



## **Connected autonomous vehicles for improving mixed traffic efficiency in unsignalized intersections with deep reinforcement learning**

Downloaded from: <https://research.chalmers.se>, 2025-12-05 03:12 UTC

Citation for the original published paper (version of record):

Peng, B., Keskin, F., Kulcsár, B. et al (2021). Connected autonomous vehicles for improving mixed traffic efficiency in unsignalized intersections with deep reinforcement learning. Communications in Transportation Research, 1. <http://dx.doi.org/10.1016/j.commtr.2021.100017>

N.B. When citing this work, cite the original published paper.



# Connected autonomous vehicles for improving mixed traffic efficiency in unsignalized intersections with deep reinforcement learning

Bile Peng<sup>a,\*</sup>, Musa Furkan Keskin<sup>b</sup>, Balázs Kulcsár<sup>b</sup>, Henk Wymeersch<sup>b</sup>

<sup>a</sup> Institute for Communications Technology, TU Braunschweig, 38 106, Braunschweig, Germany

<sup>b</sup> Department of Electrical Engineering, Chalmers University of Technology, 41 296, Gothenburg, Sweden

## ARTICLE INFO

### Keywords:

Connected vehicles  
Autonomous driving  
Intelligent transportation systems  
Deep reinforcement learning

## ABSTRACT

Human driven vehicles (HDVs) with selfish objectives cause low traffic efficiency in an un-signalized intersection. On the other hand, autonomous vehicles can overcome this inefficiency through perfect coordination. In this paper, we propose an intermediate solution, where we use vehicular communication and a small number of autonomous vehicles to improve the transportation system efficiency in such intersections. In our solution, two connected autonomous vehicles (CAVs) lead multiple HDVs in a double-lane intersection in order to avoid congestion in front of the intersection. The CAVs are able to communicate and coordinate their behavior, which is controlled by a deep reinforcement learning (DRL) agent. We design an altruistic reward function which enables CAVs to adjust their velocities flexibly in order to avoid queuing in front of the intersection. The proximal policy optimization (PPO) algorithm is applied to train the policy and the generalized advantage estimation (GAE) is used to estimate state values. Training results show that two CAVs are able to achieve significantly better traffic efficiency compared to similar scenarios without and with one altruistic autonomous vehicle.

## 1. Introduction

The connected autonomous vehicle (CAV) with the ability of communication, coordination and autonomous driving has shown significant potentials of improving transportation systems (Rios-Torres and Malikopoulos, 2016; Li et al., 2017). In recent years, it has been applied to longitudinal velocity and lane changing maneuvers (Bahram, 2017), car-following (Wei et al., 2019), traffic smoothing (Kamal et al., 2016; Keskin et al., 2020), bottleneck decongestion (Vinitsky et al., 2018b), roundabout (Zhao et al., 2018), unsignalized intersection with only CAVs (Ahn and Del Vecchio, 2016; Azimi et al., 2014; Guney and Raptis, 2020; Campos et al., 2014; Hafner et al., 2013; Zhang et al., 2017) and un-signalized intersection with mixed traffics (both CAVs and human driven vehicles (HDVs)) (Vinitsky et al., 2018a). The objective of traffic efficiency optimization is realized by combinatorial optimization (Bahram, 2017), analytical control (Zhao et al., 2018; Malikopoulos et al., 2018; Hafner et al., 2013; Zhang et al., 2017), mixed integer linear programming (Ahn and Del Vecchio, 2016), protocol-based control (Azimi et al., 2014), scheduling-based optimization (Guney and Raptis, 2020), model predictive control (MPC) (Kamal et al., 2016; Keskin et al., 2020; Wei et al., 2019) and deep reinforcement learning (DRL) (Vinitsky et al., 2018a,b). DRL makes a series of decisions based on the observation

of environment states, in order to maximize the expected long-term reward, and has been applied to games (Silver et al., 2018), robotics (Kober et al., 2013) and resource allocation for wireless networks (Ye et al., 2019; Lee et al., 2019). Relevant to our context, DRL has also been applied to CAV control in a signalized intersection (Yang et al., 2017), an unsignalized single-lane intersection with desired speed as objective (Kheterpal et al., 2018) lane changing behavior (Wang et al., 2018), roundabout (Jang et al., 2019).

In this paper, we consider the application of CAVs in mixed traffics, i.e., the traffic flow comprises both CAVs and HDVs, which is more plausible in the foreseeable future than pure CAV traffic flow. Unlike autonomous vehicles that maximize its own benefit (e.g., safety, efficiency and comfort), CAVs with altruistic objectives influence adjacent HDVs with their behavior in order to optimize the overall traffic flow. In particular, we consider a two-lane unsignalized intersection with mixed traffic. Two fleets are placed on the two lanes and each fleet is led by a CAV. The objective is to avoid queuing in front of the intersection by cooperation of the two CAVs. Since there are two CAVs leading the two fleets, cooperation between the two CAVs is required to achieve the objective and the approaches in the existing literature (e.g., desired speed) do not apply to this scenario. In this study, we expand the figure-eight scenario presented in (Vinitsky et al., 2018a) to two lanes and the

\* Corresponding author.

E-mail address: [peng@ifn.ing.tu-bs.de](mailto:peng@ifn.ing.tu-bs.de) (B. Peng).

coordinated behaviors of the two CAVs are evaluated with respect to the overall traffic efficiency. To the authors' best knowledge, multi-lane traffic optimization with cooperation between the CAVs in an unsignalized intersection with mixed traffic has not been studied before. We expect that this work would be another step towards efficient mixed-autonomy urban traffic optimized by CAVs with altruistic control strategy. DRL is applied to determine the CAVs' behavior. We show that DRL works for our scenario as well with properly defined state, action and reward that eliminate local optima during the training.

## 2. Scenario description

We consider a two-lane unsignalized intersection (one lane in each direction) as shown in Fig. 1, which appears frequently in reality. In order to make the scenario close such that vehicles stay in the scenario and pass through the intersection again and again, entries and exits are connected with ring-shaped roads, such that vehicles can cross the intersection multiple times.

Without traffic light, the human driving behavior would result in a low efficiency, as will be shown later. In this paper, we aim to improve the traffic efficiency with CAVs powered by the DRL algorithm. Unlike the "selfish" autonomous vehicles, the objective of the CAVs is to optimize the whole traffic flow instead of their own benefit. Two CAVs are available in the considered scenario. They are assumed to be able to communicate with each other and coordinate their behaviors, which are shown as the red vehicles in Fig. 1. The HDVs are depicted as white vehicles, which imitate the human driving behavior with the intelligent driver model (IDM) Treiber et al. (2000).

The scenario, vehicle dynamics, and the reinforcement learning (RL) algorithms are implemented based on the Flow project Wu et al. (2017). In our configuration, the Flow project uses simulation of urban mobility (SUMO) Behrisch et al. (2011) as traffic simulator and RLlib Liang et al. (2017) as DRL algorithm implementation.

## 3. Algorithm

### 3.1. Review of RL

DRL addresses the problem of optimal actions in a dynamic environment with the objective of maximizing the cumulative reward of all time steps. Formally, the RL problem can therefore be formulated as

$$\begin{aligned} \max_{\theta} \quad & \mathbb{E} \left( \sum_{t=0}^H \gamma^t r_t \right) \\ \text{subject to} \quad & r_t = r(s_t, a_t), \\ & s_{t+1} = f(s_t, a_t), \\ & a_t = \pi_{\theta}(s_t), \end{aligned} \quad (1)$$

where  $\gamma \in [0, 1]$  is a discounting factor,  $r_t$  is the reward at time step  $t$  and is decided by state  $s_t$  and action  $a_t$  at the same time step (the first constraint), action  $a_t$  updates state  $s_t$  to state  $s_{t+1}$  at the next time step (the second constraint),  $a_t$  is decided by policy  $\pi$  parameterized by  $\theta$  given state  $s_t$  and  $H$  is the horizon, i.e., the number of time steps in an episode.

If reward function and system dynamics are unknown to the agent, the problem is model-free RL. We choose the state-of-the-art stochastic policy optimization algorithms proximal policy optimization (PPO) and value estimation algorithm generalized advantage estimation (GAE) for actor and critic, respectively. The reason to choose PPO is that it realizes a stable training process with the trust region and requires reasonable computation effort. GAE is applied because of its flexible bias-variance tradeoff.

PPO Schulman et al. (2017) optimizes the policy in an iterative manner. In each iteration, data samples are collected with the old policy and the new policy is improved in the proximal region of the old policy (called trust region). It is proved that the expected reward is guaranteed

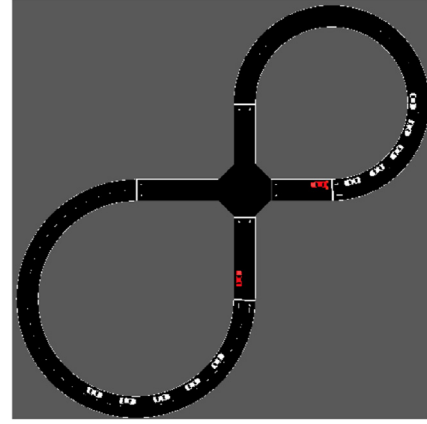


Fig. 1. Considered scenario and the initial vehicle positions. Red vehicles are CAVs whereas white ones are HDVs. We denote the bottom red car CAV 0 and the right red car CAV 1. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

to improve if the Kullback-Leibler (KL) divergence between the old and new policies is smaller than a certain threshold Schulman et al. (2015a). Based on this, trust region policy optimization (TRPO) Schulman et al. (2015a) uses Lagrange multiplier to optimize the policy subject to the KL divergence constraint to realize a stable training. On the other hand, PPO introduces a loss function, which has a nonzero gradient only in the proximal region and hence confines the optimization region. It achieves similar performance with lower computational complexity and is therefore chosen in this paper.

The GAE Schulman et al. (2015b) is applied as critic to estimate state value. The state value  $V^{\pi}(s)$  is the expected (discounted) sum of rewards given current state  $s$  and assuming the agent behaves according to policy  $\pi$ . Furthermore, the Q-value  $Q^{\pi}(s, a)$  is the expected (discounted) sum of rewards given current state  $s$  and corresponding action  $a$  and assuming the agent behaves according to policy  $\pi$  thereafter. The difference between  $Q^{\pi}(s, a)$  and  $V^{\pi}(s)$  is defined as advantage  $A^{\pi}(s, a)$ , which is used to compare action  $a$  and policy  $\pi$ . In this paper, we apply GAE Schulman et al. (2015b) to estimate the advantage.

### 3.2. Definitions of state and action

The state definition should contain all information needed to make the optimal decision. Besides, it should be directly correlated with the reward, making the learning process easier. Intuitively, the CAVs need their own positions and velocities as well as positions and velocities of vehicles in front of the intersection in both entries (from bottom and from right in Fig. 1). Other vehicles are far away from the intersection and have little impact on the CAVs' decision. Based on these considerations, we define the state as

$$s = \{x_0^A, x_1^A, v_0^A, v_1^A, x_{\text{first}}^B, x_{\text{last}}^B, v_{\text{first}}^B, v_{\text{last}}^B, x_{\text{first}}^R, x_{\text{last}}^R, v_{\text{first}}^R, v_{\text{last}}^R\} \quad (2)$$

where  $x_a^A$  and  $v_a^A$  are longitudinal coordinate and velocity of CAV  $a$ , respectively, with  $a \in \{0, 1\}$ ;  $x_m^B$  and  $v_m^B$  are longitudinal coordinate and velocity of vehicle  $m \in \{\text{first}, \text{last}\}$  within 30 m on the bottom entry of the intersection, respectively, with  $m = \text{first}$ , the first vehicle (closest to the intersection) and  $m = \text{last}$ , the last vehicle (farthest to the intersection within 30 m);  $x_n^R$  and  $v_n^R$  are state description of the right entry with similar meaning to  $x_m^B$  and  $v_m^B$ , respectively. If there is no vehicle within 30 m in front of the intersection, the coordinate is set to 30 m in front of the intersection and the velocity is set to 0. In this way, the state includes information on the CAVs (with superscript "A") and on vehicles approaching the two entries of the intersection (with superscript "B") for

bottom and “R” for right). Hence the scenario state is represented in a compact manner.

The action is the accelerations of the two CAVs. Therefore, the action space is a two dimensional continuous space  $[-3 \text{ m/s}^2, 3 \text{ m/s}^2]^2$  considering the realistic acceleration and deceleration.

### 3.3. Reward shaping

The objective is to reduce the queues in front of the intersection, similar to the one-lane, one CAV scenario from Vinitzky et al. (2018a). We cannot apply the same reward as Vinitzky et al. (2018a), to encourage vehicles to approach a desired velocity  $v_{\text{des}}$  because the initial distances from the two CAVs to the intersection are similar (see Fig. 1). A CAV needs to adapt its velocity such that there is enough time difference between the other CAV and itself to pass the intersection in order to avoid congestion. Once the proper time difference is created, the CAVs should maintain the time difference such that there is no congestion every time thereafter they pass through the intersection. Therefore, constant velocities can not realize the goal to avoid congestion. Based on the above consideration, the reward is defined as sum of the three terms with the following objectives: (i) maintaining desired velocity; (ii) minimizing queues; (iii) avoiding local optima. Based on this, we proposed the following rewards:

- 1) CAV 0 should keep its desired velocity. Therefore, the first term is defined as

$$r_1 = -\|v_0^A - v_{\text{des}}\| \quad (3)$$

where  $\|\cdot\|$  is the Euclidean distance operator. CAV 1 does not have a desired velocity.

- 2) The number of queuing vehicles should be minimized. A slow threshold velocity  $v_{\text{th}}$  is defined and vehicles moving slower than the threshold velocity are considered as queuing vehicles. The second term of the reward is thus defined as

$$r_2 = -|\{i | v_i < v_{\text{th}}\}| \quad (4)$$

where  $v_i$  is the velocity of vehicle  $i$  (the vehicle is either a CAV or an HDV, the set  $\{i | v_i < v_{\text{th}}\}$  is therefore the set of vehicles that has a speed less than  $v_{\text{th}}$ ) and  $|\cdot|$  is the cardinality operator.

- 3) Training according to  $r_1$  and  $r_2$  tends to suffer from a local optimum,<sup>1</sup> where CAV 1 tries to accelerate to its maximum velocity in order to reduce the queuing time of the fleet led by CAV 0 rather than slowing down to let the fleet led by CAV 0 pass the intersection first. To avoid this local optimum, we add a small correction term to discourage high velocities:

$$r_3 = -\|\mathbf{v}\| \quad (5)$$

where  $\mathbf{v}$  is the vector of all vehicles' velocities.

The total reward is defined as

$$r = \alpha r_1 + \beta r_2 + \kappa r_3, \quad (6)$$

where  $\alpha$ ,  $\beta$ , and  $\kappa$  are non-negative tuning parameters.

### 4. Selected results

The training process and testing results are presented in this section. We choose ADAM as the optimizer with the learning rate of  $5 \times 10^{-4}$ , the

<sup>1</sup> Due to the expectation operation of PPO, convergence to poor local optima is possible.

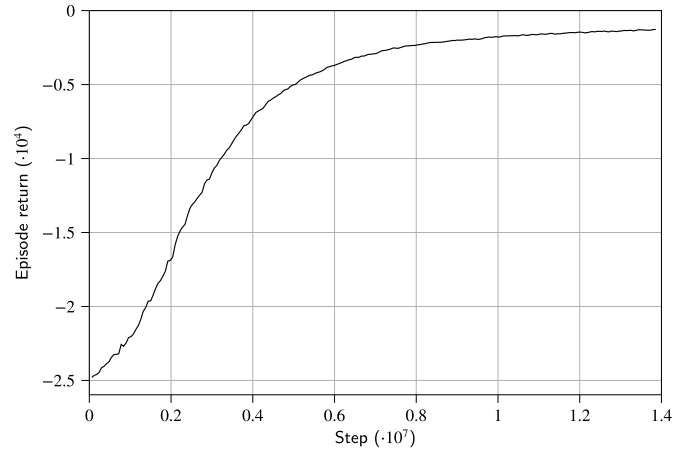


Fig. 2. Episode return in training showing monotonic improvement due to the selected reward shaping.

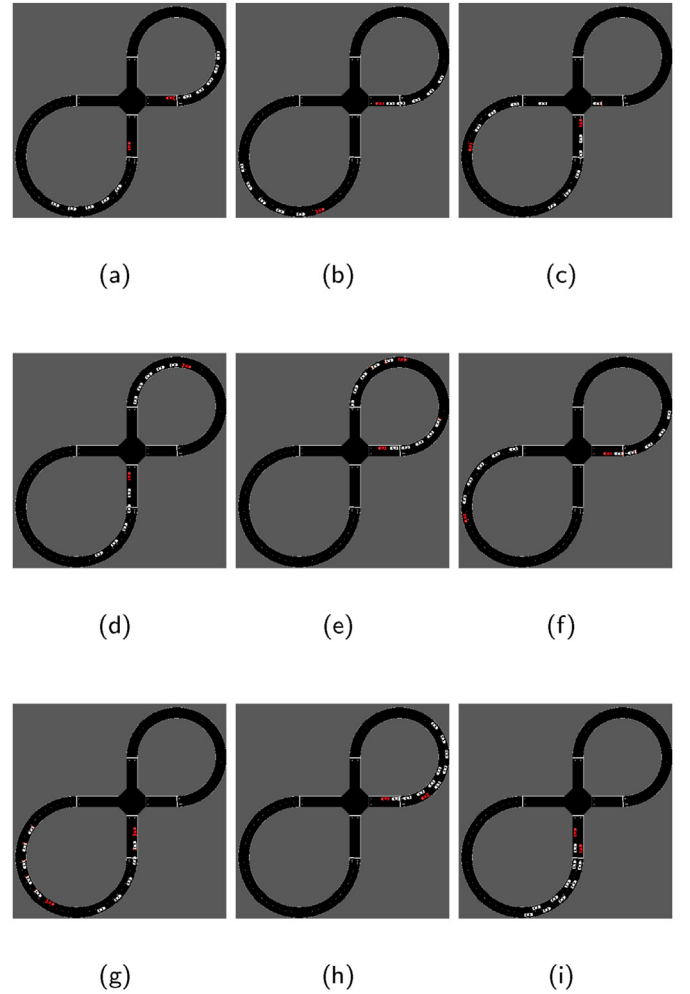


Fig. 3. Nine screenshots of the testing results when one CAV is approaching the intersection. It can be observed that the two CAVs gradually synchronize their longitudinal positions to avoid congestion.

clip of PPO  $\epsilon$  is set to 0.2. The discounting factor  $\gamma$  is set to 0.999 where as the tradeoff coefficient of GAE  $\delta$  is 0.99. In each iteration, 60 rollouts are run and the policy is optimized based on the sampled data. In our setup, desired velocity  $v_{\text{des}}$  is defined as 3 m/s, threshold velocity  $v_{\text{th}}$  is chosen as 1.1 m/s. Factors  $\alpha$ ,  $\beta$ , and  $\kappa$  are set to 10, 1, and 0.1, respectively.

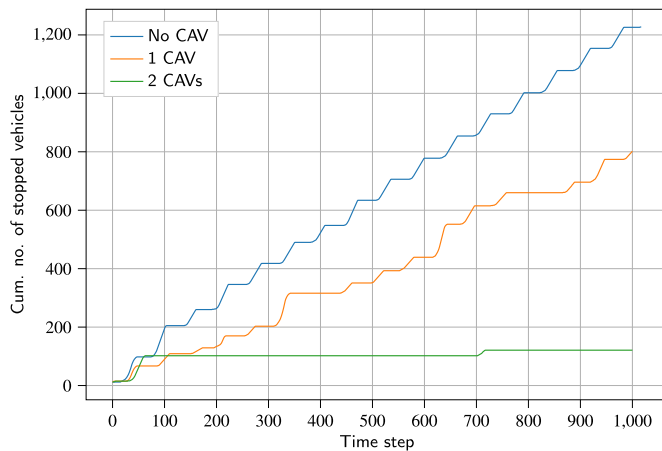


Fig. 4. Cumulative numbers of stopped vehicles in front of the intersection.

Fig. 2 shows the episode reward (sum of rewards in an episode) during training. It can be observed that the performance has been improved constantly during training. It was observed that without the  $r_3$  term, the return was much lower (results not shown).

Fig. 3 shows nine screenshots of the testing result when one CAV is approaching the intersection. While CAV 0 remains a constant velocity, CAV 1 adjusts its velocity to reach similar longitudinal position of CAV 0 and thus avoid congestion of the intersection.

Fig. 4 shows the cumulative numbers of stopped vehicles (numbers stopped vehicles up to that time) during the testing without CAV, where we define vehicles slower than 0.01 m/s as stopped, with 1 CAV and with 2 CAVs. The unique differences between the three setups are the numbers of CAVs. It is clear that 2 CAVs significantly reduce the stopped vehicles in front of the intersection and improve the traffic efficiency.

## 5. Conclusions

In this paper, DRL is applied to CAVs, which have an altruistic objective to optimize the traffic flow and the ability to communicate and coordinate their behaviors. The proposed algorithm uses the actor-critic framework with PPO as policy optimizer and GAE as value estimator. After the training, two CAVs are able to lead multiple HDVs in a closed scenario with an intersection. One CAV keeps a constant velocity whereas the other CAV adjusts its velocity to avoid congestion. A designed reward function is able to avoid local optima such that the policy converges in the desired behavior. Testing results show that the number of waiting vehicles in front of the intersection is considerably less with 2 CAVs than without and with 1 CAV. This work proves the potential of the altruistic CAVs in improving traffic efficiency. In future works, more general traffic behaviors (e.g., lane changing and turning) and situations (e.g., different numbers of HDVs) should be considered in order to make the approach closer to reality.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

The authors would like to thank Mr A. Kreidieh and Mr E. Vinitsky for their insightful suggestions.

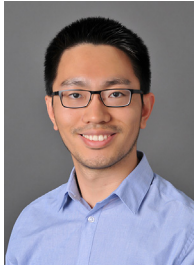
The project has been partially funded by Chalmers Transport Area of Advance under IRIS: Inverse Reinforcement-Learning and Intelligent Swarm Algorithms for Resilient Transportation Networks.

## References

- Ahn, H., Del Vecchio, D., 2016. Semi-autonomous intersection collision avoidance through job-shop scheduling. In: *Proceedings of the 19th International Conference on Hybrid Systems: Computation and Control*, pp. 185–194.
- Azimi, R., Bhatia, G., Rajkumar, R.R., Mudalige, P., 2014. Stip: spatio-temporal intersection protocols for autonomous vehicles. In: *2014 ACM/IEEE International Conference on Cyber-Physical Systems (ICCPs)*. IEEE, pp. 1–12.
- Bahram, M., 2017. Interactive Maneuver Prediction and Planning for Highly Automated Driving Functions. PhD thesis. Technische Universität München.
- Behrisch, M., Bieker, L., Erdmann, J., Krajzewicz, D., 2011. SUMO-simulation of urban MObility - an overview. In: *SIMUL 2011, the Third International Conference on Advances in System Simulation*, pp. 55–60.
- Campos, G.R., Falcone, P., Wymeersch, H., Hult, R., Sjöberg, J., 2014. Cooperative receding horizon conflict resolution at traffic intersections. In: *53rd IEEE Conference on Decision and Control*. IEEE, pp. 2932–2937.
- Guney, M.A., Raptis, I.A., 2020. Scheduling-based optimization for motion coordination of autonomous vehicles at multilane intersections. *Journal of Robotics* 2020. Article ID 6217409.
- Hafner, M.R., Cunningham, D., Caminiti, L., Del Vecchio, D., 2013. Cooperative collision avoidance at intersections: algorithms and experiments. *IEEE Trans. Intell. Transport. Syst.* 14 (3), 1162–1175.
- Jang, K., Vinitsky, E., Chalaki, B., Remer, B., Beaver, L., Malikopoulos, A.A., Bayen, A., 2019. Simulation to scaled city: zero-shot policy transfer for traffic control via autonomous vehicles. In: *Proceedings of the 10th ACM/IEEE International Conference on Cyber-Physical Systems*, pp. 291–300.
- Kamal, M.A.S., Taguchi, S., Yoshimura, T., 2016. Efficient driving on multilane roads under a connected vehicle environment. *IEEE Trans. Intell. Transport. Syst.* 17 (9), 2541–2551.
- Keskin, M.F., Peng, B., Kulcar, B., Wymeersch, H., 2020. Altruistic control of connected automated vehicles in mixed-autonomy multi-lane highway traffic. In: *2020 21st IFAC World Congress* (accepted).
- Kheterpal, N., Vinitsky, E., Wu, C., Kreidieh, A., Jang, K., Parvate, K., Bayen, A., 2018. Flow: Open Source Reinforcement Learning for Traffic Control.
- Kober, J., Bagnell, J.A., Peters, J., 2013. Reinforcement learning in robotics: a survey. *Int. J. Robot. Res.* 32 (11), 1238–1274.
- Lee, M., Yu, G., Li, Y., 2019. Learning to Branch: Accelerating Resource Allocation in Wireless Networks. *IEEE Transactions on Vehicular Technology*.
- Li, S.E., Zheng, Y., Li, K., Wang, L.-Y., Zhang, H., 2017. Platoon Control of Connected Vehicles from a Networked Control Perspective: Literature Review, Component Modeling, and Controller Synthesis. *IEEE Transactions on Vehicular Technology*.
- Liang, E., Liaw, R., Moritz, P., Nishihara, R., Fox, R., Goldberg, K., Gonzalez, J.E., Jordan, M.I., Stoica, I., 2017. Rllib: Abstractions for Distributed Reinforcement Learning arXiv preprint arXiv:1712.09381.
- Malikopoulos, A.A., Cassandras, C.G., Zhang, Y.J., 2018. A decentralized energy-optimal control framework for connected automated vehicles at signal-free intersections. *Automatica* 93, 244–256.
- Rios-Torres, J., Malikopoulos, A.A., 2016. A survey on the coordination of connected and automated vehicles at intersections and merging at highway on-ramps. *IEEE Trans. Intell. Transport. Syst.* 18 (5), 1066–1077.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., Moritz, P., 2015a. Trust region policy optimization. In: *International Conference on Machine Learning*, pp. 1889–1897.
- Schulman, J., Moritz, P., Levine, S., Jordan, M., Abbeel, P., 2015b. High-dimensional Continuous Control Using Generalized Advantage Estimation arXiv preprint arXiv: 1506.02438.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O., 2017. Proximal Policy Optimization Algorithms arXiv preprint arXiv:1707.06347.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al., 2018. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science* 362 (6419), 1140–1144.
- Treiber, M., Hennecke, A., Helbing, D., 2000. Congested traffic states in empirical observations and microscopic simulations. *Phys. Rev. E* 62 (2), 1805.
- Vinitsky, E., Kreidieh, A., Le Flem, L., Kheterpal, N., Jang, K., Wu, C., Wu, F., Liaw, R., Liang, E., Bayen, A.M., 2018a. Benchmarks for reinforcement learning in mixed-autonomy traffic. In: *Conference on Robot Learning*, pp. 399–409.
- Vinitsky, E., Parvate, K., Kreidieh, A., Wu, C., Bayen, A., 2018b. Lagrangian control through deep-RL: applications to bottleneck decongestion. In: *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, pp. 759–765.
- Wang, P., Chan, C.-Y., de La Fortelle, A., 2018. A reinforcement learning based approach for automated lane change maneuvers. In: *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, pp. 1379–1384.
- Wei, S., Zou, Y., Zhang, X., Zhang, T., Li, X., 2019. An integrated longitudinal and lateral vehicle following control system with radar and vehicle-to-vehicle communication. *IEEE Trans. Veh. Technol.* 68 (2), 1116–1127.
- Wu, C., Kreidieh, A., Parvate, K., Vinitsky, E., Bayen, A.M., 2017. Flow: Architecture and Benchmarking for Reinforcement Learning in Traffic Control arXiv preprint arXiv: 1710.05465.
- Yang, K., Tan, I., Menendez, M., 2017. A reinforcement learning based traffic signal control algorithm in a connected vehicle environment. In: *17th Swiss Transport Research Conference (STRC 2017)*. STRC.



- Ye, H., Li, G.Y., Juang, B.-H.F., 2019. Deep reinforcement learning based resource allocation for v2v communications. *IEEE Trans. Veh. Technol.* 68 (4), 3163–3173.
- Zhang, Y., Malikopoulos, A.A., Cassandras, C.G., 2017. Decentralized optimal control for connected automated vehicles at intersections including left and right turns. In: 2017 IEEE 56th Annual Conference on Decision and Control (CDC). IEEE, pp. 4428–4433.
- Zhao, L., Malikopoulos, A., Rios-Torres, J., 2018. Optimal control of connected and automated vehicles at roundabouts: an investigation in a mixed-traffic environment. *IFAC-PapersOnLine* 51 (9), 73–78.



**Bile Peng** received the B.S. degree from Tongji University, Shanghai, China, in 2009, the M.S. degree from the Technische Universität Braunschweig, Germany, in 2012, and the Ph.D. degree with distinction from the Institut für Nachrichtentechnik, Technische Universität Braunschweig in 2018. He has been a Postdoctoral researcher in the Chalmers University of Technology, Sweden from 2018 to 2019, a development engineer at IAV GmbH, Germany from 2019 to 2020. Currently, he is a Postdoctoral researcher in the Technische Universität Braunschweig, Germany. His research interests include wireless channel modeling and estimation, Bayesian inference as well as machine learning algorithms, in particular deep reinforcement learning, for resource allocation of wireless communication. Dr. Peng is a major contributor to the IEEE Standard for High Data Rate Wireless Multi-Media Networks Amendment 2: 100 Gb/s Wireless Switched Point-to-Point Physical Layer (IEEE Std 802.15.3d-2017) and received the IEEE vehicular technology society 2019 Neal Shepherd memorial best propagation paper award.



**Musa Furkan Keskin** is a researcher and a Marie Skłodowska-Curie Fellow (MSCA-IF) in the department of Electrical Engineering at Chalmers University of Technology, Gothenburg, Sweden. He obtained the B.S., M.S., and Ph.D degrees from the Department of Electrical and Electronics Engineering, Bilkent University, Ankara, Turkey, in 2010, 2012, and 2018, respectively. He received the 2019 IEEE Turkey Best Ph.D Thesis Award for his thesis on visible light positioning systems. His project "OTFS-RADCOM: A New Waveform for Joint Radar and Communications Beyond 5G" is granted by the European Commission through the H2020-MSCA-IF-2019 call. His current research interests include intelligent transportation systems, joint radar-communications, and positioning in 5G and beyond 5G systems.



**Balázs Kulcsár** received the M.Sc. degree in traffic engineering and the Ph.D. degree from Budapest University of Technology and Economics (BUTE), Budapest, Hungary, in 1999 and 2006, respectively. He has been a Researcher/Post-Doctor with the Department of Control for Transportation and Vehicle Systems, BUTE, the Department of Aerospace Engineering and Mechanics, University of Minnesota, Minneapolis, MN, USA, and with the Delft Center for Systems and Control, Delft University of Technology, Delft, The Netherlands. He is currently a Professor with the Department of Electrical Engineering, Chalmers University of Technology, Göteborg, Sweden. His main research interest focuses on traffic flow modeling and control.



**Henk Wymeersch** obtained the Ph.D. degree in Electrical Engineering / Applied Sciences in 2005 from Ghent University, Belgium. He is currently a Professor of Communication Systems with the Department of Electrical Engineering at Chalmers University of Technology, Sweden. He is also a Distinguished Research Associate with Eindhoven University of Technology. Prior to joining Chalmers, he was a postdoctoral researcher from 2005 until 2009 with the Laboratory for Information and Decision Systems at the Massachusetts Institute of Technology. Prof. Wymeersch served as Associate Editor for IEEE Communication Letters (2009–2013), IEEE Transactions on Wireless Communications (since 2013), and IEEE Transactions on Communications (2016–2018). During 2019–2021, he is a IEEE Distinguished Lecturer with the Vehicular Technology Society. His current research interests include the convergence of communication and sensing, in a 5G and Beyond 5G context.