

A Comparison Study of Different Data Resolutions for Deep Reinforcement Learning Based Adaptive Traffic Signal Control System

Mobin Yazdani, Majid Sarvi, Saeed Asadi Bagloee¹

¹Transport technologies group, Faculty of Engineering and Information Technology, University of Melbourne, Victoria

Email for correspondence: myazdani@student.unimelb.edu.au

Abstract

By recent advances in technology and Artificial Intelligence (AI), traffic signal control systems are preferably designed to have intelligence rather than rule-based structure. Deep Reinforcement Learning (RL) as a solution for sequential decision-making problems has been extensively used for adaptive traffic signal control (ATSC) systems. The deep RL-based ATSC systems have shown promising results versus current actuated (rule-based) ATSC systems. The conducted studies have employed different data resolutions either collected in vehicular network (e.g., location and speed of individual vehicles) or from camera devices (e.g., queue length, density, average speed) for proposed models. However, the impact of different data resolutions on deep RL-based ATSC systems training performance has not been studied yet. In this study, we compare the three different data resolutions in terms of computation time, training stability and results for variety of performance measurements. The Double Deep Q-Network (DDQN) algorithm is utilized as our intelligent agent. To test and evaluate the different data resolutions, a real isolated intersection is modelled in a simulation environment with real traffic volume demand. The experimental results have shown that vehicular network high resolution data can only contribute to a slight improvement versus camera data in terms of reduction in travel time, queue length etc. at the expense of more computation time in training models. Also, the camera data is more accessible compared to vehicular network data which needs sensors on plenty of vehicles in network. Hence, we recommend using camera data which provides aggregated but adequate data for deep RL-based ATSC models.

1.Introduction

Global population is increasing and Australian cities are more urbanized and congested (Audit, 2019). Traffic congestion causes significant costs in fuel consumption, emissions and traffic delays. One effective solution to address this issue is to design an efficient traffic signal control system. The system can switch and execute the traffic lights for flexible durations so to reduce the delay of transportation modes. Hence, traffic signal control systems play a pivotal role to reduce the traffic congestion costs.

With the advent of technology, intelligent transport system (ITS) devices such as camera, bluetooth, radar etc have provided practitioners with high resolution data to improve transport systems (Emami et al., 2020). Furthermore, the advanced Machine Learning (ML) algorithms have emerged as a powerful tool to improve systems's performance by intelligent agents instead of classical methods or human intervention (LeCun et al., 2015, Mnih et al., 2015).

44 Hence, ML-based methods have attracted attentions for transport applications. For traffic signal
45 systems, deep reinforcement learning (RL) is a solution to enable smart decision-making of
46 traffic lights selection only based on data in place. In fact, the traffic data collected by ITS
47 devices are fed to a deep neural network (DNN) which outputs traffic signal action values. The
48 traffic signal action selection is then taken by RL algorithm. Once the action is executed, RL
49 agent receives a scalar value (i.e., reward) that indicates how good was the selected action. As
50 ultimate goal in transportation is to reduce traffic congestion, the reward is defined to minimize
51 waiting time. By continuous interaction of deep RL agent with traffic environment, the agent
52 learns the optimal policy in action selection which leads to less traffic congestion.

53 The conducted studies in literature have shown superior performance of deep RL-based ATSC
54 over rule-based traffic signal control methods. The different variety of traffic data have been
55 used for proposed deep RL-based models. The traffic data are either assumed to be collected
56 in vehicular network (Wu et al., 2020) or from camera devices installed at intersections (Jeon
57 et al., 2018). The vehicular network data provides accurate information of individual vehicles
58 such as location and speed within the road while camera data can provide aggregated data such
59 as queue length, density and average speed of each incoming lane to the intersection. In
60 literature, vehicular network data has been extensively used as it brings information-dense data
61 samples for deep RL agent training. For real world implementation, however, it is critical to
62 examine if the traffic data is available for proposed model. Also, using the higher resolution
63 data should be justified considering the installation costs, amount of improvements and
64 computation costs. In this study, we compare and investigate the impact of different data
65 resolutions on deep RL-based ATSC performance. To represent the real world situation, we
66 test and evaluate the deep RL-based ATSC models in a microsimulation model of an
67 intersection in Melbourne, Australia. The experimental results have shown that vehicular
68 network data can only contribute to slight improvement versus camera data at the expense of
69 more training costs. Also, according to the easier accessibility of practitioners to camera data
70 versus vehicular network data, camera data is more preferable. Hence, we recommend using
71 camera data such as queue length and density for deep RL-based ATSC model.

72

73 **2. Related works**

74 Early studies in literature have applied RL on their ATSC problem (Abdoos et al., 2011,
75 Prashanth and Bhatnagar, 2010, El-Tantawy and Abdulhai, 2010, Balaji et al., 2010, Arel et
76 al., 2010, Grégoire et al., 2007, Wiering, Gao et al., 2017). The RL can learn from data collected
77 at intersection to train a smart traffic signal. However, the data used in these studies were coarse
78 and abstract because of RL agent's inability to capture problems with high dimensions. By
79 proposing Deep Q-Network (DQN) algorithm (Mnih et al., 2015) RL agents were enabled to
80 learn from high-dimensional data. Since then, the research conducted in ATSC literature aimed
81 at using data with more information such as queue length of the waiting vehicles, density or
82 speed and location of individual vehicles as presented in Table1. **Despite employing data with**
83 **different range of resolutions, the extent of improvements has not been investigated. In this**
84 **study we examine the effect of having different data resolutions in terms of computation costs,**
85 **training stability and improvements for traffic performance measurements.**

86

87

88

89

90

91

92
93
94
95
96

Table 1: Different data resolutions used in literature

Study	Data resolution	Number of intersections
(Genders and Razavi, 2016)	Speed and location of individual vehicles	1
(Van der Pol and Oliehoek, 2016)	Location of individual vehicles for state definition and speed of individual vehicles for reward	1,2,3,4
(Gao et al., 2017)	Speed and location of individual vehicles	1
(Mousavi et al., 2017)	Speed and location of individual vehicles	1
(Jeon et al., 2018)	The image of the intersection	1
(Wei et al., 2018)	Queue length, number of vehicles, updated waiting time of vehicles, vehicle's position	1
(Genders, 2019)	Queue length and density	1
(Genders and Razavi, 2019)	Queue length and density	1
(Liang et al., 2019)	Speed and location of individual vehicles	1
(Tan et al., 2019)	Queue length	24 (traffic grid)
(Wei et al., 2019)	Total number of vehicles and segment wise distribution of vehicles on each lane	4 (Arterial)
(Zheng et al., 2019)	Total number of vehicles on each lane	4 (Arterial)
(Wu et al., 2020)	Speed and location of individual vehicles, queue length and number of pedestrians	2, 6
(Wang et al., 2021)	Congestion level of the intersection	36

97
98

99 3. Deep Reinforcement Learning

100 Reinforcement learning (RL) is one of the machine learning paradigms that enables an agent
 101 to learn a specific task. The main components of RL are state (s), action (a), reward (r). In
 102 each time step t , the intelligent agent takes action a_t in state s_t , receives reward r_t and ends up
 103 in next state s_{t+1} . The r is a feedback signal for RL agent to learn the goal of the task. To
 104 explore the solution space, random actions are mostly taken at the beginning of the learning
 105 process. As RL agent is goal-oriented, unwanted actions are penalized with a reward. RL aims
 106 at maximizing the return $G_t \doteq \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ which is the cumulative of discounted rewards.
 107 The future rewards are discounted with parameter $\gamma \in [0,1)$ so to consider the importance of
 108 immediate rewards over future rewards. The most popular RL algorithm is Q-learning
 109 algorithm that each state-action pair value (Q-value) is the expected return under policy π as
 110 $Q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t | s_t = s, a_t = a]$. Based on Bellman optimality equation, the optimal Q-
 111 value under policy π is the one with best (i.e., maximum) value (Sutton and Barto, 2018). By
 112 continuous interaction of RL agent in the defined state, the RL can finally learn to take the
 113 optimal actions. To solve the problem with higher number of state-action pairs, the deep
 114 learning variant of Q-learning algorithm, Deep Q-Network, proposed by Mnih et al. (2015). In

115 this model, the deep neural network is used as a function approximator for Q-value estimations.
 116 The Q-function with optimal policy π^* is calculated as in equation1.
 117

$$Q^{\pi^*}(s, a; \theta) = \mathbb{E}_{\pi}[r_{t+1} + \gamma \max_a Q(s_{t+1}, a; \theta) | s_t = s, a_t = a] \quad \text{equation1}$$

118
 119 Where, θ is the deep neural network parameter.
 120

121 **4. Deep RL components**

122 In this study, the deep RL-based model performance is investigated with three different data
 123 resolutions as follows:

- 124 • ATSC-QLAS: lane queue length (QL) and link average speed (AS)
- 125 • ATSC-QLDS: lane queue length (QL) and lane density (DS)
- 126 • ATSC-DTSE: location and speed data for each individual vehicle popular as discrete
- 127 traffic state encoding (DTSE)
- 128
- 129

130 Our case study is a real four-leg-four-lane intersection at Rathdowne-Elgin junction in
 131 Melbourne, Australia. The simulation model has 4 zones {N,E,S,W} for origin and destination
 132 demands through which the vehicles are generated with dynamic assignment model. In this
 133 section the state, action and reward definition for models are explained.
 134

135 **4.1. state definition**

136 **4.1.1. ATSC-QLAS model**

137 The state definition for this model is queue length of the incoming lanes, **average speed** of
 138 incoming link and the current traffic signal state. Hence, the state space is $S \in \mathbb{R}^{i*j} \times$
 139 $\mathbb{R}^j \times \mathbb{T}^n$. Where, i and j are the number of intersection lanes and links respectively and, n is
 140 the number of agent actions. It is noted that queued vehicles are vehicles with speed less than
 141 5 km/hr.
 142

143 **4.1.2. ATSC-QLDS model**

144 The state definition for this model is queue length and density of vehicles in incoming lanes
 145 and, the current traffic signal state. Hence, the state space is $S \in (\mathbb{R} \times \mathbb{R})^{i*j} \times \mathbb{T}^n$.
 146

147 **4.1.3. ATSC-DTSE model**

148 In discrete traffic state encoding state definition, the individual speed and vehicles data is
 149 collected from incoming lanes. Each lane is discretized to multiple cells that fits the vehicles
 150 data as illustrated in Figure 1. The vehicle's average length in simulation is 5 meters so the cell
 151 length is selected 6 meters. The state space for this model is $S \in (\mathbb{B} \times \mathbb{R})^{(i*j)*k} \times \mathbb{T}^n$. Where,
 152 k is the number of cells along the link.
 153
 154
 155

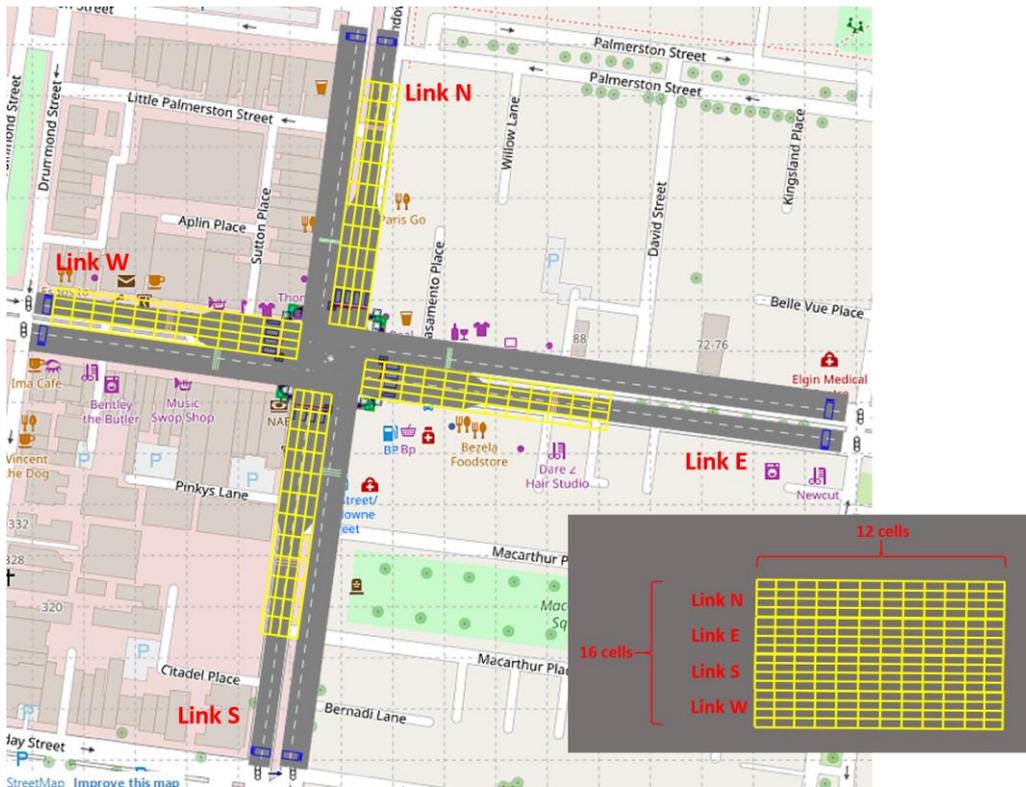


Figure 1: Discrete traffic state encoding for an intersection

156

157 **4.2. action definition**

158 The traffic light has two actions $a = \{NS, EW\}$. The RL agent either executes green light for
 159 North-South (NS stage) direction or East-West (EW stage) direction. Each action (i.e., stage)
 160 has two phases each with 5 seconds green time following the 3 seconds amber and 2 seconds
 161 red light as presented in Figure 2. The first phase includes through and left turn movement with
 162 permissive right turn movement. The next phase is a protected right turn movement, and it is
 163 only activated for transition to next action.



Figure 2: traffic signal phases of simulated intersection

164

165

166 **4.3. Reward definition**

167

168 **4.3.1. ATSC-QLAS model**

169 The reward definition must be complied with the resolution of available data. Hence, the
 170 reward function for this model is to minimize the queue length at the intersection. A vehicle is
 171 in queue if it has speed less than 5 km/hr. In each time step, the queue length data for each
 172 lane is collected as $r_t = -(q_{i,j})$. The reward for each action is the summation of rewards
 173 during the course of action execution.
 174

175 **4.3.2. ATSC-QLDS model**

176 The reward function for this model is similar to ATSC-QLAS which is penalizing the queued
 177 vehicles.
 178

179 **4.3.3. ATSC-DTSE model**

180 As vehicles speed and location data is available, reward function can be more accurate. The
 181 reward is to minimize the waiting time for each individual vehicle. In this study we assume
 182 that vehicles with speed less than 5 km/hr are in the queue and they are experiencing delay.
 183 To include the vehicles in range within [0,5) we define the reward function for each vehicle at
 184 time step t as $r_t = -(1 - 0.2v_t)$. Where vehicles with speed $v = 0$ km/hr have rewards -1
 185 and vehicles with speed $v = 5$ km/hr have reward 0. The vehicles waiting time within this
 186 range are considered with a linear function.
 187

187

188 **5. Deep RL model**

189 In this section the model specification and deep neural network architectures are explained.

190 **5.1. Double Deep Q-Network algorithm (DDQN)**

191 In this study the Double Deep Q-Network (DDQN) algorithm is used for training the agent.
 192 The algorithm leverages two networks, Primary network (θ^P) and Target network (θ^T) to
 193 update the Q-values. First, the Primary network get the states as input and outputs the Q-value
 194 of corresponding actions. As such, after each action execution, the experience memory tuple
 195 (s_t, a_t, r_t, s_{t+1}) is generated and is saved in a buffer memory with size \mathcal{B} . To slow down the
 196 learning fluctuations, Target network (θ^T) is used for target value estimation. The Target neural
 197 network is not trained. First, the (state, target) pairs are used for training the primary network
 198 and Target weights are then updated by soft update rate β as $\theta^T = \beta\theta^P + (1 - \beta)\theta^T$. To
 199 address the correlation between consecutive memory tuples, the experience tuples are
 200 uniformly sampled with minibatch size \mathcal{M} from a buffer memory β . The DDQN algorithm can
 201 address the maximization bias of the DQN algorithm as follows:

$$Q^{\pi^*}(s, a; \theta^P) = \mathbb{E}_{\pi} [r_{t+1} + \gamma Q(s_{t+1}, \underset{a_t}{\operatorname{argmax}} Q(s_{t+1}, a_t; \theta^P); \theta^T) | s_t = s, a_t = a] \quad \text{equation 2}$$

202

203 In this study, the action selection for DDQN agent is based on decaying ϵ -greedy strategy
 204 where enables agent to do the exploration-exploitation process. To clarify, the agent mostly
 205 takes random actions with probability ϵ_{init} in the beginning of training episodes. The ϵ is

206 decreased with decay ratio η until it reaches minimum epsilon ϵ_{min} . To learn the task, DDQN
 207 agent run is repeated for N number of episodes. In each episode (e) the model has a warm up
 208 period (W) before starting the process. The whole algorithm is presented in table 3. The DDQN
 209 ATSC model hyperparameters are as Table 2.
 210

Table 2: DDQN ATSC model hyperparameters

Variable	Hyperparameters	Value
γ	Discounted factor	0.95
α	Learning rate	0.0001
β	Target update rate	0.001
\mathcal{B}	Buffer memory size	200
\mathcal{M}	Mini batch size	32
ϵ_{min}	Minimum epsilon	0.01
ϵ_{init}	Initial epsilon	0.9
η	Decay ratio	0.05

211

212 **5.2. Deep neural network architecture**

213 **5.2.1. ATSC-QLAS model**

214 This model has a deep neural network with 4 layers. The first layer is a vector with size 16 for
 215 queue length input and another vector with size 4 for average speed input. The input data is fed
 216 to 3 fully connected layers with 64, 64 and 32 neurons respectively. The ReLU activation
 217 function is also applied after each layer. The last layer is the RL actions which is 2 in this study.
 218

219 **5.2.2. ATSC-QLDS model**

220 This model has a deep neural network with 3 layers. The first layer is a vector with size 16 for
 221 queue length input and another vector with size 16 for density input. The input data is fed to 2
 222 fully connected layers with 128 and 64 neurons respectively. The ReLU activation function is
 223 also applied after each layer. The last layer is the RL actions.
 224

225 **5.2.3. ATSC-DTSE model**

226 In this model, the state inputs are fed to the CNN with 5 layers. The first layer consists of two
 227 image-like input (speed and location) both with 16×12 dimensions. It is convolved with 16
 228 filters each with 4×4 size and stride of 2. Similarly, the next layer is convolution but with 32
 229 filters each with 2×2 size and stride of 1. Then the flattened images are fed to a fully connected
 230 layers first with 128 neurons and second with 64 neurons. It is noted that the ReLU activation
 231 function is applied after each layer. The last layer is the output of the deep CNN network with
 232 the number of RL actions. In this study, RL has 2 actions to select.

233 For all models, the RMSprop optimizer is used for deep neural networks (Tieleman and Hinton,
 234 2012) to minimize the loss function mean squared error (MSE).
 235
 236
 237
 238

239

Algorithm 1 Double Deep Q-Network Pseudocode for Adaptive Traffic Signal Control

Initialize: Primary Network $Q_P(s, a)$ with random weights θ^P
Initialize: Target Network $Q_T(s, a)$ with random weights θ^T
Initialize: Experience replay buffer memory \mathcal{B} for DDQN agent

- 1: **for** $e = 1, 2, \dots, N$ **do**
- 2: **for** $step = 1, 2, \dots, W$ **do**
- 3: Initialize traffic state S_t matrices
- 4: Start Exploration-Exploitation strategy
 decay = $N \times \eta$
 $\epsilon = 1 - e/\text{decay}$
 $\epsilon \leftarrow \epsilon \times \epsilon_{init}$
 $\epsilon \leftarrow clip(\epsilon, \epsilon_{min}, \epsilon_{init})$
- 5: Action selection for observed traffic state S_t

$$Q_{\pi^*} = \begin{cases} \operatorname{argmax}_a Q_n & \text{if } \epsilon_{\text{random}} > \epsilon \\ \text{select random actions} & \text{otherwise} \end{cases}$$
- 6: Execute action a_t in s_t , receives reward r_t and ends up in s_{t+1}
- 7: store the experience tuple (s_t, a_t, r_t, s_{t+1}) in buffer memory \mathcal{B}
- 8: derive \mathcal{M} random minibatch tuple samples from memory \mathcal{B}
 Set $target(s, a) = \begin{cases} \text{reward} & \text{for terminal } s_t \\ \text{equation2} & \text{for non-terminal } s_t \end{cases}$
- 9: Update Target network as $\theta^T = \beta\theta^P + (1 - \beta)\theta^T$
 Update Primary network by minimizing loss function MSE
- 11: **end for**
- 12: **end for**

240

241 **6. Numerical case-study**

242 In this study, the Vissim Component Object Model (COM) is used to read, evaluate and change
 243 the simulation objects with python.

244 **6.1. Simulation setting and parameters**

245 The DDQN ATSC model for each data resolution is executed for 1000 episodes. Each episode
 246 is 1800 second which is the two consecutive 15 minutes volume data collected at the
 247 intersection for morning rush hour. The origin-destination demand matrices are as follows:

248

249

Table 3-1: Demand matrix number 1

Date: 08-01-2019 , time: 08:00-08:15				
	N	E	S	W
N	0	53	217	13
E	6	0	46	120
S	37	12	0	15
W	10	49	11	0

250

251 **Table 3-2: Demand matrix number 2**

252 **Date: 08-01-2019 , time: 08:15-08:30**

	N	E	S	W
N	0	67	261	14
E	5	0	54	113
S	40	12	0	16
W	10	69	12	0

253 **6.2. Experimental results**

254 To assess the performance of the deep RL models, data measurement and queue estimators are
 255 added in Vissim simulation model. The experimental results are set out both in tables and plots.
 256 The Table 4 represents the mean value the last 100 episodes with the 95% confidence interval.
 257 To better visualization, the plots are a simple moving average over 50 episodes for the whole
 258 training episodes.

259
 260 **6.2.1. cumulative reward**

261 The reward for ATSC-DTSE is the summation of $r_t = -(1 - 0.2v_t)$ for each time step during
 262 the action execution which consists of vehicles individual speed data. However, the ATSC-
 263 QLDS and ATSC-QLAS reward is the summation of queued vehicles $r_t = -(q_{i,j})$. As the
 264 ATSC-DTSE data resolution is higher than ATSC-QLAS, the rewards are more accurate for
 265 RL agent to take the optimal actions.

266 **6.2.2. Performance measures**

267 According to Table1, the ATSC-DTSC is superior in all performance measures. The ATSC-
 268 DTSE outperforms ATSC-QLAS (with queue length and average speed data) for more than
 269 4%. However, the improvement over the ATSC-QLDS is less than 0.7% which is quite
 270 negligible. Also, the ATSC-DTSE computation time is 30% more than ATSC-QLSD model
 271 training time.

272
 273 **Table 4: Performance metrics results of last 100 training episodes with 95% confidence interval**

model	ATSC-QLAS (M1)	ATSC-QLDS (M2)	ATSC-DTSE (M3)	M3	M3
				Improvement over M1 (%)	Improvement over M2 (%)
Cumulative reward	-22156.50 ± 163.56	-21352.31 ± 165.08	-20691.06 ± 276.89	6.61	3.1
Average travel time (s)	31.60 ± 0.13	30.90 ± 0.15	30.76 ± 0.20	2.66	0.45
Average delay (s)	21.56 ± 0.13	20.84 ± 0.15	20.71 ± 0.20	3.94	0.62
Average queue length (m)	4.93 ± 0.04	4.73 ± 0.05	4.70 ± 0.07	4.67	0.63
Average queue length maximum (m)	35.63 ± 0.45	34.80 ± 0.45	34.90 ± 0.68	2.05	-0.29
Average number of queue stops	72.99 ± 0.51	71.13 ± 0.37	69.46 ± 0.43	4.84	2.35

274
 275
 276
 277

278
279
280

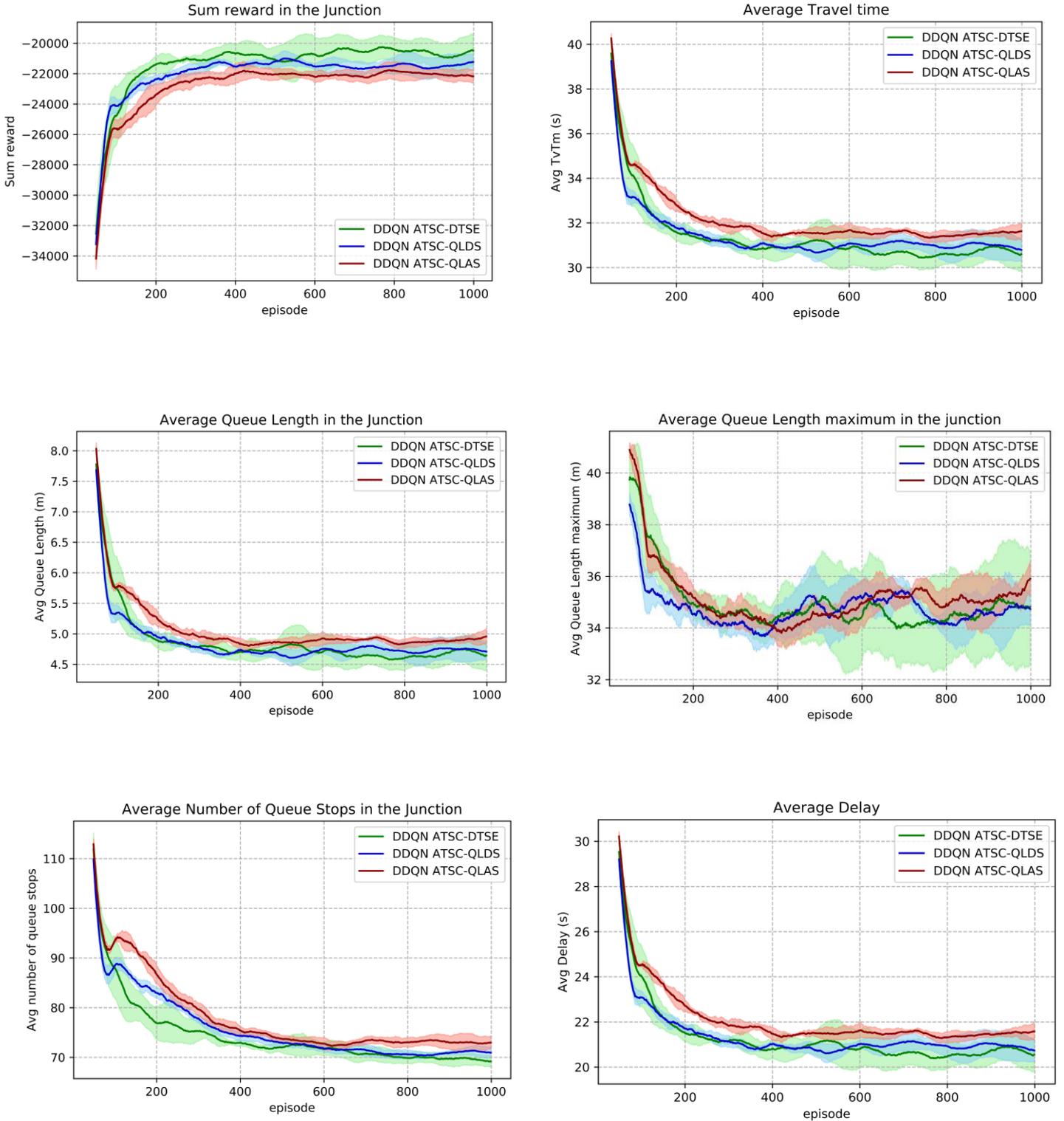


Figure 3: experimental results are the simple moving average over 50 episodes

281

282

283 7. Discussion

284 The simulation experiments have shown that ATSC-QLDS model (with queue length and
 285 density data) and ATSC-DTSE model (with information-dense data) have almost similar
 286 results in terms of travel time and delay reduction, queue length shortening. The ATSC-DTSE
 287 data is only available in a vehicular network which is hard and expensive when it comes to
 288 real-world implementation. The ATSC-QLDS data, in turn, can be collected by cameras at the
 289 intersection. Thus, it is recommended to use queue length and density data for the real-world
 290 applications of deep RL-based ATSC.

291 8. Conclusion

292 In this paper, we evaluate the performance of deep RL-based ATSC model with three different
 293 types of data resolutions. The first data is vehicles queue length of each incoming lane and
 294 vehicles average speed of the incoming links (ATSC-QLAS model). The second data is the
 295 vehicles queue length and density for incoming lanes which is provided by camera at the
 296 intersection (ATSC-QLDS model). The third and highest data resolution is speed and location
 297 of each individual vehicle in a vehicular network (ATSC-DTSE). The comprehensive
 298 experiments on an isolated intersection with DDQN algorithm agent has shown the superiority
 299 of ATSC-DTSE in all performance measures (over 4% versus ATSC-QLAS and less than 0.7%
 300 versus ATSC-QLDS). However, the ATSC-DTSE training time is 30% more than ATSC-
 301 QLDS. Also, the ATSC-QLDS model camera data is more accessible for real-world
 302 implementation while ATSC-DTSE must be collected in a vehicular network with high
 303 communication rates between vehicles. Hence, it is recommended to use the vehicles queue
 304 length and density data instead of information-dense speed and location data of each individual
 305 vehicle.

306 9. Reference

- 307 ABDOOS, M., MOZAYANI, N. & BAZZAN, A. L. Traffic light control in non-stationary
 308 environments based on multi agent Q-learning. 14th International IEEE conference
 309 on intelligent transportation systems (ITSC), 2011. 1580-1585.
- 310 AREL, I., LIU, C., URBANIK, T. & KOHLS, A. 2010. Reinforcement learning-based multi-
 311 agent system for network traffic signal control. *IET Intelligent Transport Systems*, 4,
 312 128-135.
- 313 AUDIT, T. A. I. 2019. Urban Transport Crowding and Congestion.
- 314 BALAJI, P., GERMAN, X. & SRINIVASAN, D. 2010. Urban traffic signal control using
 315 reinforcement learning agents. *IET Intelligent Transport Systems*, 4, 177-188.
- 316 EL-TANTAWY, S. & ABDULHAI, B. An agent-based learning towards decentralized and
 317 coordinated traffic signal control. 13th International IEEE Conference on Intelligent
 318 Transportation Systems, 2010. 665-670.
- 319 EMAMI, A., SARVI, M. & ASADI BAGLOEE, S. 2020. A review of the critical elements
 320 and development of real-world connected vehicle testbeds around the world.
 321 *Transportation Letters*, 1-26.
- 322 GAO, J., SHEN, Y., LIU, J., ITO, M. & SHIRATORI, N. 2017. Adaptive traffic signal
 323 control: Deep reinforcement learning algorithm with experience replay and target
 324 network. *arXiv preprint arXiv:1702.02755*.
- 325 GENDERS, W. 2019. Policy Analysis of Adaptive Traffic Signal Control Using
 326 Reinforcement Learning. *Journal of Computing in Civil Engineering*, 34, 04019046.

- 327 GENDERS, W. & RAZAVI, S. 2016. Using a deep reinforcement learning agent for traffic
328 signal control. *arXiv preprint arXiv:1611.01142*.
- 329 GENDERS, W. & RAZAVI, S. 2019. Asynchronous n-step Q-learning adaptive traffic signal
330 control. *Journal of Intelligent Transportation Systems*, 23, 319-331.
- 331 GRÉGOIRE, P.-L., DESJARDINS, C., LAUMÔNIER, J. & CHAIB-DRAA, B. Urban traffic
332 control based on learning agents. *IEEE Intelligent Transportation Systems*
333 *Conference, 2007*. 916-921.
- 334 JEON, H., LEE, J. & SOHN, K. 2018. Artificial intelligence for traffic signal control based
335 solely on video images. *Journal of Intelligent Transportation Systems*, 22, 433-445.
- 336 LECUN, Y., BENGIO, Y. & HINTON, G. J. N. 2015. Deep learning. 521, 436-444.
- 337 LIANG, X., DU, X., WANG, G. & HAN, Z. 2019. A Deep Reinforcement Learning Network
338 for Traffic Light Cycle Control. *IEEE Transactions on Vehicular Technology*, 68,
339 1243-1253.
- 340 MNIH, V., KAVUKCUOGLU, K., SILVER, D., RUSU, A. A., VENESS, J., BELLEMARE,
341 M. G., GRAVES, A., RIEDMILLER, M., FIDJELAND, A. K. & OSTROVSKI, G.
342 2015. Human-level control through deep reinforcement learning. *Nature*, 518, 529.
- 343 MOUSAVI, S. S., SCHUKAT, M. & HOWLEY, E. 2017. Traffic light control using deep
344 policy-gradient and value-function-based reinforcement learning. *IET Intelligent*
345 *Transport Systems*, 11, 417-423.
- 346 PRASHANTH, L. & BHATNAGAR, S. 2010. Reinforcement learning with function
347 approximation for traffic signal control. *IEEE Transactions on Intelligent*
348 *Transportation Systems*, 12, 412-421.
- 349 SUTTON, R. S. & BARTO, A. G. 2018. *Reinforcement learning: An introduction*.
- 350 TAN, T., BAO, F., DENG, Y., JIN, A., DAI, Q. & WANG, J. 2019. Cooperative deep
351 reinforcement learning for large-scale traffic grid signal control *IEEE transactions on*
352 *cybernetics*, ePub.
- 353 TIELEMAN, T. & HINTON, G. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running
354 average of its recent magnitude. *COURSERA: Neural networks for machine learning*.
- 355 VAN DER POL, E. & OLIEHOEK, F. A. 2016. Coordinated deep reinforcement learners for
356 traffic light control. *Proceedings of Learning, Inference Control of Multi-Agent*
357 *Systems*.
- 358 WANG, T., CAO, J. & HUSSAIN, A. 2021. Adaptive Traffic Signal Control for large-scale
359 scenario with Cooperative Group-based Multi-agent reinforcement learning.
360 *Transportation research part C: emerging technologies*, 125, 103046.
- 361 WEI, H., CHEN, C., WU, K., ZHENG, G., YU, Z., GAYAH, V. & LI, Z. 2019. Deep
362 Reinforcement Learning for Traffic Signal Control along Arterials.
- 363 WEI, H., ZHENG, G., YAO, H. & LI, Z. Intellilight: A reinforcement learning approach for
364 intelligent traffic light control. *Proceedings of the 24th ACM SIGKDD International*
365 *Conference on Knowledge Discovery & Data Mining, 2018*. 2496-2505.
- 366 WIERING, M. Multi-agent reinforcement learning for traffic light control. *Machine*
367 *Learning: Proceedings of the Seventeenth International Conference (ICML'2000)*.
368 1151-1158.
- 369 WU, T., ZHOU, P., LIU, K., YUAN, Y., WANG, X., HUANG, H. & WU, D. O. 2020.
370 Multi-agent deep reinforcement learning for urban traffic light control in vehicular
371 networks. *IEEE Transactions on Vehicular Technology*, 69, 8243-8256.
- 372 ZHENG, G., ZANG, X., XU, N., WEI, H., YU, Z., GAYAH, V., XU, K. & LI, Z. 2019.
373 Diagnosing Reinforcement Learning for Traffic Signal Control. *arXiv preprint*
374 *arXiv:1904.04716*.
- 375