

---

# Reinforcement Learning Benchmarks for Traffic Signal Control

---

James Ault

Texas A&M University  
College Station, TX  
jault@tamu.edu

Guni Sharon

Texas A&M University  
College Station, TX  
guni@tamu.edu

## Abstract

We propose a toolkit for developing and comparing reinforcement learning (RL)-based traffic signal controllers. The toolkit includes implementation of state-of-the-art deep-RL algorithms for signal control along with benchmark control problems that are based on realistic traffic scenarios. Importantly, the toolkit allows a first-of-its-kind comparison between state-of-the-art RL-based signal controllers while providing benchmarks for future comparisons. Consequently, we compare and report the relative performance of current RL algorithms. The experimental results suggest that previous algorithms are not robust to varying sensing assumptions and non-stylized intersection layouts. When more realistic signal layouts and advanced sensing capabilities are considered, a distributed deep Q-learning approach is shown to outperform previously reported state-of-the-art algorithms in many cases.

## 1 Introduction

Travel time studies in urban areas show that 12–55% of commute travel time is due to delays induced by signalized intersections (stopped or approach delay) [14, 24]. Hence, optimized signal control has the potential of reducing commute time, traffic congestion, emissions, and fuel consumption, while requiring minimal infrastructure changes.

A signalized intersection is composed of incoming and outgoing roads where each road is affiliated with one or more lanes. A signal controller must assign right of passage to phases, where each phase corresponds to a specific traffic movement through the intersection (incoming to outgoing roads/lanes). Two phases are defined to be in *conflict* if they cannot be enabled simultaneously (their affiliated traffic movement is intersecting). Each intersection serves vehicles which are assumed to continuously arrive on incoming roads. Each vehicle is associated with a specific, outgoing target road. At each time step, the signal controller is tasked with assigning right-of-passage (green signal) to a set of non-conflicting phases such that some utility measurement is optimized. The utility to be optimized is commonly defined as the sum of vehicles’ delay imposed by the intersection [20, 3].

Given that signalized intersections vary with regards to their layout and demand profile, optimized control policies may differ and are instance dependent. Consequently, signal controllers usually require to be optimized based on the observed state of the environment. Such online optimization of signal controllers requires: (a) sensing the state of approaching traffic (e.g., number and position of approaching vehicles, approaching speeds, queue length, accumulated delay) aggregated by approaching roads/lanes, and (b) defining a control policy that takes the current state of traffic as input and outputs the next phases to be enabled (which is translated to a green, yellow, and red assignment for each light box).

A line of publications have attempted to harness deep reinforcement-learning (RL) techniques towards this control optimization problem. While several approaches claim state-of-the-art performance [7,

36 [18], they fail to compare performance between themselves using a standard testbed. Consequently, it  
37 is challenging to determine which algorithm results in state-of-the-art performance and under what  
38 circumstances. This paper attempts to address this gap by establishing:

- 39 1. An RL testbed environment for traffic signal control that is based on the well-established  
40 Simulation of Urban Mobility traffic simulator (SUMO) [5].
- 41 2. Benchmark single- and multiagent-signal control tasks which are based on realistic traffic  
42 scenarios from SUMO.
- 43 3. An OpenAI GYM interface [6] which allows easy deployment of standard RL algorithms.
- 44 4. A standardized implementation of state-of-the-art RL-based signal control algorithms.

45 The presented testbed and benchmark environment denoted *REinforced Signal COntrol* or RESCO  
46 is used to provide a first-of-its-kind comparison between state-of-the-art RL-based signal control  
47 algorithms. The reported comparative study is performed using the aforementioned, SUMO simulator  
48 and scenarios which are inspired by real-world cities and calibrated demands. The reported results  
49 paint a picture where algorithms that claim state-of-the-art performance struggle in realistic traffic  
50 scenarios and are often outperformed by a decentralized deep Q-learning approach [19].

## 51 2 Background

52 Recent publications [15, 25] proposed to utilize state-of-the-art reinforcement learning (RL) algo-  
53 rithms for online optimization of signal controllers. In this approach, the state of the intersection is  
54 usually defined by the set of incoming vehicles (incoming lane, speed, waiting time, queue position)  
55 and the current signal (right-of-passage) assignment. An RL agent is tasked with optimizing a policy  
56 which maps states to signal assignment. Such an approach showed a potential reduction of up to 73%  
57 in vehicle delays when compared to fixed-time actuation [20].

### 58 2.1 Reinforcement learning

59 In reinforcement learning (RL) an agent is assumed to learn through interactions with the environment.  
60 The environment is commonly modeled as a Markov decision process (MDP) which is defined by:  
61  $\mathcal{S}$  – the state space,  $\mathcal{A}$  – the action space,  $\mathcal{P}(s_t, a, s_{t+1})$  – the transition function of the form  
62  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ ,  $R(s, a)$  – the reward function of the form  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , and  $\gamma$  – the  
63 discount factor. The agent is assumed to follow an internal policy  $\pi$  which maps states to actions, i.e.,  
64  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ . The agent's chosen action ( $a_t$ ) at the current state ( $s_t$ ) affects the environment such that a  
65 new state emerges ( $s_{t+1}$ ) as well as some reward ( $r_t$ ) representing the immediate utility gained from  
66 performing action  $a$  at state  $s$ , given by  $R(s, a)$ . The observed reward is used to tune the policy such  
67 that the expected sum of discounted reward,  $J_\pi = \sum_t \gamma^t r_t$ , is maximized. The policy  $\text{argmax}_\pi[J_\pi]$   
68 is the optimal policy and is denoted by  $\pi^*$ .

69 Common approaches for training a policy using RL include, value-based, policy-gradient, and  
70 actor-critic approaches. A value-based approach attempts to learn the expected future utility from  
71 states (*state value*) or from action-state pairs (*action value* or *q-value*). The control policy is then  
72 directed towards actions/states that maximize the expected utility ( $J_\pi$ ). A prominent example of  
73 a value-based approach is the model-free deep Q-learning algorithm [19]. In the policy-gradient  
74 approach [30] a policy is defined through a parameterized differential equation, where the parameters  
75 are gradually updated, following the policy gradient, towards favorable outcomes (as experienced  
76 through the reward function). Using state or action value estimations for defining favorable outcomes  
77 for policy-gradient updates is usually referred to as an actor-critic approach. A prominent example  
78 of a state-of-the-art actor-critic approach is the proximal-policy optimization (PPO) algorithm [22]  
79 which provides some guarantees regarding monotonicity in the policy improvement over training  
80 iterations.

### 81 2.2 Traffic signal control as an MDP

82 A signalized intersection is composed of incoming and outgoing roads where each road is assembled  
83 from one or more lanes. The intersection is assigned a set of phases,  $\Phi$ . Each phase,  $\varphi \in \Phi$ , is  
84 affiliated with a specific traffic movement through the intersection, as illustrated in Figure 1. Two

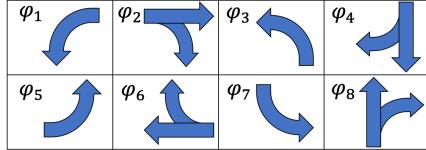


Figure 1: Typical phases, or traffic movements, in a 4-way intersection.

85 phases are defined to be in *conflict* if they cannot be enabled simultaneously (their affiliated traffic  
86 movement is intersecting).

87 For example, in the phase allocation presented in Figure 1,  $\varphi_2$  and  $\varphi_1$  are conflicting phases. At each  
88 time step, a signal controller is responsible to enable some combination of non-conflicting phases  
89 such that a long-term objective function is optimized. When considering RL-based controllers, the  
90 signalized intersection environment is commonly modeled through the following MDP.

91 **State Space ( $S$ ):** the state space is defined by the state of incoming traffic and the currently enabled  
92 phases. Specifically, the state of incoming traffic is defined through the assumed sensing capabilities.  
93 These assumptions vary between publications. Some published work [3] assume state-of-the-art  
94 traffic sensing technology [13] which allows high-resolution data regarding incoming traffic. For  
95 example, real-time observations regarding the number of approaching vehicles, stopped vehicles  
96 accumulated waiting time, number of stopped vehicles, and average speed of approaching vehicles.  
97 Other publications assume less informative sensing capabilities. For instance by assuming that only  
98 the stopped queue length per lane is visible [7], or by assuming that only the waiting time of the first  
99 vehicle in the queue is visible [18].

100 Previous work also vary in the assumed sensing radius. While some assume a sensing radius that  
101 covers the entirety of the incoming roads [7],<sup>1</sup> others assume a more realistic 50 meter [18] or 200  
102 meter [3] sensing radius.

103 **Action Space ( $A$ ):** at each time-step, the controller chooses a set of non-conflicting phases to be  
104 assigned the right-of-passage (green light). If the chosen phases differ from the currently enabled  
105 phases then a mandatory yellow phase is enforced for a predefined time duration. Note that assigning  
106 yellow phases is not a part of the action space but a constraint imposed on the control sequence by  
107 the environment.

108 **Transition function ( $P$ ):** the transition function is defined by the traffic progression following the  
109 signal assignment. This progression can be defined within a simulated environment following a  
110 specific traffic model (as in this paper) or by real-world traffic progression as part of a real-world  
111 implementation (out of scope for this paper).

112 **Reward function ( $R$ ):** previous publications commonly used (minus) queue length summed over  
113 all incoming lanes as their reward function [29]. Such a reward function is simple to implement  
114 and is relevant to congestion alleviation. On the other hand, it fails to normalize the benefits from  
115 optimized signal operation over entire trips. Consequently, other reward functions were suggested.  
116 Among them, (minus) total delays imposed by the intersection [23], (minus) waiting time at the  
117 intersection [3], and (minus) traffic pressure [18].

### 118 2.2.1 Multiagent control

119 In many real-world scenarios traffic flow is sought to be optimized over a road network which includes  
120 multiple intersections. In such scenarios, signal control is assumed to be centrally coordinated over  
121 all the intersections. These control problems are denoted *multiagent signal control* and are defined by  
122 the Cartesian product of the composing intersections' state spaces and action spaces.

123 Following the curse of dimensionality [4] introduced by a multiagent control task, some researchers  
124 take a hierarchical control approach where intersections are partitioned into local control groups [18,  
125 7].

<sup>1</sup>Chen et al. 2020 does not explicitly specify the sensing radius. However, the referenced codebase uses an unbounded sensing radius.

126 **2.3 Related work**

127 Next, we review several RL algorithms that claim state-of-the-art performance for traffic signal  
128 control.

129 **Deep Q-Learning:** deep Q-learning [19] has been proposed for the control of signalized intersections  
130 in a number of works [3, 25, 23, 16, 20, 12, 27, 15]. Here we adopt the implementation of [3], which  
131 reported up to a 19.4% improvement over an actuated controller. Ault et al. defined a convolutional  
132 Q-network where queue length, number of approaching vehicles, total approaching speed, and total  
133 waiting time are aggregated in convolutional layers over lanes composing the same incoming road.  
134 This implementation uses a (minus) total waiting time reward function.

135 **MPLight:** MPLight [7] utilizes the concept of pressure to coordinate multiple intersections. Pressure  
136 is the difference in queue lengths from incoming lanes of an intersection and the queue length on a  
137 downstream intersection's receiving lane. Chen et al. used pressure as both the state and reward for a  
138 DQN agent shared over all intersections on top of the FRAP [32] model. The authors reported up to a  
139 19.2% improvement in travel times over the next best compared method, PressLight [28].

140 **FMA2C:** FMA2C [8] builds on the prior work of MA2C. MA2C enabled cooperation between  
141 signal control agents (one per intersection) denoted workers. Adjacent workers are coordinated  
142 through a local discounted neighborhood reward, discounted neighborhood-appended states, and  
143 action fingerprint sharing on otherwise independent advantage actor-critic agents. FMA2C extends  
144 this to a hierarchy of managing agents on top of the workers. The managing agents are trained to  
145 optimize flow within their assigned region. The workers are then trained to incorporate the high-level  
146 goals of their managing agent. The authors reported up to a 6% improvement in average delays over  
147 the Greedy controller described in Section 3.4.

148 **2.3.1 Evaluation environments for RL-based signal controllers**

149 For their experimental evaluation, previous publications relied on custom made scenarios that were, in  
150 many cases, tailored for the evaluated RL algorithm. Jinming and Feng 2020 used the well-established  
151 simulation of urban mobility (SUMO) environment. SUMO is widely accepted in the transportation  
152 community and—as we also note—is a reasonable testbed choice. Jinming and Feng did provide  
153 results on a scenario that is based on a real-world city (Monaco). However, the reported scenario  
154 was based on a modified Monaco scenario which contains an addition of 18 synthetic traffic signals  
155 beyond the official “MoST” scenario [9] and also includes non-validated, inflated traffic demands.  
156 We extend on that work by introducing benchmark signal control tasks that are based on validated  
157 traffic scenarios while providing a flexible interface that is suitable for a variety of RL algorithms.  
158 Moreover, our codebase is publicly available allowing researchers to perform meaningful comparative  
159 evaluations.

160 Zhang et al. 2019 presented their own simulation testbed denoted *CityFlow*. This evaluation testbed  
161 suffers from two main drawbacks, (1) as oppose to SUMO, CityFlow is not rigorously calibrated and  
162 evaluated by the general transportation community. CityFlow is claimed to produce equivalent output  
163 as SUMO. However, those claims are based on results from simplified grid network scenarios; (2) a  
164 common benchmark scenario in CityFlow is the Manhattan, NY network. This scenario is claimed to  
165 represent real-world city layout and demand. However, support for this claim is limited.

166 Other relevant publications [21] presented evaluations that are based on the autonomous intersection  
167 management (AIM) simulator [10]. The main drawback of the AIM simulator is the lack of traffic  
168 scenarios which are based on real-world cities. AIM commonly produces a simple grid network with  
169 symmetric intersections. Again, one might claim that such a grid network is akin to the road layout in  
170 Manhattan, NY, yet a deeper analysis of traffic trends is required to support such claims and their  
171 relevancy to the real world.

172 **3 The Reinforced Signal Control (RESCO) toolkit**

173 We present a standard RL traffic signal control testbed denoted reinforced signal control or RESCO.  
174 The main goals of this standard testbed are:

- 175 1. Provide benchmark single- and multiagent-signal control tasks which are based on well-  
176 established traffic scenarios.

- 177        2. An OpenAI GYM interface [6] within the testbed environment which allows easy deploy-  
 178        ment of state-of-the-art RL algorithms.  
 179        3. A standardized implementation of state-of-the-art RL-based signal control algorithms.

180        RESCO is open source and free to use/modify under the creative commons BY-NC-SA license. The  
 181        code is built on top of SUMO-RL [1] and is available on Github at [github.com/jault/RESCO](https://github.com/jault/RESCO). The  
 182        embedded traffic scenarios are distributed with their own licenses. Cologne based scenarios are under  
 183        creative commons BY-NC-SA and Ingolstadt based, with the GNU General Public License 3.

### 184        3.1 State and action Space

185        In order to allow a variety of sensing assumptions (as described in Section 2.2, “State space”),  
 186        RESCO assumes the most advanced sensing capabilities [3]. The user can pull any subset of the state  
 187        features based on specific sensing assumptions. Specifically, per state, per intersection, per lane, the  
 188        system provides a dictionary with the following features: stopped vehicles queue length, number  
 189        of approaching vehicles (not stopped), total waiting time for stopped vehicles, sum of approaching  
 190        vehicles<sup>1</sup> (not stopped) speed, maximum waiting time (over stopped vehicles), number of arrivals  
 191        (added vehicles during last time step), number of departures (vehicles removed during last time step).  
 192        On top of that, the user can also define the effective sensing distance on initialization.

193        For each intersection the provided benchmark control task defines the set of non-conflicting phases.  
 194        The action space is the set of non-conflicting phase combinations as defined in Section 2.2, “Action  
 195        space”. By default acts are chosen for the next 10 seconds of simulation, with the first 3 seconds  
 196        reserved for yellow signals, if necessary, following Ma and Wu [18]. The interface accepts a unique  
 197        intersection ID and an index representing the chosen phase combination to be enabled.<sup>2</sup>

### 198        3.2 Reward metrics

199        In order to allow maximal flexibility for the user, the interface allows designating any of the reward  
 200        metrics defined in Section 2.2 “Reward function”, as well as a custom weighted combination of these  
 201        metrics. When initializing a control task, the user can pass a weight vector as an argument where  
 202        each entry designates the weight of one of the metrics in the reward function. The weight vector  
 203        entries are defined as follows. 1: system travel time, 2: approximated signal induced delays,<sup>3</sup> 3: total  
 204        waiting time at intersections, 4: average queue length, 5: traffic pressure.

### 205        3.3 Benchmark control tasks

206        The presented signal control benchmark tasks are based off two well-established SUMO scenarios,  
 207        namely, “TAPAS Cologne” [26] and “InTAS” [17]. Both scenarios describe traffic within a real-world  
 208        city, Cologne and Ingolstadt (Germany) respectively, including road network layout and calibrated  
 209        demands. The road network for these scenarios is shown in Figure 2. These scenarios were chosen as  
 210        they are well-accepted by the transportation community and include a congested downtown zone with  
 211        multiple signalized intersections. Note that RESCO can easily be fitted to other SUMO scenarios.

212        Three benchmark control tasks are considered per traffic scenario, namely, (a) controlling a single  
 213        main intersection, (b) coordinated control of multiple intersections along an arterial corridor, and  
 214        (c) coordinated control of multiple intersections within a congested area (downtown). The affiliated  
 215        intersections are depicted in Figure 2.

### 216        3.4 Benchmark algorithms

217        RESCO defines three baseline controllers and several RL-based controllers.

#### 218        Baseline controllers:

<sup>2</sup>Control tasks in RESCO provide the set of valid (non-conflicting) phase assignments per intersection, each affiliated with a unique index. Note that this set might differ between intersections based on their layout, unique features, and safety considerations.

<sup>3</sup>Signal induced delays can only be computed in hindsight. RESCO follows previous work [2] which suggested real-time approximation as the difference between the vehicle’s current speed and the maximum speed limit over all vehicles.

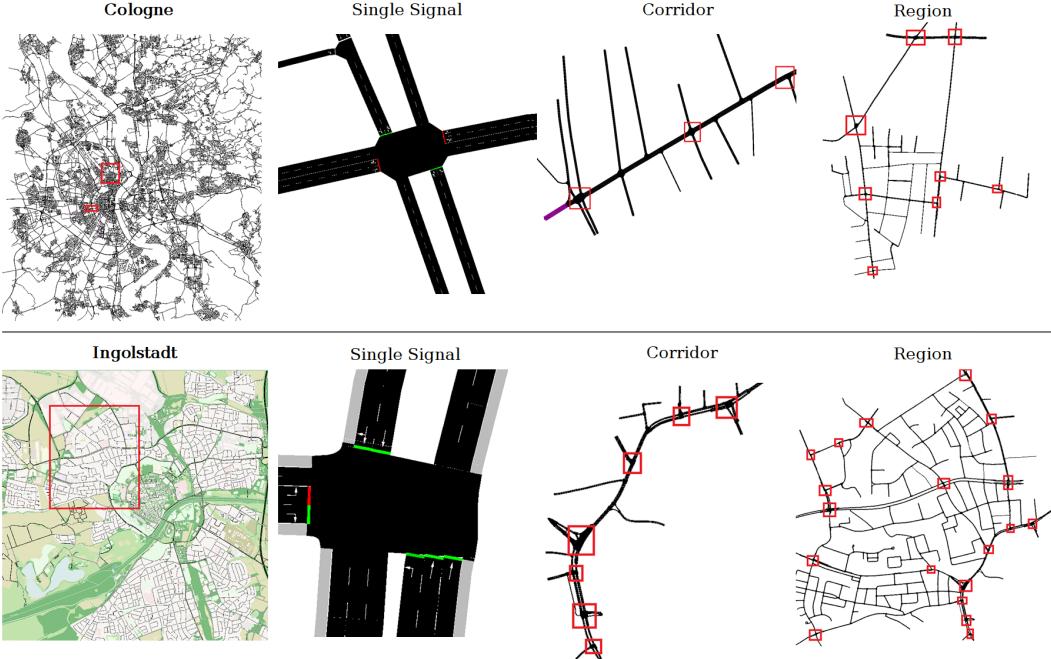


Figure 2: Road networks extracted from the full-city SUMO scenarios for benchmarking. In the full-city maps the areas bounded in red are the locations of the extracted networks. In the corridor and region networks the bounded areas mark signalized intersections. All other intersections in the extracted areas are controlled by traffic priority laws. Traffic demands for the extracted areas were taken from routes in the full-city scenarios using the standard SUMO tool for creating sub-scenarios.

219 (1) **Fixed-time** (or Pre-timed) control where each phase combination is enabled for a fixed duration  
 220 following a fixed cycle. The intervals are defined per intersection as part of the SUMO traffic scenario;  
 221 (2) **Max-pressure** control where the phase combination with the maximal joint pressure is enabled  
 222 as described in Chen et al. [7];  
 223 (3) **Greedy** control where the phase combination with the maximal joint queue length and approaching  
 224 vehicle count is enabled as described in Ma and Wu [18].

225 **RL controllers:**

- 226 (1) **IDQN** – independent DON agents, one per intersection, each with convolution-layers for lane  
 227 aggregation as described by Ault et Al. 2020. Hyper-parameters are left as the default values in the  
 228 Preferred RL library with the exception of the target network update interval, which was adjusted to  
 229 the Atari environment settings of 500 steps per update;
- 230 (2) **IPPO** – the same deep neural network is used as IDQN, coming from [3]. Hyper-parameters were  
 231 set following the Preferred RL defaults for the Atari environments;
- 232 (3) **MPLight** – implementation is based on the FRAP open source implementation [32] along with  
 233 the ChainerRL DQN implementation and pressure sensing. Hyper-parameters are identical to that of  
 234 IDQN;
- 235 (4) **Extended MPLight** – Denoted MPLight\*, similar to the MPLight implementation with the  
 236 addition of sensing information matching IDQN appended to the existing pressure state;
- 237 (5) **FMA2C** – built on top of MA2C open source implementation [8]. Hyperparameters were set  
 238 according to the open source implementation by Chu et al..

239 In each of IDQN, IPPO, and MPLight the learning algorithm implementation is called directly from  
 240 the ChainerRL [11] python library successor, Preferred RL.

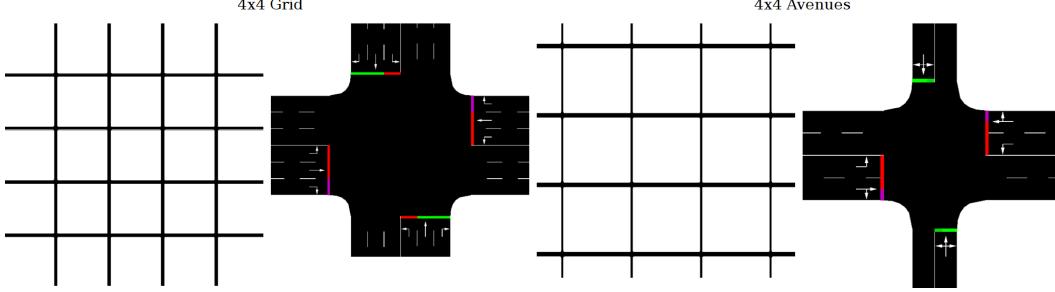


Figure 3: A  $4 \times 4$  Grid validation scenario based on that presented by Chen et al. 2020 (left), and a  $4 \times 4$  Avenue scenario based on the scenario presented by Ma and Wu 2020 (right).

## 241 4 Experiments

242 The experimental section is divided into three. First, we present results that validate our benchmark  
 243 algorithms implementation against performance trends reported in previous publications. Second,  
 244 we provide a comparative study between state-of-the-art approaches on the realistic traffic scenarios.  
 245 Finally, we draw general conclusions that characterize the algorithms and their performance.

### 246 4.1 Validation

247 For validating the benchmark algorithms implementation, we compare the learning curves and final  
 248 performance of the RL controllers from Section 3.4 against the baseline controllers. The traffic  
 249 scenarios are chosen to be similar to those presented in previous publications including the sensing  
 250 assumptions which might change between the algorithms. Consequently, the affiliated results cannot  
 251 be used for a comparative study between the RL controllers.

252 For validating the MPLight implementation, a traffic scenario was configured to be similar to the  
 253 synthetic  $4 \times 4$  symmetric network that was presented in the original MPLight publication [7].  
 254 The RESCO network is depicted in Figure 3,  $4 \times 4$  Grid, and can be compared with the original  
 255 network (Figure 5 in Chen et al. [7]). Demand was set according to Config. 4 in Chen et al. [7].  
 256 Figure 4 presents the learning curves for MPLight and IDQN. It also includes the baseline controllers  
 257 performance (dotted lines) and the final performance for FMA2C and IPPO. FMA2C and IPPO  
 258 require much more training episodes to converge (about 1,400 as appose to 100 by IDQN and  
 259 MPLight). As a result, only the final performance is included for them. The full training curves for  
 260 FMA2C and IPPO are available in the appendix. Chen et al. 2020 reported a 13% improvement for  
 261 MPLight over Max-pressure. The RESCO results show a similar trend with an 11% improvement.

262 For validating the FMA2C implementation, a traffic scenario was configured to be similar to the  
 263 synthetic  $4 \times 4$  traffic grid presented in the original FMA2C publication [18]. This grid is formed  
 264 by two-lane arterial streets and one-lane avenues. The RESCO induced intersections is depicted in  
 265 Figure 3 and can be compared with the original network (Figure 3(a) in Ma and Wu [18]). Demand  
 266 was set according to the original scenario definition. Figure 4 presents the learning curve for the same  
 267 set of controllers. Again, The full training curves for FMA2C and IPPO are available in the appendix.

268 On  $4 \times 4$  Avenues validation scenario Ma and Wu 2020 reported a 4% improvement for FMA2C  
 269 over Greedy and 19% improvement over IDQN. The RESCO results show a more pronounced trend  
 270 for improvement over the Greedy control with a 20% improvement and a 44% improvement over  
 271 IDQN. In general, both the MPLight and FMA2C implementations in RESCO present trends that  
 272 follow the originally published results for each.

### 273 4.2 Comparative study

274 For the full comparative study all of the benchmark RL controllers were compared on each of the  
 275 benchmark control tasks. Hyper-parameters are constant throughout and not tuned to each scenario.  
 276 Following Ault et al., All sensing capabilities are assumed to be available within a 200 meter radius.  
 277 The state and reward choices made by the algorithms are the primary ways in which coordination is  
 278 proposed. Significantly adjusting the state or reward would therefore no longer be representative of

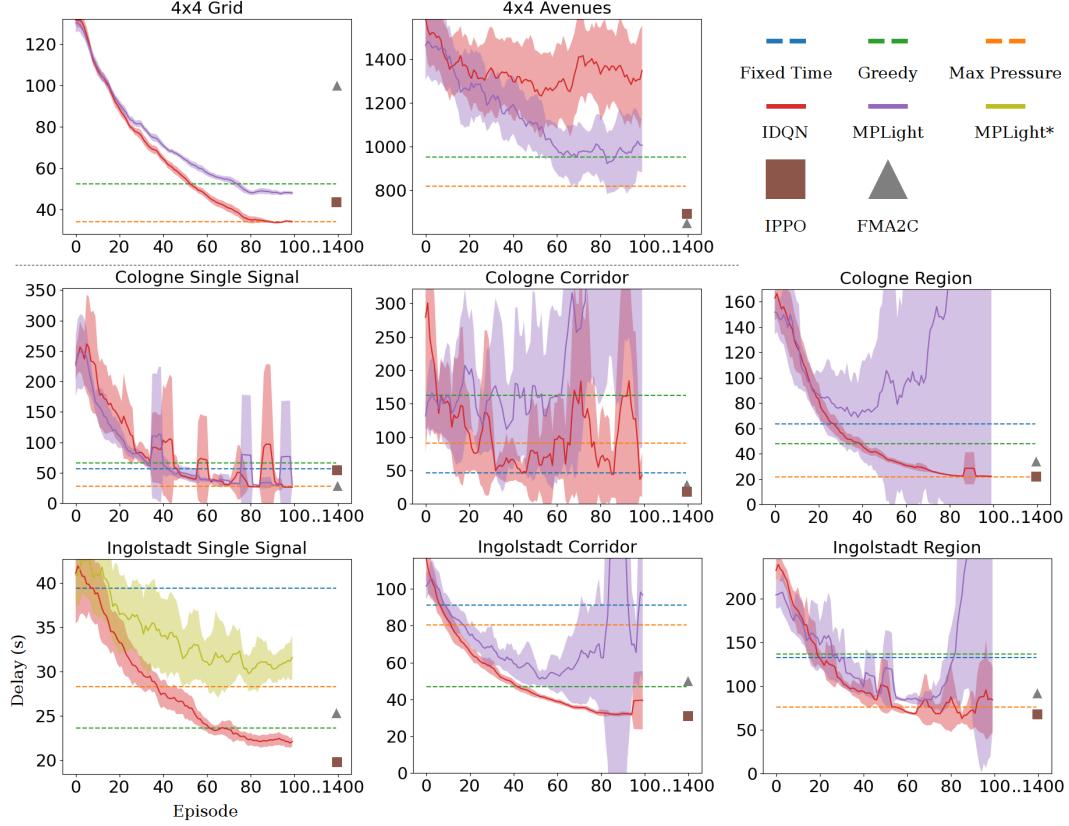


Figure 4: Learning curves over 5 random seeds, with a sliding window average of 5 episodes. Error margins display one-standard deviation from the mean delay in each episode. Baseline algorithms are marked with dashed lines, while solid lines indicate RL algorithms. Square and triangle mark the best performance of their associated algorithms, which require many times the number of episodes of the others. MPLight\* denotes the extended state version of MPLight.

279 the algorithms. Therefore the state and reward functions were set according to the definitions of each  
 280 algorithm (see Section 2.3). Expanding the state to include more measurements generally decreased  
 281 performance on coordination tasks. Future work should fully explore how best to augment relevant  
 282 algorithms when assuming increased sensing capabilities.

283 Figure 4 presents the learning curves for each algorithm in every scenario while Table 1 reports the  
 284 best performing episode (not necessarily the final performance) of each algorithm averaged over five  
 285 random seeds. For example, consider the “Ingolstadt Region” task, Figure 4 paints a picture where  
 286 MPLight diverges during training and eventually results in an average delay that is  $> 200$  seconds.  
 287 Table 1, by contrast, presents a delay of 78 seconds for the same task. The discrepancy is because  
 288 Table 1 corresponds to the best episode, i.e., training episode 64 (before the training divergence).  
 289 Table 1 also includes commonly reported metrics of delay, trip time (duration), waiting time, and  
 290 queue length. However, we observe that in all but one case any one metric is sufficient to indicate  
 291 superior performance in the others. Namely, in the “Ingolstadt Corridor” scenario IPPO reports  
 292 improved performance over IDQN in delay, trip time, and waiting time while IDQN improves in  
 293 queue length, but not to a statistically significant degree. As a result, learning curves are presented  
 294 only for the delay metric (in Figure 4).

295 Both Table 1 and Figure 4 report either extended state MPLight or the standard MPLight, the better  
 296 performing of the two. Extended state MPLight is denoted with an asterisk. In most cases the added  
 297 sensing information in extended state MPLight is not beneficial, however in the “Ingolstadt Single  
 298 Signal” scenario the additional state information is beneficial and allows for convergence where  
 299 standard MPLight failed to converge. MPLight, in our experiments, works well in scenarios with  
 300 similarly structured intersections. However, in scenarios with varying irregular intersections we

Table 1: Performance on benchmark scenarios

<b>IDQN</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	21.48	31.19	<b>59.64</b>	<b>26.05</b>	23.99	22.06
Avg. Trip Time	35.29	68.69	<b>197.23</b>	<b>43.59</b>	59.0	86.02
Avg. Wait	3.93	8.71	<b>20.19</b>	<b>7.98</b>	8.5	5.46
Avg. Queue	0.43	<b>0.67</b>	<b>0.8</b>	<b>2.09</b>	0.87	0.38
<b>IPPO</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	<b>19.85</b>	<b>30.7</b>	67.65	55.07	<b>22.13</b>	<b>21.49</b>
Avg. Trip Time	<b>34.19</b>	<b>68.34</b>	205.44	67.7	<b>57.45</b>	<b>85.54</b>
Avg. Wait	<b>3.21</b>	<b>8.2</b>	26.45	26.15	<b>7.37</b>	<b>5.01</b>
Avg. Queue	<b>0.39</b>	0.71	1.15	8.88	<b>0.76</b>	<b>0.35</b>
<b>MPLight</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	*28.31	48.21	78.16	28.74	83.65	60.42
Avg. Trip Time	*41.07	76.58	215.72	45.85	102.3	123.93
Avg. Wait	*8.27	15.05	34.57	8.61	46.25	30.34
Avg. Queue	*0.61	1.34	1.48	2.45	5.4	2.33
<b>FMA2C</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	25.36	48.99	90.42	30.12	25.37	33.28
Avg. Trip Time	39.4	85.03	226.5	47.31	61.68	97.53
Avg. Wait	7.27	21.9	44.16	11.23	11.3	14.19
Avg. Queue	0.94	1.79	1.74	3.11	1.68	0.98

301 observed divergence in learning. We suspect that this phenomena is due to shared control parameters  
 302 between significantly different control tasks (one controller per intersection).

303 IDQN achieves the best performance in the “Ingolstadt Region” scenario and “Cologne Single Signal”  
 304 scenario, while IPPO achieves the best performance in all other scenarios. However, this is misleading  
 305 as IPPO demonstrates significant instability. Failing to converge in the “Ingolstadt Region” and  
 306 “Cologne Corridor” scenarios, and converging to significantly worse solutions in the “Cologne Single  
 307 Signal”. When examining the performance of each algorithm in the final 10 episodes of training,  
 308 IDQN outperforms all other algorithms in all tasks except the “Cologne Regional”, which IPPO  
 309 still performs well in. Both FMA2C and IPPO have drastically worse sample efficiency than the  
 310 DQN-based methods of MPLight and IDQN. MPLight reaches its best performing episode in 30-80%  
 311 of the time IDQN does, but the performance reached is generally worse. The full training results are  
 312 provided in the appendix.

### 313 4.3 Conclusions

314 Synthetic scenarios, even if made to be challenging coordination tasks, relate poorly to more realistic  
 315 scenarios from the perspective of gauging the performance of reinforcement learning algorithms.  
 316 Irregular and sparsely distributed signals require more general solutions than previously proposed.  
 317 Deep RL methods should give increased attention to hyper-parameter sensitivity as failing to do so  
 318 may result in unstable or poor performance when deployed to previously unseen scenarios.

319 Our experiments suggest that decentralized control algorithms are more robust compared to coordi-  
 320 nated control algorithms when considering realistic traffic scenarios. We speculate that this is, in-part,  
 321 due to the efficient data aggregation of the independent learners utilizing convolutional layers as  
 322 described by Ault et al. 2020. Future work should examine merging this efficient data representation  
 323 with state-of-the-art coordinated control approaches.

324 Coordinated approaches (FMA2C, MPLight) work well under the specific sensing assumptions that  
 325 they made, but have limited applicability when advanced sensing capabilities are considered. The  
 326 decentralized algorithms, IDQN and IPPO, do not share this limitation as they can effectively learn  
 327 instance dependent features.

328 The time elapsed before an algorithm reaches its best performance is important if the ultimate goal is  
329 real-world implementation. To this end, the MPLight algorithm of Chen et al. 2020 has a considerable  
330 advantage, presenting a trade-off of reliability and converged performance for learning speed. Ma and  
331 Wu 2020’s FMA2C may have an advantage when presented with challenging synthetic coordination  
332 problems, but requires a large amount of time to achieve the feat and it is questionable if scenarios  
333 fashioned from real networks and traffic demands would resemble the studied synthetic scenarios.

## 334 **References**

- 335 [1] L. N. Alegre. Sumo-rl. <https://github.com/LucasAlegre/sumo-rl>, 2019.
- 336 [2] I. Arel, C. Liu, T. Urbanik, and A. G. Kohls. Reinforcement learning-based multi-agent system  
337 for network traffic signal control. *IET Intelligent Transport Systems*, 4(2):128–135, 2010.
- 338 [3] J. Ault, J. Hanna, and G. Sharon. Learning an interpretable traffic signal control policy. In  
339 *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent  
340 Systems (AAMAS 2020)*. International Foundation for Autonomous Agents and Multiagent  
341 Systems, May 2020.
- 342 [4] A. G. Barto and S. Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete  
343 event dynamic systems*, 13(1):41–77, 2003.
- 344 [5] M. Behrisch, L. Bieker, J. Erdmann, and D. Krajzewicz. Sumo—simulation of urban mobility:  
345 an overview. In *Proceedings of SIMUL 2011, The Third International Conference on Advances  
346 in System Simulation*. ThinkMind, 2011.
- 347 [6] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.  
348 Openai GYM. *arXiv preprint arXiv:1606.01540*, 2016.
- 349 [7] C. Chen, H. Wei, N. Xu, G. Zheng, M. Yang, Y. Xiong, K. Xu, and Z. Li. Toward a thousand  
350 lights: Decentralized deep reinforcement learning for large-scale traffic signal control. In  
351 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3414–3421,  
352 2020.
- 353 [8] T. Chu, J. Wang, L. Codecà, and Z. Li. Multi-agent deep reinforcement learning for large-  
354 scale traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 21(3):  
355 1086–1095, 2019.
- 356 [9] L. Codeca and J. Härrí. Monaco SUMO Traffic (MoST) Scenario: A 3D Mobility Scenario  
357 for Cooperative ITS. In *SUMO 2018, SUMO User Conference, Simulating Autonomous and  
358 Intermodal Transport Systems, May 14-16, 2018, Berlin, Germany*, Berlin, GERMANY, 05  
359 2018.
- 360 [10] K. Dresner and P. Stone. A multiagent approach to autonomous intersection management.  
361 *Journal of artificial intelligence research*, 31:591–656, 2008.
- 362 [11] Y. Fujita, P. Nagarajan, T. Kataoka, and T. Ishikawa. Chainerrl: A deep reinforcement learning  
363 library. *Journal of Machine Learning Research*, 22(77):1–14, 2021. URL <http://jmlr.org/papers/v22/20-376.html>.
- 365 [12] W. Genders and S. Razavi. Using a deep reinforcement learning agent for traffic signal control.  
366 *arXiv preprint arXiv:1611.01142*, 2016.
- 367 [13] L. A. Klein. *ITS sensors and architectures for traffic management and connected vehicles*. CRC  
368 Press, 2017.
- 369 [14] D. M. Levinson. Speed and delay on signalized arterials. *Journal of Transportation Engineering*,  
370 124(3):258–263, 1998.
- 371 [15] L. Li, Y. Lv, and F.-Y. Wang. Traffic signal timing via deep reinforcement learning. *IEEE/CAA  
372 Journal of Automatica Sinica*, 3(3):247–254, 2016.
- 373 [16] X. Liang, X. Du, G. Wang, and Z. Han. Deep reinforcement learning for traffic light control in  
374 vehicular networks. *arXiv preprint arXiv:1803.11115*, 2018.

- 375 [17] S. C. Lobo, S. Neumeier, E. M. Fernandez, and C. Facchi. InTAS—the ingolstadt traffic scenario  
 376 for SUMO. *arXiv preprint arXiv:2011.11995*, 2020.
- 377 [18] J. Ma and F. Wu. Feudal multi-agent deep reinforcement learning for traffic signal control.  
 378 In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent  
 379 Systems (AAMAS)*, pages 816–824, 2020.
- 380 [19] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,  
 381 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-  
 382 forcement learning. *nature*, 518(7540):529–533, 2015.
- 383 [20] S. S. Mousavi, M. Schukat, and E. Howley. Traffic light control using deep policy-gradient and  
 384 value-function-based reinforcement learning. *IET Intelligent Transport Systems*, 11(7):417–423,  
 385 2017.
- 386 [21] T. T. Pham, T. Brys, M. E. Taylor, T. Brys, M. M. Drugan, P. Bosman, M.-D. Cock, C. Lazar,  
 387 L. Demarchi, D. Steenhoff, et al. Learning coordinated traffic light control. In *Proceedings of  
 388 the Adaptive and Learning Agents workshop (at AAMAS-13)*, volume 10, pages 1196–1201.  
 389 IEEE, 2013.
- 390 [22] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
 391 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 392 [23] S. M. A. Shabestary and B. Abdulhai. Deep learning vs. discrete reinforcement learning for adap-  
 393 tive traffic signal control. In *2018 21st International Conference on Intelligent Transportation  
 394 Systems (ITSC)*, pages 286–293. IEEE, 2018.
- 395 [24] A. Tirachini. Estimation of travel time and the benefits of upgrading the fare payment technology  
 396 in urban bus services. *Transportation Research Part C: Emerging Technologies*, 30:239–256,  
 397 2013.
- 398 [25] E. Van der Pol and F. A. Oliehoek. Coordinated deep reinforcement learners for traffic light  
 399 control. *Proceedings of Learning, Inference and Control of Multi-Agent Systems (at NIPS 2016)*,  
 400 2016.
- 401 [26] C. Varschen and P. Wagner. Mikroskopische modellierung der personenverkehrsnachfrage auf  
 402 basis von zeitverwendungstagebüchern. *Integrierte Mikro-Simulation von Raum- und Verkehrsne-  
 403 wicklung. Theorie, Konzepte, Modelle, Praxis*, 81:63–69, 2006.
- 404 [27] H. Wei, G. Zheng, H. Yao, and Z. Li. Intellilight: A reinforcement learning approach for intelli-  
 405 gent traffic light control. In *Proceedings of the 24th ACM SIGKDD International Conference  
 406 on Knowledge Discovery & Data Mining*, pages 2496–2505. ACM, 2018.
- 407 [28] H. Wei, C. Chen, G. Zheng, K. Wu, V. Gayah, K. Xu, and Z. Li. Presslight: Learning max  
 408 pressure control to coordinate traffic signals in arterial network. In *Proceedings of the 25th ACM  
 409 SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1290–1298,  
 410 2019.
- 411 [29] M. A. Wiering. Multi-agent reinforcement learning for traffic light control. In *Machine Learning:  
 412 Proceedings of the Seventeenth International Conference (ICML'2000)*, pages 1151–1158, 2000.
- 413 [30] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforce-  
 414 ment learning. *Machine learning*, 8(3-4):229–256, 1992.
- 415 [31] H. Zhang, S. Feng, C. Liu, Y. Ding, Y. Zhu, Z. Zhou, W. Zhang, Y. Yu, H. Jin, and Z. Li.  
 416 Cityflow: A multi-agent reinforcement learning environment for large scale city traffic scenario.  
 417 In *The World Wide Web Conference*, pages 3620–3624, 2019.
- 418 [32] G. Zheng, Y. Xiong, X. Zang, J. Feng, H. Wei, H. Zhang, Y. Li, K. Xu, and Z. Li. Learning  
 419 phase competition for traffic signal control. In *Proceedings of the 28th ACM International  
 420 Conference on Information and Knowledge Management*, pages 1963–1972, 2019.

421    **Checklist**

- 422    1. For all authors...
- 423       (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
424       contributions and scope? **[Yes]**
- 425       (b) Did you describe the limitations of your work? **[Yes]**
- 426       (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 427       (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
428       them? **[Yes]**
- 429    2. If you are including theoretical results...
- 430       (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 431       (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 432    3. If you ran experiments (e.g. for benchmarks)...
- 433       (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
434       mental results (either in the supplemental material or as a URL)? **[Yes]**
- 435       (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
436       were chosen)? **[Yes]** See Section 3.4
- 437       (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
438       ments multiple times)? **[Yes]** See Figure 4
- 439       (d) Did you include the total amount of compute and the type of resources used (e.g., type  
440       of GPUs, internal cluster, or cloud provider)? **[No]**
- 441    4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 442       (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 443       (b) Did you mention the license of the assets? **[Yes]**
- 444       (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
- 445       (d) Did you discuss whether and how consent was obtained from people whose data you're  
446       using/curating? **[No]** Attribution and licensing is all that was required from other's data
- 447       (e) Did you discuss whether the data you are using/curating contains personally identifiable  
448       information or offensive content? **[N/A]**
- 449    5. If you used crowdsourcing or conducted research with human subjects...
- 450       (a) Did you include the full text of instructions given to participants and screenshots, if  
451       applicable? **[N/A]**
- 452       (b) Did you describe any potential participant risks, with links to Institutional Review  
453       Board (IRB) approvals, if applicable? **[N/A]**
- 454       (c) Did you include the estimated hourly wage paid to participants and the total amount  
455       spent on participant compensation? **[N/A]**

Table 2: Final Performance on benchmark scenarios

<b>IDQN</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	<b>21.48</b>	<b>31.19</b>	<b>61.11</b>	<b>26.05</b>	<b>24.55</b>	22.06
Avg. Trip Time	<b>35.29</b>	<b>68.69</b>	<b>198.72</b>	<b>43.59</b>	<b>59.0</b>	86.02
Avg. Wait	<b>3.93</b>	<b>8.71</b>	<b>21.36</b>	<b>7.98</b>	<b>8.5</b>	5.46
Avg. Queue	<b>0.43</b>	<b>0.67</b>	<b>0.82</b>	<b>2.09</b>	<b>0.87</b>	0.38
<b>IPPO</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	21.85	32.0	369.34	209.49	58.98	<b>21.62</b>
Avg. Trip Time	35.68	69.38	509.28	224.56	93.7	<b>85.71</b>
Avg. Wait	4.2	8.86	315.31	193.66	43.5	<b>5.19</b>
Avg. Queue	0.49	0.74	5.45	21.86	5.73	<b>0.36</b>
<b>MPLight</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	*29.93	50.81	275.21	29.0	579.69	164.91
Avg. Trip Time	*42.82	79.95	406.27	46.17	595.36	225.14
Avg. Wait	*9.78	16.0	230.88	8.61	548.42	142.35
Avg. Queue	*0.65	1.35	4.69	2.49	23.77	5.31
<b>FMA2C</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	27.74	51.78	96.67	30.82	32.58	34.44
Avg. Trip Time	41.4	87.75	233.04	48.07	66.83	98.68
Avg. Wait	8.74	23.32	48.88	11.87	15.84	14.7
Avg. Queue	1.07	1.86	1.84	3.2	2.01	1.03

## 456 A Appendix

### 457 A.1 License

458 RESCO is open source and free to use/modify under the creative commons BY-NC-SA license. The  
 459 code is available on Github at [<https://github.com/jault/RESCO>]. Included benchmark scenarios are  
 460 distributed with their licenses. Cologne based scenarios are under creative commons BY-NC-SA and  
 461 Ingolstadt based, with the GNU General Public License 3. All experiment can be reproduced from  
 462 the source code, which includes all hyper-parameters and configuration.

### 463 A.2 Supplementary results

464 Figure 5 displays the full learning curves of the IPPO and FMA2C algorithms.

465 Table 2 shows the final, post-training performance for each RL algorithm on the benchmark control  
 466 tasks. These results show the learning divergence for MPLight in the Cologne Corridor, Cologne  
 467 Region, Ingolstadt Corridor, and Ingolstadt Region scenarios as compared to the best performance.

468 Table 3 gives the episode number in which the best performance from Table1 were reached. It  
 469 supports the conclusion that MPLight, followed by IDQN, have the best sample efficiency of the  
 470 compared algorithms.

471 To produce Table 4, the standard deviation of each episode over the 5 random seeds is calculated for  
 472 each metric. The average over all episodes is then displayed in Table 4 per algorithm, allowing for a  
 473 notion of stability. FMA2C displays better stability than any of IDQN, IPPO, and MPLight.

474 Table 5 extends Table 1 by including standard deviations. Table 1 reports the best performing episode  
 475 of each algorithm averaged over five random seeds.

476 Table 6 provides the parameters of each algorithm, note that IDQN and MPLight share the same  
 477 parameters.

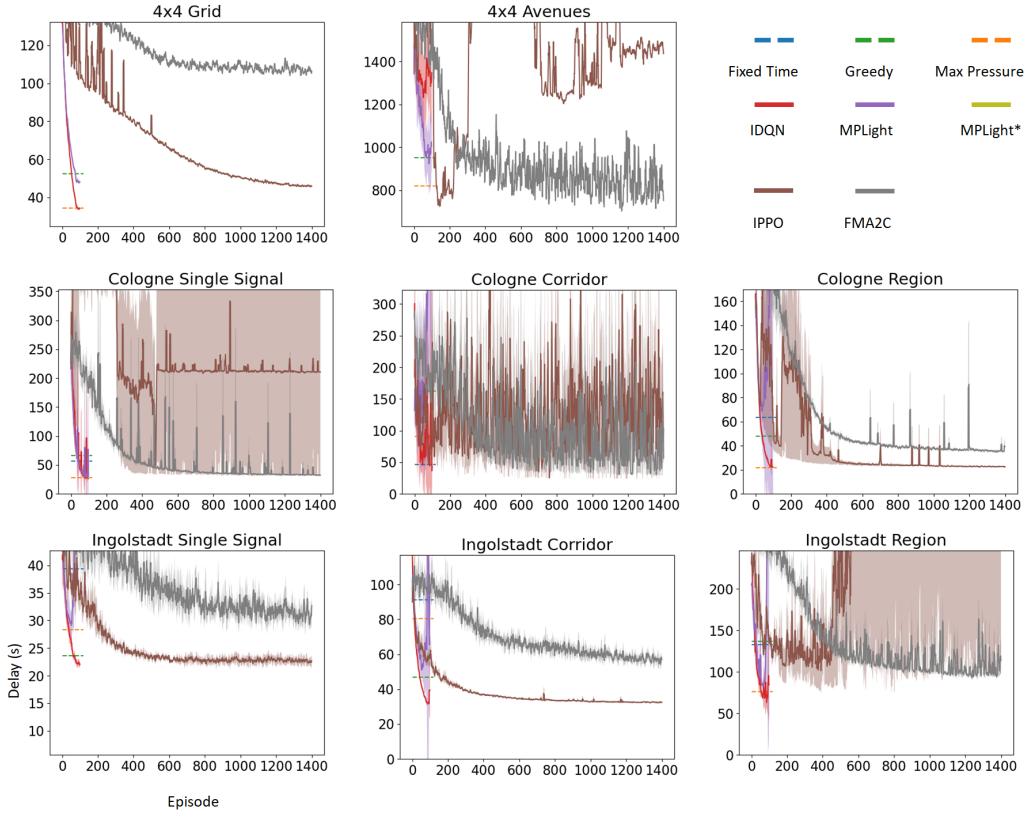


Figure 5: Full Learning curves over 5 random seeds

Table 3: Episodes before best performance

<i>Algorithm</i>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
IDQN	94	93	85	94	82	98
IPPO	1052	1373	386	475	744	1394
MPLight	<b>*76</b>	<b>59</b>	<b>71</b>	<b>83</b>	<b>34</b>	<b>35</b>
FMA2C	1361	1339	1372	1230	1375	1380

Table 4: Average standard deviation over all trials on benchmark scenarios

<b>IDQN</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	1.56	3.27	17.82	37.03	77.93	6.14
Avg. Trip Time	1.22	2.59	18.4	31.34	75.68	5.64
Avg. Wait	0.91	2.08	15.83	30.28	77.35	4.11
Avg. Queue	0.11	0.1	0.36	3.25	2.71	0.26
<b>IPPO</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	1.14	0.88	196.31	280.04	72.07	13.27
Avg. Trip Time	0.9	0.74	191.44	277.68	71.8	13.95
Avg. Wait	0.68	0.56	190.07	282.48	73.2	12.41
Avg. Queue	0.06	0.03	1.84	16.46	4.2	0.75
<b>MPLight</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	*3.5	32.03	70.11	25.54	189.9	153.38
Avg. Trip Time	*2.64	22.43	68.51	20.78	183.75	146.09
Avg. Wait	*2.13	21.61	67.21	19.97	185.82	141.03
Avg. Queue	*0.13	0.65	1.05	1.86	6.5	4.23
<b>FMA2C</b>	Ing. Single	Ing. Corr.	Ing. Reg.	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	3.27	4.34	12.49	11.67	45.28	5.35
Avg. Trip Time	2.4	3.54	12.84	9.37	43.48	4.94
Avg. Wait	2.09	2.68	11.01	8.82	44.43	4.02
Avg. Queue	0.19	0.1	0.23	1.06	1.67	0.23

Table 5: Performance on benchmark scenarios with standard deviations

<b>IDQN</b>	Ing. Single	Ing. Corr.	Ing. Reg.
Avg. Delay	$21.48 \pm 0.56$	$31.19 \pm 0.97$	<b><math>59.64 \pm 2.13</math></b>
Avg. Trip Time	$35.29 \pm 0.48$	$68.69 \pm 0.72$	<b><math>197.23 \pm 2.18</math></b>
Avg. Wait	$3.93 \pm 0.25$	$8.71 \pm 0.56$	<b><math>20.19 \pm 1.48</math></b>
Avg. Queue	$0.43 \pm 0.01$	<b><math>0.67 \pm 0.03</math></b>	<b><math>0.8 \pm 0.05</math></b>
<b>IPPO</b>	Ing. Single	Ing. Corr.	Ing. Reg.
Avg. Delay	<b><math>19.85 \pm 0.21</math></b>	<b><math>30.7 \pm 0.98</math></b>	$67.65 \pm 19.49$
Avg. Trip Time	<b><math>34.19 \pm 0.26</math></b>	<b><math>68.34 \pm 0.77</math></b>	$205.44 \pm 8.4$
Avg. Wait	<b><math>3.21 \pm 0.28</math></b>	<b><math>8.2 \pm 0.52</math></b>	$26.45 \pm 15.12$
Avg. Queue	<b><math>0.39 \pm 0.02</math></b>	$0.71 \pm 0.0$	$1.15 \pm 0.56$
<b>MPLight</b>	Ing. Single	Ing. Corr.	Ing. Reg.
Avg. Delay	$28.31 \pm 0.8$	$48.21 \pm 4.24$	$78.16 \pm 2.06$
Avg. Trip Time	$41.07 \pm 1.01$	$76.58 \pm 1.42$	$215.72 \pm 2.21$
Avg. Wait	$8.27 \pm 0.97$	$15.05 \pm 1.73$	$34.57 \pm 1.78$
Avg. Queue	$0.61 \pm 0.03$	$1.34 \pm 0.06$	$1.48 \pm 0.04$
<b>FMA2C</b>	Ing. Single	Ing. Corr.	Ing. Reg.
Avg. Delay	$25.36 \pm 1.5$	$48.99 \pm 0.54$	$90.42 \pm 1.69$
Avg. Trip Time	$39.4 \pm 1.11$	$85.03 \pm 1.63$	$226.5 \pm 1.47$
Avg. Wait	$7.27 \pm 0.64$	$21.9 \pm 0.15$	$44.16 \pm 2.0$
Avg. Queue	$0.94 \pm 0.0$	$1.79 \pm 0.02$	$1.74 \pm 0.04$
<b>IDQN</b>	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	<b><math>26.05 \pm 0.57</math></b>	$23.99 \pm 1.11$	$22.06 \pm 0.36$
Avg. Trip Time	<b><math>43.59 \pm 0.52</math></b>	$59.0 \pm 0.87$	$86.02 \pm 0.39$
Avg. Wait	<b><math>7.98 \pm 0.35</math></b>	$8.5 \pm 0.59$	$5.46 \pm 0.2$
Avg. Queue	<b><math>2.09 \pm 0.1</math></b>	$0.87 \pm 0.02$	$0.38 \pm 0.02$
<b>IPPO</b>	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	$55.07 \pm 11.83$	<b><math>22.13 \pm 0.41</math></b>	<b><math>21.49 \pm 0.13</math></b>
Avg. Trip Time	$67.7 \pm 9.27$	<b><math>57.45 \pm 0.26</math></b>	<b><math>85.54 \pm 0.17</math></b>
Avg. Wait	$26.15 \pm 6.62$	<b><math>7.37 \pm 0.36</math></b>	<b><math>5.01 \pm 0.17</math></b>
Avg. Queue	$8.88 \pm 2.74$	<b><math>0.76 \pm 0.1</math></b>	<b><math>0.35 \pm 0.0</math></b>
<b>MPLight</b>	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	$28.74 \pm 1.67$	$83.65 \pm 27.96$	$60.42 \pm 20.17$
Avg. Trip Time	$45.85 \pm 1.14$	$102.3 \pm 21.53$	$123.93 \pm 20.43$
Avg. Wait	$8.61 \pm 0.65$	$46.25 \pm 19.27$	$30.34 \pm 15.48$
Avg. Queue	$2.45 \pm 0.24$	$5.4 \pm 1.94$	$2.33 \pm 0.97$
<b>FMA2C</b>	Col. Single	Col. Corr.	Col. Reg.
Avg. Delay	$30.12 \pm 0.3$	$25.37 \pm 0.17$	$33.28 \pm 0.49$
Avg. Trip Time	$47.31 \pm 0.17$	$61.68 \pm 0.21$	$97.53 \pm 0.43$
Avg. Wait	$11.23 \pm 0.01$	$11.3 \pm 0.03$	$14.19 \pm 0.88$
Avg. Queue	$3.11 \pm 0.05$	$1.68 \pm 0.04$	$0.98 \pm 0.06$

Table 6: Algorithm Parameters

<b>IDQN</b>			
Optimizer: Adam	Adam $\beta_1$ : 0.9	Adam $\beta_2$ : 0.999	Adam $\epsilon$ : 1e-8
Step Size: 1e-3	Replay Size: 10,000	Minibatch Size: 32	$\gamma$ : 0.99
Target Update Interval: 500 Steps			
<b>MPLight</b>			
Optimizer: Adam	Adam $\beta_1$ : 0.9	Adam $\beta_2$ : 0.999	Adam $\epsilon$ : 1e-8
Step Size: 1e-3	Replay Size: 10,000	Minibatch Size: 32	$\gamma$ : 0.99
Target Update Interval: 500 Steps			
<b>FMA2C</b>			
Optimizer: RMSProp	RMSProp $\alpha$ : 0.99	RMSProp $\epsilon$ : 1e-5	Step Size: 2.5e-4
Entropy Coef: 1e-3	Entropy Ratio: 0.5	$\gamma$ : 0.96	Value Coef: 0.5
Batch Size: 120	Reward Norm 2,000	Reward Clip: 2	
<b>IPPO</b>			
Optimizer: Adam	Adam $\beta_1$ : 0.9	Adam $\beta_2$ : 0.999	Adam $\epsilon$ : 1e-5
Step Size: 2.5e-4	$\gamma$ : 0.99	Entropy Coef: 1e-3	Max L2 Norm: 0.5
Update Interval: 1024	Minibatch Size: 256	Epochs 4	$\lambda$ : 0.95
$\epsilon$ Clip: 0.1			