
Deep Reinforcement Learning for Intelligent Transportation Systems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Intelligent Transportation Systems (ITSs) are envisioned to play a critical role in
2 improving traffic flow and reducing congestion, which is a pervasive issue impact-
3 ing urban areas around the globe. Rapidly advancing vehicular communication
4 and edge cloud computation technologies provide key enablers for smart traffic
5 management. However, operating viable real-time actuation mechanisms on a
6 practically relevant scale involves formidable challenges, e.g., Markov Decision
7 Processes (MDP) and conventional Reinforcement Learning (RL) techniques suffer
8 from poor scalability due to state space explosion. Motivated by these issues,
9 we explore the potential for Deep Q-Networks (DQN) to optimize traffic light
10 control policies. As an initial benchmark, we establish that the DQN algorithms
11 yield the “thresholding” policy in a single-intersection. Next, we examine the
12 scalability properties of DQN algorithms and their performance in a linear network
13 topology with several intersections along a main artery. We demonstrate that DQN
14 algorithms produce intelligent behavior, such as the emergence of “greenwave”
15 patterns, reflecting their ability to learn favorable traffic light actuations.

16 1 Introduction

17 Emerging Intelligent Transportation Systems (ITSs) [1–6] are expected to play an instrumental role
18 in improving traffic flow, thus optimizing fuel efficiency, reducing delays and enhancing the overall
19 driving experience. Today traffic congestion is an exceedingly complex and vexing issue faced by
20 metropolitan areas around the world. In particular, street intersections in dense urban traffic zones
21 (e.g., Times Square in Manhattan) can act as severe bottlenecks.

22 Current traffic light control policies typically involve preprogrammed cycles that may be optimized
23 based on historical data and adapted according to daily patterns. The options for adaptation to
24 real-time conditions, e.g. through detection wires in the pavement, tend to be fairly rudimentary.
25 Evolving vehicular communication technologies offer a crucial capability to obtain more fine-grained
26 knowledge of the positions and speeds of vehicles. Such comprehensive real-time information
27 can be leveraged, in conjunction with edge cloud computation, for significantly improving traffic
28 flow through more agile traffic light control policies, or in the longer term, via direct actuation
29 instructions for fully automated driving scenarios [7]. While the potential benefits are immense,
30 so are the technical challenges that evidently arise in solving such real-time actuation problems on
31 an unprecedented scale in terms of intrinsic complexity, geographic range, and number of objects
32 involved.

33 Under suitable assumptions, the problem of optimal dynamic traffic light control may be formulated
34 as a Markov decision process (MDP) [8–10]. The MDP framework provides a rigorous notion of
35 optimality along with a basis for computational techniques such as value iteration, policy iteration
36 or linear programming. However, an MDP formulation involves strong model assumptions, which

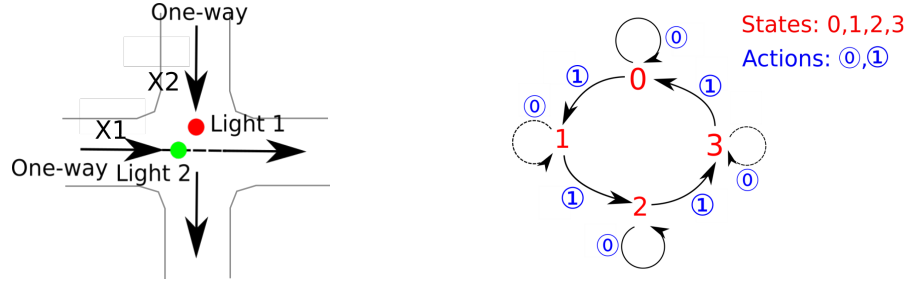


Figure 1: One intersection with two traffic flows (left) where X_1 and X_2 are the queue lengths, and the state transition diagram (right).

37 may not always be satisfied in reality, and knowledge of relevant system parameters, which may
 38 not be readily available. Owing to these issues, an MDP approach tends to be vulnerable to model
 39 mis-specification and inaccurate parameter estimation. Moreover, in terms of computational aspects,
 40 an MDP approach suffers from the curse of dimensionality, resulting in excessively large state spaces
 41 in realistic problem instances and exceedingly slow convergence.

42 Reinforcement Learning (RL) techniques, such as Q-learning, overcome some of these limitations
 43 [11–14] and have been previously considered in the context of optimal dynamic traffic light control
 44 [15–18]. However, conventional RL techniques are still prone to prohibitively large state spaces and
 45 extremely sluggish convergence, implying poor scalability beyond a single-intersection scenario.

46 Motivated by the above issues, we explore in the present paper the potential for deep learning
 47 algorithms, particularly Deep Q-Networks (DQN) [19], to optimize real-time traffic light control
 48 policies in large-scale transportation systems. As an initial validation benchmark, we analyze a
 49 single-intersection scenario and corroborate that the DQN algorithms match the provably optimal
 50 performance achieved by an MDP approach and exhibit a similar threshold structure. Next, we
 51 consider a linear network topology with several intersections to examine the scalability properties
 52 of DQN algorithms and their performance in the presence of highly complex interactions created
 53 by the flow of vehicles along the main artery. As mentioned above, the use of an MDP approach or
 54 standard RL techniques involves an excessive computational burden in these scenarios; hence the
 55 optimal achievable performance cannot be easily quantified. As a relevant qualitative feature, we
 56 demonstrate that DQN algorithms produce intelligent behavior, such as the emergence of “greenwave”
 57 patterns [20, 21], even though such structural features are not explicitly prescribed in the optimization
 58 process. This emergent intelligence confirms the capability of the DQN algorithms to learn favorable
 59 structural properties solely from observations.

60 The remainder of the paper is organized as follows. In Section 2, we present a detailed model
 61 description and problem statement. In Section 3, we provide a specification of the DQN algorithms
 62 for a single intersection as well as a linear network with several intersections. Section 4 discusses the
 63 computational experiments conducted to evaluate the performance of the proposed DQN algorithms
 64 and illustrate the emergence of “greenwave” patterns. In Section 5, we conclude with a few brief
 65 remarks and some suggestions for further research.

66 2 Model Description and Problem Statement

67 We model the road intersections and formulate our optimization problem. For the sake of transparency,
 68 we consider an admittedly stylized model that only aims to capture the most essential features that
 69 govern the dynamics of contending traffic flows at road intersections. We throughout adopt a discrete-
 70 time formulation to simplify the description and allow direct application of MDP techniques for
 71 comparison, but the methods and results naturally extend to continuous-time operation.

72 2.1 Single Road Intersection

73 As mentioned earlier, we start with a single-intersection scenario to facilitate the validation of the
 74 DQN algorithms by comparing it with the MDP approach. We consider the simplest meaningful
 75 setup with two intersecting unidirectional traffic flows as schematically depicted in the left side of
 76 Fig. 1. The state $S(t)$ of the system at the beginning of time slot t may be described by the three-tuple

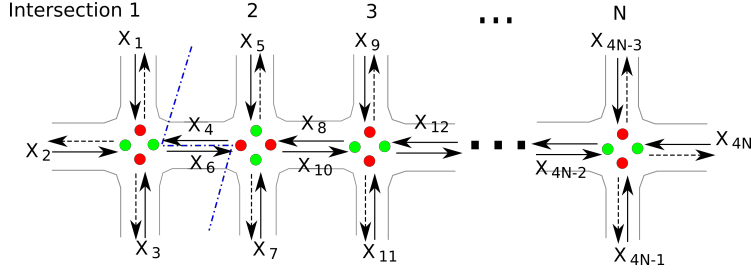


Figure 2: Linear bidirectional road network.

77 $(X_1(t), X_2(t); Y(t))$, with $X_i(t)$ denoting the number of vehicles of traffic flow i waiting to cross
 78 the intersection and $Y(t) \in \{0, 1, 2, 3\}$ indicating the configuration of the traffic lights:

- 79 • “0”: green light for direction 1 and hence red light for direction 2;
- 80 • “1”: yellow light for direction 1 and hence red light for direction 2;
- 81 • “2”: green light for direction 2 and hence red light for direction 1;
- 82 • “3”: yellow light for direction 2 and hence red light for direction 1.

83 Each configuration k can either simply be continued in the next time slot or must otherwise be
 84 switched to the natural subsequent configuration $(k + 1) \bmod 4$. This is determined by the action
 85 $A(t)$ selected at the end of time slot t , which is represented by a binary variable as follows: “0” for
 86 continue, and “1” for switch:

$$Y(t + 1) = (Y(t) + A(t)) \bmod 4. \quad (1)$$

87 These rules give rise to a strictly cyclic control sequence as illustrated in the right side of Fig. 1.

88 The evolution of the queue state over time is governed by the recursion

$$(X_1(t + 1), X_2(t + 1)) = (X_1(t) + C_1(t) - D_1(t), X_2(t) + C_2(t) - D_2(t)), \quad (2)$$

89 with $C_i(t)$ denoting the number of vehicles of traffic flow i appearing at the intersection during time
 90 slot t and $D_i(t)$ denoting the number of departing vehicles of traffic flow i crossing the intersection
 91 during time slot t . While not essential for our analysis, we make the simplifying assumption that if
 92 one of the two traffic flows is granted the green light, then exactly one waiting vehicle of that traffic
 93 flow, if any, will cross the intersection during that time slot, i.e.,

$$D_1(t) = \min\{1, X_1(t)\} \text{ if } Y(t) = 0; D_2(t) = \min\{1, X_2(t)\} \text{ if } Y(t) = 2; \quad (3)$$

94 and $D_1(t) = 0$ if $Y(t) \neq 0$ and $D_2(t) = 0$ if $Y(t) \neq 2$.

95 2.2 Linear Road Topology

96 To examine the performance and scalability properties of the DQN algorithms in more complex
 97 large-scale scenarios, we will consider a linear road topology. Specifically, we investigate a linear
 98 artery network topology with N intersections and bidirectional traffic flows, representing a main
 99 artery with cross streets as schematically depicted in Fig 2. We do not account for any traffic
 100 flows making left or right turns, but the analysis could easily be generalized to accommodate that.
 101 The state $S(t)$ of the system at the beginning of time slot t may be described by a $(5N)$ -tuple
 102 $(X_{n1}(t), X_{n2}(t), X_{n3}(t), X_{n4}(t); Y_n(t))_{n=1, \dots, N}$, with directions 1 and 2 corresponding to the east-
 103 west direction of the main artery and the north-south direction of the cross streets, and thus

$$Y_n(t + 1) = (Y_n(t) + A_n(t)) \bmod 4, \quad (4)$$

104 with $A_n(t)$ denoting the action selected for the n -th intersection at the end of time slot t .

105 The evolution of the various queue states is governed by the recursion

$$X_{ni}(t + 1) = X_{ni}(t) + C_{ni}(t) - D_{ni}(t), \quad (5)$$

106 with $C_{ni}(t)$ denoting the number of vehicles in direction i appearing at the n th intersection during
 107 time slot t and $D_{ni}(t)$ denoting the number of vehicles crossing the n -th intersection in direction i

108 during time slot t , $i = 1, \dots, 4$, $n = 1, \dots, N$. While $C_{11}(t)$, $C_{N3}(t)$, and $C_{n2}(t)$, $C_{n4}(t)$,
 109 $n = 1, \dots, N$, correspond to vehicles approaching the intersection from the external environment,
 110 we have $C_{n+1,1}(t+u) = D_{n1}(t)$ and $C_{n3}(t+u) = D_{n+1,3}(t)$, $n = 1, \dots, N-1$. This reflects
 111 that the vehicles crossing the n -th intersection in eastern direction during time slot t appear at the
 112 $(n+1)$ -th intersection u time slots later; likewise, vehicles passing through the $(n+1)$ -th intersection
 113 in western direction during time slot t arrive at the n -th intersection u time slots later. In this manner,
 114 the vehicles that travel along the main artery create highly complex interactions among the various
 115 intersections, which present additional challenges in optimizing the control policy.

116 Note that

$$D_{n1}(t) = \min\{1, X_{n1}(t)\}, D_{n3}(t) = \min\{1, X_{n3}(t)\}, \text{ if } Y_n(t) = 0, \quad (6)$$

117 $D_{n1}(t), D_{n3}(t) = 0$ if $Y_n(t) \neq 0$, and similarly for $D_{n2}(t)$ and $D_{n4}(t)$ depending on whether
 118 $Y_n(t) = 2$ or not.

119 2.3 Optimization Goal

120 We assume that the ‘‘congestion cost’’ in time slot t may be expressed as a function $F(X(t))$
 121 of the queue state, with $X(t) = (X_1(t), X_2(t))$ in the single-intersection scenario and $X(t) =$
 122 $(X_{n1}(t), X_{n2}(t), X_{n3}(t), X_{n4}(t))_{n=1, \dots, N}$ in the linear topology with N intersections. The goal is
 123 to find a dynamic control policy which selects actions over time so as to minimize the long-term
 124 expected discounted cost $\mathbb{E}[\sum_{t=1}^{\infty} \gamma^t F(X(t))]$, with $\gamma \in (0, 1)$ representing a discount factor.

125 3 Algorithm Design

126 We provide a detailed specification of the DQN algorithms for the scenarios of a single intersection
 127 or a linear topology with several intersections as described in the previous section.

128 First of all, let $Q(s, a)$ be the maximum achievable expected discounted reward (or minimum negative
 129 congestion cost in our context) under the optimal policy starting from state $s = (X; Y)$ when action a
 130 is taken. The $Q(s, a)$ values satisfy the equations

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, s'; a) \max_{a' \in \mathcal{A}} Q(s', a') = r(s, a) + \gamma \mathbb{E} \left[\max_{a' \in \mathcal{A}} Q(s', a') \right], \quad (7)$$

131 with $r(s, a) = F(X)$ denoting the congestion cost in queue state X , and $p(s, s'; a)$ denoting the
 132 transition probability from state s to state s' when action a is taken. Observe that the values
 133 $V(s) = \max_{a' \in \mathcal{A}} Q(s, a')$ satisfy the Bellman optimality equations

$$V(s) = \max_{a' \in \mathcal{A}} \{r(s, a') + \gamma \sum_{s' \in \mathcal{S}} p(s, s'; a') V(s')\} = \max_{a' \in \mathcal{A}} \{r(s, a') + \gamma \mathbb{E}[V(s')]\}. \quad (8)$$

134 The system state S serves as the input for both the target network and the evaluate network
 135 in the DQN algorithms, with $S = (X_1, X_2; Y)$ in the single-intersection scenario and $S =$
 136 $(X_{n1}, X_{n2}, X_{n3}, X_{n4}; Y_n)_{n=1, \dots, N}$ in the linear topology with N intersections. Equation (7) pro-
 137 vides the basis for deriving the target Q-values at each time step, while the Q-learning update for the
 138 neural network approximator in the i -th iteration is calculated based on

$$\text{Loss}(\theta_i) = \mathbb{E}_{s, a, r, s' \sim \text{memory}} \left[\left(r + \gamma \max (q_target(s', a'; \theta'_i)) - q_eval(s, a; \theta_i) \right)^2 \right], \quad (9)$$

139 where r is reward (negative cost) in the current step, s' and a' are the state and action in the next step,
 140 θ_i are parameters of the evaluate Q-network in the i -th iteration and θ'_i are parameters of the target
 141 Q-network with delayed update following the evaluate network.

142 The DQN algorithms sample from and train on data collected in memory. The online samples are
 143 stored in memory for further learning. A warm-up period of k_0 time steps is applied before the
 144 learning operations are initiated. The evaluate network is updated with the AdamOptimizer [22]
 145 gradient-descent and ϵ -greedy policy, whereas the update of the target network is slightly later.

Algorithm 1 DQN for single intersection or linear road topology with N intersections

- 1: Initialize queue and control states: either $X_1, X_2 = 0; Y = 0$ [single intersection] or $X_{n1}, X_{n2}, X_{n3}, X_{n4} = 0; Y_n = 0$ for all $n = 1, \dots, N$ [linear topology];
 - 2: For steps $k = 1, \dots, K$ do:
 - 3: $s = [X_1, X_2; Y]$ [single intersection] or $s = [X_{11}, \dots, X_{N4}; Y_1, \dots, Y_N]$ [linear topology];
 - 4: Select $a^* = \arg \max_{a \in \mathcal{A}} q_eval(s; a)$, using `eval_net` to evaluate the Q -value for each action;
 - 5: Generate random variables C_1, C_2 [single intersection] or C_{11}, C_{N3} and C_{n2}, C_{n4} for all $n = 1, \dots, N$;
 - 6: Given a^* , determine new queue and control states X'_1, X'_2, Y' according to Eq. (1)-(3) [single intersection] or $X'_{11}, \dots, X'_{N4}, Y'_1, \dots, Y'_N$ according to Eq. (4)-(6) [linear topology];
 - 7: $r = -((X'_1)^2 + (X'_2)^2)$ [single intersection] or $r = -\sum_{n=1}^N \sum_{i=1}^4 (X'_{ni})^2$ [linear topology];
 - 8: $s' = [X'_1, X'_2; Y']$ [single intersection] or $s' = [X'_{11}, \dots, X'_{N4}; Y'_1, \dots, Y'_N]$ [linear topology];
 - 9: Store transition $[s, a^*, r, s']$ in memory;
 - 10: Perform learning operation if $k > k_0$:
 - 11: Sample a minibatch of samples from memory;
 - 12: Update target network: $\theta'_i = \theta_i$;
 - 13: Calculate the target Q-value: $q_target(s, a^*) = r + \gamma \max_{a' \in \mathcal{A}} q_target(s', a')$;
 - 14: Update evaluate network with gradient descent (using AdamOptimizer) and ϵ -greedy policy: $Loss(\theta_i) = \mathbb{E}[(q_target - q_eval)^2]$.
-

146 Based on the above outline, we provide the specification of the DQN algorithm for the single-
 147 intersection scenario Fig. 1 and a linear topology with N intersections in Fig. 2. It is worth observing
 148 that even in the latter case we adopt a ‘single-agent’ DQN algorithm which has access to the global
 149 state of the network, as opposed to the ‘multiple-agent’ method with one agent for each individual
 150 intersection as considered in [17, 23]. While the single-agent approach involves a larger state space, it
 151 allows more intelligent control and coordination on a global level, which manifests itself for example
 152 in the emergence of greenwave patterns as we will demonstrate in the next section.

153 4 Performance Evaluation

154 We present simulations to evaluate the performance of the DQN algorithm in Alg. 1, and in particular
 155 illustrate the emergence of “greenwave” patterns in linear topology networks.

156 4.1 Single Road Intersection

157 As an initial validation benchmark, we first consider a single-intersection scenario as described in
 158 Subsection 2.1. The reason for considering this toy scenario is that the state space is sufficiently small
 159 for the optimal policy to be computed using the baseline MDP approach. We assume the numbers of
 160 arriving vehicles of both traffic flows in each time step as represented by the random variables C_1
 161 and C_2 to be independent and Bernoulli distributed with parameter $p = 1/4$. We use a quadratic
 162 congestion cost function $F(X_1, X_2) = X_1^2 + X_2^2$ and a discount factor $\gamma = 0.99$.

163 Inspection of the results in Fig. 3 shows that the DQN policy as obtained using Alg. 1 coincides with
 164 the optimal MDP policy. In particular, it matches the optimal performance and exhibits a similar
 165 threshold structure. This structural property was also reported in [24] for a strongly related two-queue
 166 dynamic optimization problem (with switch-over costs rather than switch-over times).

167 4.2 Linear Road Topology

168 We now turn to the scenario in Fig. 2. This is a more challenging scenario which serves to examine
 169 the scalability properties of our algorithm and its performance in the presence of highly complex
 170 interactions arising from the flow of vehicles along the main east-west arterial road.

171 Assume the numbers of externally arriving vehicles in eastern and western directions in each time
 172 step, represented by the random variables C_{11} and C_{N3} , to be independent and Bernoulli distributed

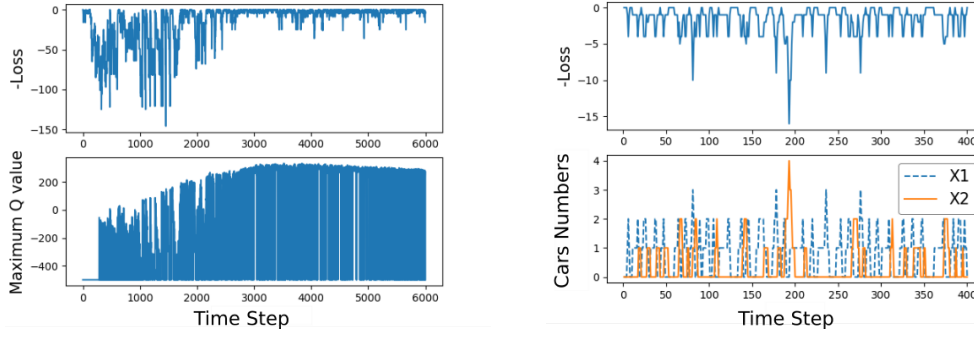


Figure 3: Learning curve for DQN policy (left) and thresholding property of DQN policy (right).

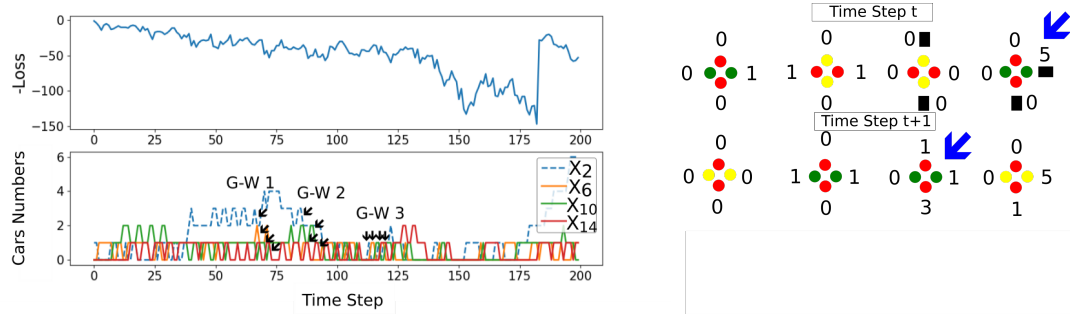


Figure 4: Learning curve for the linear network of size 4×1 (left) and Greenwave traffic lights (right). “Greenwave” is abbreviated to be “G-W”. In the right, numbers indicate the number of vehicles waiting on each road. Black rectangles indicate incoming vehicles from peripheral roads.

173 with parameter $p_1 = 1/4$. The numbers of arriving vehicles in southern and northern directions on
 174 each of the N cross streets in each time step, represented by the random variables C_{n2} and C_{n4} ,
 175 $n = 1, \dots, N$, are also independent and Bernoulli distributed with parameter $p_2 = 1/8$. We use a
 176 quadratic congestion cost function $F(X) = \sum_{n=1}^N \sum_{i=1}^4 X_{ni}^2$ and a discount factor $\gamma = 0.99$. In
 177 simulations, the evaluate and target networks used in Alg. 1 have both 4 fully-connected layers of
 178 size 200, 100, 40 and 2, respectively. We use ReLu as activation functions and squared difference
 179 loss.

180 The use of an MDP approach is computationally infeasible in this case due to the state space
 181 explosion, and hence the degree of optimality of our algorithm cannot be assessed in a quantitative
 182 manner. Instead we have therefore examined qualitative features to validate the intelligent behavior
 183 of our algorithm and evaluate its performance merit. In particular, we observed the emergence of
 184 “greenwave” patterns as shown in Fig. 4, even though such structural features are not explicitly
 185 prescribed in the optimization process. Specifically, the “greenwave” phenomenon is reflected as
 186 consecutive reduction of car numbers in each road. This emergent intelligence confirms the capability
 187 of our algorithm to learn favorable structural properties solely from observations.

188 5 Conclusion

189 We have explored the scope for Deep Q-Networks (DQN) to optimize real-time traffic light control
 190 policies in emerging large-scale Intelligent Transportation Systems. As an initial benchmark, we
 191 established that DQN algorithms deliver the optimal performance achieved by an MDP approach in a
 192 single-intersection scenario. We subsequently evaluated the scalability properties of DQN algorithms
 193 in a linear topology with several intersections, and demonstrated the emergence of intelligent behavior
 194 such as “greenwave” patterns, confirming their ability to learn desirable structural features. In future
 195 research we intend to investigate locality properties and analyze how these can be exploited in the
 196 design of distributed coordination schemes for wide-scale deployment scenarios.

References

- 197
- 198 [1] MM Vazifeh, P Santi, G Resta, SH Strogatz, and C Ratti, “Addressing the minimum fleet problem in
199 on-demand urban mobility,” *Nature*, vol. 557, no. 7706, pp. 534, 2018.
- 200 [2] Ming Zhu, Xiao-Yang Liu, Feilong Tang, Meikang Qiu, Ruimin Shen, Wei Wennie Shu, and Min-You Wu,
201 “Public vehicles for future urban transportation,” *IEEE Trans. Intelligent Transportation Systems*, vol. 17,
202 no. 12, pp. 3344–3353, 2016.
- 203 [3] Ming Zhu, Xiao-Yang Liu, and Xiaodong Wang, “An online ride-sharing path-planning strategy for public
204 vehicle systems,” *IEEE Transactions on Intelligent Transportation Systems*, 2018.
- 205 [4] Muhammad Alam, Joaquim Ferreira, and José Fonseca, “Introduction to intelligent transportation systems,”
206 in *Intelligent Transportation Systems*. Springer, 2016, pp. 1–17.
- 207 [5] Yisheng Lv, Yanjie Duan, Wenwen Kang, Zhengxi Li, Fei-Yue Wang, et al., “Traffic flow prediction with
208 big data: A deep learning approach,” *IEEE Trans. Intelligent Transportation Systems*, vol. 16, no. 2, pp.
209 865–873, 2015.
- 210 [6] Rajeshwari Sundar, Santhosh Hebbar, and Varaprasad Golla, “Implementing intelligent traffic control
211 system for congestion control, ambulance clearance, and stolen vehicle detection,” *IEEE Sensors Journal*,
212 vol. 15, no. 2, pp. 1109–1113, 2015.
- 213 [7] Liang Qi, MengChu Zhou, and WenJing Luan, “A two-level traffic light control strategy for preventing
214 incident-based urban traffic congestion,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 19,
215 no. 1, pp. 13–24, 2018.
- 216 [8] Simona Onori, Lorenzo Serraio, and Giorgio Rizzoni, “Dynamic programming,” in *Hybrid Electric*
217 *Vehicles*, pp. 41–49. Springer, 2016.
- 218 [9] Martin L Puterman, *Markov decision processes: discrete stochastic dynamic programming*, John Wiley &
219 Sons, 2014.
- 220 [10] Sheldon M Ross, *Introduction to stochastic dynamic programming*, Academic Press, 2014.
- 221 [11] Andrew Gehret Barto, Steven J Bradtke, and Satinder P Singh, *Real-time learning and control using*
222 *asynchronous dynamic programming*, University of Massachusetts at Amherst, Department of Computer
223 and Information Science, 1991.
- 224 [12] Andrew W Moore and Christopher G Atkeson, “Prioritized sweeping: Reinforcement learning with less
225 data and less time,” *Machine Learning*, vol. 13, no. 1, pp. 103–130, 1993.
- 226 [13] Richard S Sutton, “Learning to predict by the methods of temporal differences,” *Machine Learning*, vol. 3,
227 no. 1, pp. 9–44, 1988.
- 228 [14] Christopher JCH Watkins and Peter Dayan, “Q-learning,” *Machine Learning*, vol. 8, no. 3-4, pp. 279–292,
229 1992.
- 230 [15] Elise van der Pol, Frans A Oliehoek, T Bosse, and B Bredeweg, “Video demo: Deep reinforcement
231 learning for coordination in traffic light control,” *AmsterdamVrije Universiteit, Department of Computer*
232 *Sciences*, 2016.
- 233 [16] Elise Van der Pol and Frans A Oliehoek, “Coordinated deep reinforcement learners for traffic light control,”
234 *Proceedings of Learning, Inference and Control of Multi-Agent Systems (at NIPS 2016)*, 2016.
- 235 [17] MA Wiering, “Multi-agent reinforcement learning for traffic light control,” in *Proceedings of the*
236 *Seventeenth International Conference Machine Learning (ICML)*, 2000, pp. 1151–1158.
- 237 [18] Hua Wei, Guanjie Zheng, Huaxiu Yao, and Zhenhui Li, “Intellilight: A reinforcement learning approach
238 for intelligent traffic light control,” in *Proceedings of the 24th ACM SIGKDD International Conference on*
239 *Knowledge Discovery & Data Mining*. ACM, 2018, pp. 2496–2505.
- 240 [19] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
241 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al., “Human-level control
242 through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529, 2015.
- 243 [20] Stefan Lämmer and Dirk Helbing, “Self-control of traffic lights and vehicle flows in urban road networks,”
244 *Journal of Statistical Mechanics: Theory and Experiment*, vol. 4, pp. 04019, 2008.
- 245 [21] Stefan Lämmer and Dirk Helbing, “Self-stabilizing decentralized signal control of realistic, saturated
246 network traffic,” Santa Fe Institute, 2010.
- 247 [22] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint*
248 *arXiv:1412.6980*, 2014.
- 249 [23] Sergey Satunin and Eduard Babkin, “A multi-agent approach to intelligent transportation systems modeling
250 with combinatorial auctions,” *Expert Systems with Applications*, vol. 41, no. 15, pp. 6622–6633, 2014.
- 251 [24] Micha Hofri and Keith W Ross, “On the optimal control of two queues with server setup times and its
252 analysis,” *SIAM Journal on Computing*, vol. 16, no. 2, pp. 399–420, 1987.