may become marginally better at the slipstream pass and marginally worse at the chicane. Choosing which policy to race against humans became a complex, multi-objective optimization problem.

Extended Data Figure 3 illustrates the policy-selection process. Agent policies were saved at regular intervals during training. Each saved policy then competed in a single-race scenario against other AI agents, and various metrics, such as lap times and car collisions, were gathered and used to filter the saved policies to a smaller set of candidates. These candidates were then run through an *n*-athlon—a set of pre-specified evaluation scenarios—testing their lap speed and performance in certain tactically important scenarios, such as starting and using the slipstream. The performance on each scenario was scored and the results of each policy on each scenario were combined in a single ranked spreadsheet. This spreadsheet, along with various plots and videos, was then reviewed by a human committee to select a small set of policies that seemed the most competitive and the best behaved. From this set, each pair of policies competed in a multi-race, round-robin, policy-versus-policy tournament. These competitions were scored using the same team scoring as that in the exhibition event and evaluated on collision metrics. From these results, the committee chose policies that seemed to have the best chance of winning against the human drivers while minimizing

penalties. These final candidate policies were then raced against test drivers at Polyphony Digital and the subjective reports of test drivers were factored into the final decision.

The start of Sarthe posed a particularly challenging problem for policy selection. Because the final chicane is so close to the starting line, the race was configured with a stationary grid start. From that standing start, all eight cars quickly accelerated and entered the first chicane. Although a group of eight GT Sophy agents might get through the chicane fairly smoothly, against human drivers, the start was invariably chaotic and a fair amount of bumping occurred. We tried many variations of our reward functions to find a combination that was deemed an acceptable starter by our test drivers while not giving up too many positions. In the October 2021 Sarthe race, we configured GT Sophy to use a policy that started well, and—after 2,100 metres—switch to a slightly more competitive policy for the rest of the race. Despite the specialized starter, the instance of GT Sophy that began the race in pole position was involved in a collision with a human driver in the first chicane, slid off the course and fell to last place. Despite that setback, it managed to come back and win the race.

Immediately after the official race, we ran a friendly rematch against the same drivers but used the starter policy for the whole track. The results were similar to the official race.

Fairness versus humans

Competitions between humans and AI systems cannot be made entirely fair; computers and humans think in different ways and with different hardware. Our objective was to make the competition fair enough, while using technical approaches that were consistent with how such an agent could be added to the game. The following list compares some of the dimensions along which GT Sophy differs from human players:

First, perception. GT Sophy had a map of the course with precise *x*, *y* and *z* information about the points that defined the track boundaries. Humans perceived this information less precisely by means of vision. However, the course map did not have all of the information about the track and humans have an advantage in that they could see the kerbs and surface material outside the boundaries, whereas GT Sophy could only sense these by driving on them.

Second, opponents. GT Sophy had precise information about the location, velocity and acceleration of the nearby vehicles. However, it represented these vehicles as single points, whereas humans could perceive the whole vehicle. GT Sophy has a distinct advantage in that it can see vehicles behind it as clearly as it can see those in front, whereas humans have to use the mirrors or the controller to look to the sides and behind them. GT Sophy never practiced against opponents that didn't have full visibility, so it didn't intentionally take advantage of human blind spots.

Third, vehicle state. GT Sophy had precise information about the load on each tyre, slip angle of each tyre and other vehicle state. Humans learn how to control the car with less precise information about these state variables.

Fourth, vehicle controls. There are certain vehicle controls that the human drivers had access to that GT Sophy did not. In particular, expert human drivers often use the traction control system in grid starts and use the transmission controls to change gears.

Fifth, action frequency. GT Sophy took actions at 10 Hz, which was sufficient to control the car but much less frequent than human actions in GT. Competitive GT drivers use steering and pedal systems that give them 60 Hz control. Whereas a human can't take 60 distinct actions per second, they can smoothly turn a steering wheel or press on a brake pedal. Extended Data Figure 2b, c contrasts GT Sophy's 10-Hz control pattern to Igor Fraga's much smoother actions in a corner of Sarthe.

Sixth, reaction time. GT Sophy's asynchronous communication and inference takes around 23–30 ms, depending on the size of the network. Although evaluating performance in professional athletes and gamers is a complex field^{34,35}, an oft-quoted metric is that professional athletes have a reaction time of 200–250 ms. To understand how the

Article

performance of GT Sophy would be affected if its reaction time were slowed down, we ran experiments in which we introduced artificial delays to its perception pipeline. We retrained our agent with delays of 100 ms, 200 ms and 250 ms in the Maggiore time-trial setting, using the same model architecture and algorithm as our time-trial baseline. All three of these tests achieved a superhuman lap time.

Tests versus top GT drivers

The following competitive GT drivers participated in the time-trial evaluations:

- Emily Jones: 2020 FIA Gran Turismo Manufacturers Series, Team Audi.
- Valerio Gallo: 2nd place 2020 FIA Gran Turismo Nations Cup; winner 2021 Olympic Virtual Series Motor Sport Event; winner 2021 FIA Gran Turismo Nations Cup.
- Igor Fraga: winner 2018 FIA Gran Turismo Nations Cup; winner 2019 Manufacturer Series championship; winner 2020 Toyota Racing Series (real racing).

GT Sophy won all of the time-trial evaluations as shown in Fig. 3f and was reliably superhuman on all three tracks, as shown in Fig. 1d, e and Extended Data Fig. 1a–d. Notably, the only human with a time within the range of GT Sophy's 100 lap times on any of the tracks was Valerio Gallo on Maggiore. It is worth noting that the data in Fig. 1d, e was captured by Polyphony Digital after the time-trial event in July 2021. Valerio was the only participant represented in the data that had seen the trajectories of GT Sophy on Maggiore, and—between those two events—Valerio's best time improved from 114.466 to 114.181 s.

It is also interesting to examine what behaviours give GT Sophy such an advantage in time trials. Extended Data Figure 2a shows an analysis of Igor's attempt to match GT Sophy on Sarthe, showing the places on the course where he fell farther behind. Not surprisingly, the hardest chicanes and corners are the places where GT Sophy has the biggest performance gains. In most of these corners, Igor seems to catch up a little bit by braking later, but is then unable to take the corner itself as fast, resulting in him losing ground overall.

The following competitive GT drivers participated in the team racing event:

- Takuma Miyazono: winner 2020 FIA Gran Turismo Nations Cup; winner 2020 FIA Gran Turismo Manufacturer Series; winner 2020 GR Supra GT Cup.
- Tomoaki Yamanaka: winner 2019, 2021 Manufacturer Series.
- Ryota Kokubun: winner 2019 FIA Gran Turismo Nations Cup, Round 5, Tokyo; 3rd place 2020 FIA Gran Turismo Nations Cup.
- Shotaro Ryu: 2nd place Japan National Inter-prefectural Esports Championship (National Athletic Meet) 2019 Gran Turismo Division (Youth).

Driver testimonials

The following quotes were captured after the July 2021 events:

"I think the AI was very fast turning into the corner. How they approach into it, as well as not losing speed on the exit. We tend to sacrifice a little bit the entry to make the car be in a better position for the exit, but the AI seems to be able to carry more speed into the corner but still be able to have the same kind of exit, or even a faster exit. The AI can create this type of line a lot quicker than us,... it was not a possibility before because we never realized it. But the AI was able to find it for us." – Igor Fraga

"It was really interesting seeing the lines where the AI would go, there were certain corners where I was going out wide and then cutting back in, and the AI was going in all the way around, so I learned a lot about the lines. And also knowing what to prioritize. Going into turn 1 for example, I was braking later than the AI, but the AI would get a much better exit

than me and beat me to the next corner. I didn't notice that until I saw the AI and was like 'Okay, I should do that instead'" – Emily Jones

"The ghost is always a reference. Even when I train I always use someone else's ghost to improve. And in this case with such a very fast ghost,... even though I wasn't getting close to it, I was getting closer to my limits." – Valerio Gallo

"I hope we can race together more, as I felt a kind of friendly rivalry with [GT Sophy]." (translated from Japanese) – Takuma Miyazono

"There is a lot to learn from [GT Sophy], and by that I can improve myself. [GT Sophy] does something original to make the car go faster, and we will know it's reasonable once we see it." (translated from Japanese) – Tomoaki Yamanaka

Data availability

There are no static data associated with this project. All data are generated from scratch by the agent each time it learns. Videos of the races are available at https://sonyai.github.io/gt_sophy_public.

Code availability

Pseudocode detailing the training process and algorithms used is available as a supplement to this article. The agent interface in GT is not enabled in commercial versions of the game; however, Polyphony Digital has provided a small number of universities and research facilities outside Sony access to the API and is considering working with other groups.

44. Xia, L., Zhou, Z., Yang, J. & Zhao, Q. DSAC: distributional soft actor critic for risk-sensitive reinforcement learning. Preprint at <https://arxiv.org/abs/2004.14547> (2020).
45. Fujimoto, S., van Hoof, H. & Meger, D. In *Proc. 35th International Conference on Machine Learning* 1587–1596 (PMLR, 2018).
46. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
47. Liu, Z., Li, X., Kang, B. & Darrell, T. In *International Conference on Learning Representations* (ICLR, 2021).
48. Kingma, D. P. & Ba, J. In *International Conference on Learning Representations* (ICLR, 2015).
49. Cassirer, A. et al. Reverb: a framework for experience replay. Preprint at <https://arxiv.org/abs/2102.04736> (2021).
50. Narvekar, S., Sinapov, J., Leonetti, M. & Stone, P. In *Proc. 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2016)* (2016).

Acknowledgements We thank K. Yamauchi, S. Takano, A. Hayashi, C. Ferreira, N. Nozawa, T. Teramoto, M. Hakim, K. Yamada, S. Sakamoto, T. Ueda, A. Yago, J. Nakata and H. Imanishi at Polyphony Digital for making the Gran Turismo franchise, providing support throughout the project and organizing the Race Together events on 2 July 2021 and 21 October 2021. We also thank U. Gallizzi, J. Beltran, G. Albowicz, R. Abdul-ahad and the staff at CGEI for access to their PlayStation Now network to train agents and their help building the infrastructure for our experiments. We benefited from the advice of T. Grossenbacher, a retired competitive GT driver. Finally, we thank E. Kato Marcus and E. Ohshima of Sony AI, who managed the partnership activities with Polyphony Digital and Sony Interactive Entertainment.

Author contributions P.R.W. managed the project. S.B., K.K., P.K., J.M., K.S. and T.J.W. led the research and development efforts. R.C., A.D., F.E., F.F., L.G., V.K., H.L., P.M., D.O., C.S., T.S. and M.D.T. participated in the research and the development of GT Sophy and the AI libraries. H.A., L.B., R.D. and D.W. built the research platform that connected to CGEI's PlayStation network. P.S. provided executive support and technical and research advice and P.D. provided executive support and technical advice. H.K. and M.S. conceived and set up the project, provided executive support, resources and technical advice and managed stakeholders.

Competing interests P.R.W. and other team members have submitted US provisional patent application 63/267,136 covering aspects of the scenario training techniques described in this paper.

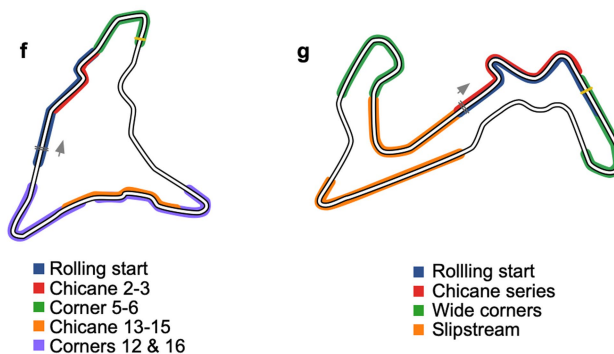
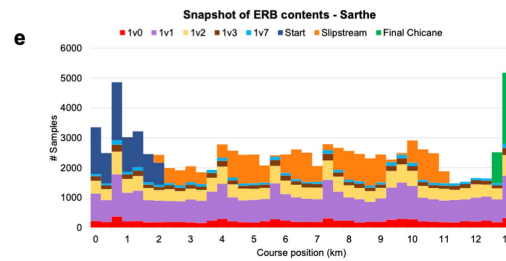
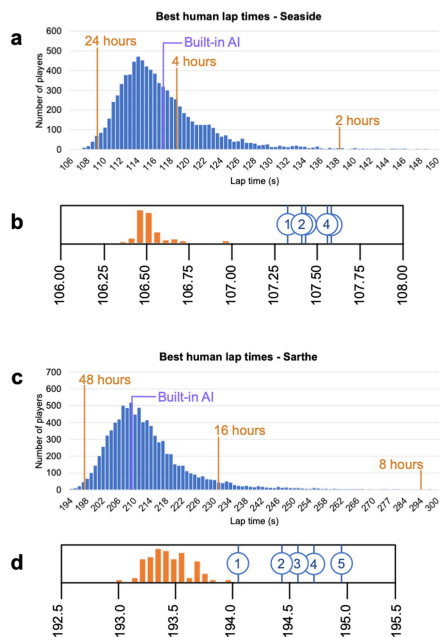
Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41586-021-04357-7>.

Correspondence and requests for materials should be addressed to Peter R. Wurman.

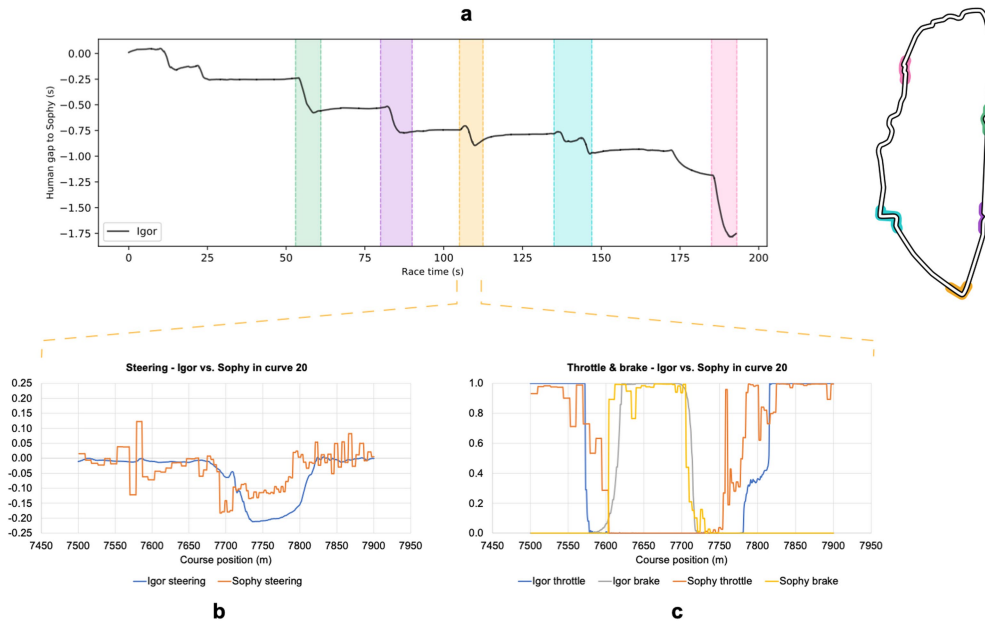
Peer review information Nature thanks the anonymous reviewers for their contribution to the peer review of this work.

Reprints and permissions information is available at <http://www.nature.com/reprints>.



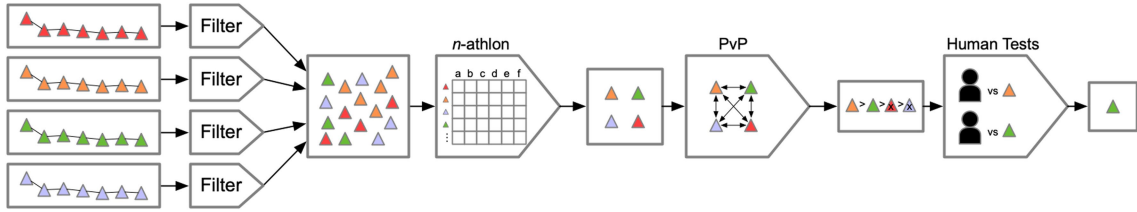
Extended Data Fig. 1 | Seaside and Sarthe training. Kudos Prime data from global time-trial challenges on Seaside (a and b) and Sarthe (c and d), with the cars used in the competition. Note that these histograms represent the single best lap time for more than 12,000 individual players on Seaside and almost 9,000 on Sarthe. In both cases, the secondary diagrams compare the top five human times to a histogram of 100 laps by the 2 July 2021 time-trial version of GT Sophy. In both cases, the data show that GT Sophy was reliably superhuman, with all 100 laps better than the best human laps. Not surprisingly, it takes

longer for the agent to train on the much longer Sarthe course, taking 48 h to reach the 99th percentile of human performance. e, Histogram of a snapshot of the ERB during training on Sarthe on the basis of the scenario breakdown in Fig. 1f. The x axis is the course position and the stacked colours represent the number of samples that were collected in that region from each scenario. In a more condensed format than Fig. 1f, f and g show the sections of Seaside and Maggiore that were used for skill training.



Extended Data Fig. 2 | Time trial on Sarthe. An analysis of Igor Fraga’s best lap in the time-trial test compared with GT Sophy’s lap. **a**, Areas of the track where Igor lost time with respect to GT Sophy. Corner 20, highlighted in yellow, shows an interesting effect common to the other corners in that Igor seems to catch up a little by braking later, but then loses time because he has to brake longer

and comes out of the corner slower. Igor’s steering controls (**b**) and Igor’s throttle and braking (**c**) compared with GT Sophy on corner 20. Through the steering wheel and brake pedal, Igor is able to give smooth, 60-Hz signals compared with GT Sophy’s 10-Hz action rate.



Extended Data Fig. 3 | Policy selection. An illustration of the process by which policies were selected to run in the final race. Starting on the left side of the diagram, thousands of policies were generated and saved during the experiments. They were first filtered in the experiment to select the subset on the Pareto frontier of a simple evaluation criteria trading off lap time versus off-course and collision metrics. The selected policies were run through a series of tests evaluating their overall racing performance against a common set of opponents and their performance on a variety of hand-crafted skill tests.

The results were ranked and human judgement was applied to select a small number of candidate policies. These policies were matched up in round-robin, policy-versus-policy competitions. The results were again analysed by the human committee for overall team scores and collision metrics. The best candidate policies were run in short races against test drivers at Polyphony Digital. Their subjective evaluations were included in the final decisions on which policies to run in the October 2021 event.

Article

Extended Data Table 1 | Reward weights

Course	R_{cp}	R_{soc}	R_{loc}	R_w	R_{ts}	R_{ps}	R_c	R_r	R_{uc}
Seaside	1	0.01	0	0.01	0.25	0.5	5	0.1	0
Maggiore	1	0.01	0	0.01	0.25	0.5	4	0.1	0
Sarthe	1	0	5	0.01	0	0.5	5	0.1	5

Reward weights for each track.

Supplementary information

**Outracing champion Gran Turismo drivers
with deep reinforcement learning**

In the format provided by the
authors and unedited

This document is supplementary material to “Outracing champion Gran Turismo drivers with deep reinforcement learning.”

Approximate Python Code

Here we provide approximate Python code that replicates the most critical aspects of our work. For the purposes of illustration, the code presented here is implemented in less computationally efficient ways than our actual implementation; has abstracted away many of the system implementation details such as remote process communication and check-pointing; and is written more directly than our actual implementation which supported wide configuration capabilities that otherwise detract from readability. Our implementation uses TensorFlow [1]; all references to ‘tf.x’ refer to an operation/class in the TensorFlow API, and references to ‘keras.x’ refers to an operation/class in the TensorFlow Keras API. A high-level overview of the training and rollout processes is provided in Figure 1.

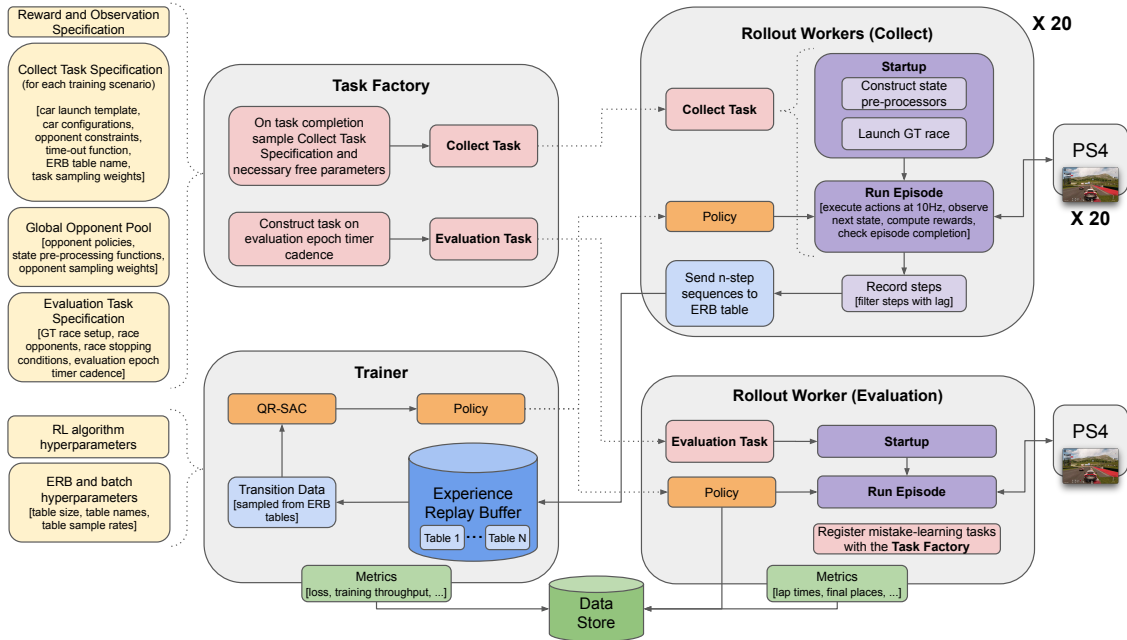


Figure 1: Diagram of modules in the GT Sophy training process

```

1 def train_loop():
2     models = Models(
3         policy=make_policy_network(),
4         af1=make_qr_q_network(),
5         af2=make_qr_q_network(),
6         af1_target=make_qr_q_network(),
7         af2_target=make_qr_q_network(),
8     )
9     # Sync prediction and target networks before start
10    copy_weights(models.af1.trainable_params, models.af1_target.trainable_params)
11    copy_weights(models.af2.trainable_params, models.af2_target.trainable_params)
12
13    opts = Optimizers(
14        policy=keras.optimizers.Adam(learning_rate=2.5e-5),
15        critic=keras.optimizers.Adam(learning_rate=5.0e-5, global_clipnorm=10.0),
16    )
17
18    param_server = start_parameter_server(models.policy)
19    replay_server = start_experience_replay_server()
20    table_datasets = [
21        make_dataset_for_table(replay_server, table_name)
22        for table_name in TABLE_NAMES
23    ]
24    task_server = start_task_factory_server(warm_up=True)
25
26    # Wait for MINIMUM_STEPS to be recorded to replay before we start any training
27    replay_server.block_and_wait(MINIMUM_STEPS)
28
29    while True:
30        for _ in range(BATCHES_PER_EPOCH):
31            transition_data = sample_data(table_datasets)
32            # see QR-SAC code listing for implementation
33            losses = step_qr_sac(models, opts, transition_data)
34            write_metrics(losses)
35            param_server.update_params(models.policy)
36            # warmup over after first epoch
37            task_server.warm_up = False
38

```

```

39 def sample_data(table_datasets: List[tf.data.Dataset]) -> Transition:
40 # For simplicity, this code omits handling the edge case when only a subset of
41 # the tables have data to sample.
42 # Before concat, each dataset has a different batch size based on the table
43 # weights
44 sub_traj_batches = [next(iter(dataset)) for dataset in table_datasets]
45 sub_traj_batches = nested_concat(sub_traj_batches, axis=0)
46 # sub_traj_batches has fields obs, action, reward, done
47 # Each tensor has two batch dimensions NT and 1 channel dimension, where N is
48 # batch and T is time. The obs T dim is one longer than other T dims.
49
50 discounting = tf.pow(DISCOUNT, tf.range(N_STEP, dtype=tf.float32))
51 discounting = tf.reshape(discounting, (1, N_STEP, 1))
52 discounted_reward = tf.math.reduce_sum(
53     sub_traj_batches.reward * discounting,
54     axis=1,
55 ) # (N, 1)
56 # Resulting transition removes T dimension, selecting relevant positions
57 return Transition(
58     obs=sub_traj_batches.obs[:, 0, :], # first obs
59     action=sub_traj_batches.action[:, 0, :], # first action
60     reward=discounted_reward, # discounted cumulative reward
61     next_obs=sub_traj_batches.obs[:, -1, :], # last obs
62     done=sub_traj_batches.done[:, -1, :], # last done
63 )

```

Listing 1: Training Loop Code

```

1 def step_qr_sac(m: Models, opts: Optimizers, t: Transition,) -> Losses:
2
3     critic_params = m.af1.trainable_params + m.af2.trainable_params
4     target_critic_params = (
5         m.af1_target.trainable_params + m.af2_target.trainable_params
6     )
7
8     # these values do not need gradients traced through them
9     policy_head_next = m.policy.policy_head(t.next_obs)
10    actions_next = m.policy.sample_from_head(policy_head_next)
11    log_prob_next = m.policy.log_prob(policy_head_next, actions_next)
12
13    # update policy
14    with tf.GradientTape() as tape:
15        policy_head = m.policy.policy_head(t.obs)
16        sampled_actions = m.policy.sample_from_head(policy_head)
17        log_prob = m.policy.log_prob(policy_head, sampled_actions)
18        policy_loss = policy_loss(
19            log_prob_samples=log_prob,
20            q1_sampled=m.af1.action_value(t.obs, sampled_actions),
21            q2_sampled=m.af2.action_value(t.obs, sampled_actions),
22        )
23    opts.policy.minimize(policy_loss, m.policy.trainable_params, tape)
24
25    # update critic networks
26    with tf.GradientTape() as tape:
27        critic_loss = critic_loss(
28            q1_quantile_observed=m.af1.quantile(t.obs, t.action),
29            q2_quantile_observed=m.af2.quantile(t.obs, t.action),
30            q1_quantile_sampled_next=m.af1_target.quantile(t.next_obs, actions_next),
31            q2_quantile_sampled_next=m.af2_target.quantile(t.next_obs, actions_next),
32            log_prob_next=log_prob_next,
33            reward=t.reward,
34            done=t.done,
35        )
36    opts.critic.minimize(critic_loss, critic_params, tape)
37
38    # move target network params toward prediction networks
39    for pred_param, target_param in zip(critic_params, target_critic_params):
40        target_param.assign(
41            SMOOTH_FACTOR * pred_param + (1.0 - SMOOTH_FACTOR) * target_param
42        )
43
44    return Losses(policy_loss, critic_loss)
45

```

```

46 def policy_loss(
47     log_prob_samples: tf.Tensor, # (N, 1)
48     q1_sampled: tf.Tensor, # (N, 1)
49     q2_sampled: tf.Tensor, # (N, 1)
50 ) -> tf.Tensor:
51     min_q = tf.minimum(q1_sampled, q2_sampled)
52     return tf.math.reduce_mean(ALPHA * log_prob_samples - min_q)
53
54
55 def critic_loss(
56     q1_quantile_observed: tf.Tensor, # (N, M)
57     q2_quantile_observed: tf.Tensor, # (N, M)
58     q1_quantile_sampled_next: tf.Tensor, # (N, M)
59     q2_quantile_sampled_next: tf.Tensor, # (N, M)
60     log_prob_next: tf.Tensor, # (N, 1)
61     reward: tf.Tensor, # (N, 1)
62     done: tf.Tensor, # (N, 1)
63 ) -> tf.Tensor:
64     target = tf.stop_gradient(qr_sac_bootstrap_target(
65         q1_quantile_sampled_next,
66         q2_quantile_sampled_next,
67         log_prob_next,
68         reward,
69         done,
70         DISCOUNT**N_STEP, # adjust for N-step transition
71         ALPHA,
72     ))
73
74     q1_loss = tf.math.reduce_mean(quantile_loss(q1_quantile_observed, target))
75     q2_loss = tf.math.reduce_mean(quantile_loss(q2_quantile_observed, target))
76     return q1_loss + q2_loss
77
78
79 def quantile_loss(
80     quantile: tf.Tensor, # (N, M)
81     quantile_target: tf.Tensor, # (N, M)
82 ) -> tf.Tensor:
83     batch_size = quantile.shape[0]
84     num_quantiles = quantile.shape[1]
85
86     y = tf.reshape(quantile_target, [batch_size, -1, 1]) # (N, M, 1)
87     x = tf.reshape(quantile, [batch_size, 1, -1]) # (N, 1, M)
88
89     # compute huber loss
90     diff = y - x # (N, M, M)
91     huber_loss = tf.where(
92         tf.math.abs(diff) < 1.0, 0.5 * diff ** 2, tf.math.abs(diff) - 0.5
93     ) # (N, M, M)
94
95     # sample taus
96     steps = tf.range(num_quantiles, dtype=tf.float32) # (M,)
97     taus = tf.reshape((steps + 1) / num_quantiles, [1, 1, -1]) # (1, 1, M)
98     taus_minus = tf.reshape((steps / num_quantiles), [1, 1, -1]) # (1, 1, M)
99     taus_hat = (taus + taus_minus) / 2.0 # (1, 1, M)
100
101     # compute quantile loss
102     delta = tf.cast(tf.stop_gradient(diff) < 0.0, dtype=tf.float32) # (N, M, M)
103     element_wise_loss = tf.math.abs(taus_hat - delta) * huber_loss # (N, M, M)
104
105     return tf.reduce_mean(tf.reduce_sum(element_wise_loss, axis=2), axis=1) # scalar

```

```

106 def qr_sac_bootstrap_target(
107     q1_quantile_sample_next: tf.Tensor, # (... , M)
108     q2_quantile_sample_next: tf.Tensor, # (... , M)
109     next_log_prob_samples: tf.Tensor, # (... , 1)
110     reward: tf.Tensor, # (... , 1)
111     done: tf.Tensor, # (... , 1)
112     discount: float, # scalar
113     alpha: float, # scalar
114 ) -> tf.Tensor:
115     quantile_next = minimum_quantile(
116         q1_quantile_sample_next, q2_quantile_sample_next
117     )
118     entropy = alpha * next_log_prob_samples
119     return reward + discount * (1.0 - done) * (quantile_next - entropy)
120
121
122 def minimum_quantile(
123     quantile1: tf.Tensor, # (N, M)
124     quantile2: tf.Tensor, # (N, M)
125 ) -> tf.Tensor:
126     batch_size = quantile1.shape[0]
127
128     # compute Q value
129     q1 = tf.reduce_mean(quantile1, axis=1, keepdims=True) # (N, 1)
130     q2 = tf.reduce_mean(quantile2, axis=1, keepdims=True) # (N, 1)
131
132     # pick quantiles based on minimum prediction
133     stacked_q = tf.concat([q1, q2], axis=1) # (N, 2)
134     min_indices = tf.cast(tf.math.argmax(stacked_q, axis=1), dtype=tf.int32) # (N,)
135     indices = tf.stack([tf.range(batch_size), min_indices], axis=1) # (N, 2)
136
137     stacked_quantile = tf.stack([quantile1, quantile2], axis=1) # (N, 2, M)
138     minimum_quantile = tf.gather_nd(stacked_quantile, indices) # (N, M)
139
140     return minimum_quantile

```

Listing 2: QR-SAC Code

```

1 def get_task() -> Task:
2     # The returned task defines which cars are agent controlled and for each opponent
3     # whether to control it with a pre-trained policy, PID, or the built-in AI
4     task_server = get_task_server_info()
5     if task_server.warm_up:
6         # During warm up, run 20 cars across the track at once all selecting
7         # uniformly random actions
8         return Task(
9             cars_launch=uniform_jittered_launch(num_cars=20, mph_range=(0, 60)),
10            random_policy=True,
11            agent_controlled_car_ids=list(range(20)),
12            car_id_to_opp_policy={},
13            car_id_to_opp_pid={},
14            builtin_ai_car_ids=[],
15            table_name="1v0",
16            time_out_fn=lambda steps, obs: steps >= MAX_STEPS,
17            is_eval=False,
18            builtin_ai_params=None,
19        )
20     elif task_server.needs_eval:
21         # Time for periodic evaluation
22         task_server.needs_eval = False
23         return Task(
24             cars_launch=EVAL_LAUNCH,
25             random_policy=False,
26             agent_controlled_car_ids=[0], # one car controlled by current policy
27             car_id_to_opp_policy=EVAL_OPP_POLICIES,
28             car_id_to_opp_pid={},
29             builtin_ai_car_ids=EVAL_OPP_BUILTIN_IDS,
30             table_name=None, # Do not collect training data in evals
31             time_out_fn=lambda steps, obs: steps >= MAX_STEPS,
32             is_eval=True,
33             builtin_ai_params=None,
34        )
35     elif (
36         # limited slots for mistake learning to prevent regular tasks from starving
37         len(task_server.mistake_learning_tasks) > 0
38         and task_server.mistake_slots_available()
39     ):
40         return task_server.pop_mistake_learning_task()
41     else:
42         # Standard data collection setting
43         # First we sample a task specification (e.g., 1v0, slipstream, etc.)
44         # The TaskSpec is almost a Task, but the launch, opponent policies, and
45         # builtin-ai params need to be sampled from a distribution
46         task_spec = np.random.choice(TASK_SPECS, p=TASK_SPEC_PROBS)
47         return Task(
48             cars_launch=task_spec.sample_launch(),
49             random_policy=False,
50             agent_controlled_car_ids=task_spec.agent_ids,
51             car_id_to_opp_policy=task_spec.sample_opp_policies(),
52             car_id_to_opp_pid=task_spec.pids,
53             builtin_ai_car_ids=task_spec.builtin_ids,
54             table_name=task_spec.table_name,
55             time_out_fn=task_spec.time_out_fn,
56             is_eval=False,
57             builtin_ai_params=task_spec.sample_builtin_ai_params(),
58        )

```

Listing 3: The task factory is run in a separate process and is remotely queried by rollout workers.

```

1 def rollout_loop():
2     policy = make_policy_network()
3     param_server = get_parameter_server()
4     ps_server = get_playstation_server()
5     replay_server = get_experience_replay_server()
6     task_server = get_task_server()
7     while True:
8         task = get_task()
9         if task.random_policy:
10            agent_policy = UniformRandomPolicy()
11        else:
12            load_weights(policy, param_server.get_params()) # Get latest policy
13            agent_policy = policy
14
15        ps_server.launch_race( # Start the race on the PlayStation
16            TRACK_NAME, CAR_INFO, task.cars_launch, task.builtin_ai_car_ids
17        )
18        obs = ps_server.wait_for_race_start()
19        p_obs = preprocess(obs) # Preprocess creates an obs batch for all cars
20        trajs = [
21            start_trajectory(initial_obs=p_obs[agent_id])
22            for agent_id in task.agent_car_ids()
23        ]
24
25        for step in range(MAX_STEPS):
26            actions = compute_all_actions(
27                task, agent_policy, p_obs, is_exploit=task.is_eval
28            )
29            ps_server.set_and_hold(actions)
30            # Get the next observation 6 frames later = 0.1 seconds @60fps
31            next_obs = ps_server.wait_for_obs_at_frame(obs.frame_count + 6)
32            rewards = compute_reward(obs, next_obs) # list of rewards for all cars
33            dones = [False] * task.cars_launch.num_cars # Continuing task
34
35            obs = next_obs
36            p_obs = preprocess(next_obs)
37
38            for agent_id, traj in zip(task.agent_car_ids(), trajs):
39                traj.append(
40                    obs=p_obs[agent_id],
41                    action=actions[agent_id],
42                    reward=rewards[agent_id],
43                    done=dones[agent_id],
44                )
45                if task.is_eval and should_retry(p_obs[agent_id]):
46                    task_server.add_mistake_learning_task(traj)
47
48            # create N-step subtrajectory if this task collects data
49            if step + 1 >= N_STEP and task.table_name:
50                for traj in trajs:
51                    # this writes one additional observation time step
52                    replay_server.write(task.table_name, traj[-N_STEP:])
53
54            if task.time_out_fn(step, obs):
55                break
56
57        if task.is_eval:
58            write_metrics(trajs)
59
60

```

```

61 def compute_all_actions(
62     task: Task,
63     agent_policy: Policy,
64     p_obs: tf.Tensor,
65     is_exploit: bool,
66 ) -> List[Optional[tf.Tensor]]:
67     # Actions for built-in AI cars remain None
68     all_actions = [None for _ in range(task.cars_launch.num_cars)]
69     # Compute actions for agent-controlled cars
70     for car_id in task.agent_car_ids():
71         if is_exploit:
72             a = agent_policy.exploit_action(p_obs[car_id])
73         else:
74             a = agent_policy.sample_action(p_obs[car_id])
75         all_actions[car_id] = a
76     # Compute actions for opponent cars using previously trained policies
77     for car_id, opp_policy in task.car_id_to_opp_policy.items():
78         if is_exploit:
79             a = opp_policy.exploit_action(p_obs[car_id])
80         else:
81             a = opp_policy.sample_action(p_obs[car_id])
82         all_actions[car_id] = a
83     # Compute actions for opponent cars using a PID controller
84     for car_id, pid in task.car_id_to_opp_pid.items():
85         # PID controllers always deterministically select actions
86         a = pid.exploit_action(p_obs[car_id])
87         all_actions[car_id] = a
88
89     return all_actions

```

Listing 4: Rollout worker code. 21 instances, each with their own PlayStation connection, are run in parallel.

Hyperparameters

Here we provide a set of tables that cover the hyperparameters used for the training experiments. We focus on the parameters used to train GT Sophy for competitive racing.

RL Algorithm	
n-step	7
discount factor	0.9896
entropy regularization coefficient, α	0.01
number of quantiles	32
parameter smoothing factor	0.005
policy learning rate	2.5e-5
critic learning rate	5.0e-5
optimizer	Adam
global clip norm (critic optimizer only)	10.0

Table 1: Table of parameters for QR-SAC RL algorithm.

Neural Network	
number of hidden layers	4
number of units per hidden layer	2048
activation	relu
dropout (policy network only)	0.1

Table 2: Table of parameters for the neural network architecture.

Trainer	
mini-batch size	1024
total replay buffer capacity	10e+6
number of buffer tables	8 to 10
training steps per epoch	6000
minimum transitions to take before training	40000

Table 3: Table of parameters for the trainer.

References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.