

1 **REINFORCEMENT LEARNING-BASED TRAFFIC SIGNAL CONTROL IN SPECIAL**  
2 **SCENARIO**

3 **Dingyi Zhuang**

4 Department of Civil Engineering and Applied Mechanics  
5 McGill University  
6 Montreal, Quebec, H3A 0C3, Canada  
7 Email: dingyi.zhuang@mail.mcgill.ca  
8 ORCID: 0000-0003-3208-6016

9 **Zhenyuan Ma**

10 Department of Civil Engineering and Applied Mechanics  
11 McGill University  
12 Montreal, Quebec, H3A 0C3, Canada  
13 Email: zhenyuan.ma@mail.mcgill.ca

14 **Lijun Sun, Corresponding Author**

15 Department of Civil Engineering and Applied Mechanics  
16 McGill University  
17 Montreal, Quebec, H3A 0C3, Canada  
18 Email: lijun.sun@mcgill.ca  
19 ORCID: 0000-0001-9488-0712

20 Word count: 5482 words text + 0 table(s) x 250 words (each) = 5482 words

21 Submission Date: August 1, 2020

1 **ABSTRACT**

2 Traffic signal control (TSC) is an important task to optimize the performance, both economic and  
3 environmental, of the transportation system. Recently reinforcement learning (RL) techniques  
4 have shown promising results in optimizing control policies for basic and standard intersections,  
5 which only receive traffic flow from connected arterials. However, the typical urban traffic net-  
6 works are composed of freeways and arterial. The exchange flow between those two always causes  
7 congestion. This paper studies the effects of this type exchange flow on reinforcement learning-  
8 based TSC controllers. The first studied traffic network is built upon standard intersections. Then  
9 we build another network by adding an freeway off-ramp to the first network. Several up-to-date  
10 deep RL algorithms are used to learn optimal TSC signal for those two networks. Experimental  
11 results show the exchange flow between freeways and arterials greatly influence the learning re-  
12 wards and safety. Positions and length settings of the ramp are also discussed, showing that RL  
13 results are sensitive to the infrastructure settings. We appeal for more attention on the practical  
14 traffic scenario discussion when conducting RL-based traffic control.

15 *Keywords:* Reinforcement learning, Traffic light control, Freeway off-ramp scenario

## 1 INTRODUCTION

2 The steadily increasing number of population and economic activities has led to the rise of urban  
3 congestion. Better solutions for TSC problem on the road will result in less traffic congestion and  
4 safer driving environment, which improves the efficiency and safety of the transportation system.  
5 Several advanced control methods have been widely used for reducing congestion. One promising  
6 strategy is the use of advanced RL-based TSC methods. In the RL context, the control problem will  
7 be modelled as an Markov Decision Process (MDP) with sequentially state-action-reward tuples.  
8 Observations change according to the selected action under a certain policy, with reward received  
9 to judge the policy. Back to transportation field, how to formulate the MDP framework is one of  
10 the most important tasks. Intuitively, we will consider the traffic lights as agents, and their phases  
11 as the states (1, 2). Whether the traffic lights change the phase and which phase they change to will  
12 be the action set. The objective is to find a policy that that maximize the expected return, e.g. the  
13 delay of vehicles.

14 However, current researches focus on improving the design of RL formulation and agent  
15 modeling, and they tend to apply models on larger road network to test the coordination perfor-  
16 mance. For instance, Wei et al. (1) consider a traffic road network with 196 intersections and design  
17 the cooperation algorithm among them. However, this large-scale road network only consists of  
18 arterial blocks and grids, where most of the intersections are simple ordinary 4-way intersections  
19 (i.e. crossroads) connected with typical arterial. This kind of road network is common in cities, but  
20 the real-world urban network is composed of freeways and arterial roads. The existing RL studies  
21 only consider the typical arterial roads, the coordination of traffic flow between the freeway and ar-  
22 terial are not fully studied. Lim et al. (3) have demonstrated that the existence of freeway off-ramp  
23 generates the spillover to the mainstream traffic flow, which decreases the freeway vehicle capac-  
24 ities and influence the TSC of the overall functions of urban freeways. Therefore, proper signal  
25 control in the neighboring intersection is more desired, which can save average delay between 15%  
26 and 40%, and average travel time up to 25% according to U.S. Federal Highway Administration  
27 (4).

28 Inspired by that, our goal of this study is to test RL models on a special traffic scenario in  
29 which the arterial network is connected with a freeway off-ramp. We will analyze the impact on RL  
30 model performance and road safety given changes in traffic flow or infrastructure influence (like  
31 the length of ramps) to discuss whether small changes in the environment setting will influence  
32 RL greatly. Our work is dedicated to examine whether RL models trained on simplified scenarios  
33 without discussion of infrastructures is reasonable. One important finding is that simply adding a  
34 freeway off-ramp will greatly influence the final rewards, which is much higher than increasing the  
35 same amount of inflow vehicles. We hope our work can appeal for more attention on the practical  
36 traffic scenario discussion rather than more fancy algorithm design on RL-based TSC.

37 The remainder of this paper is as follows. In section 2 we review some literature about  
38 classical signal control methods and reinforcement learning based ones. Section 3 describes our  
39 special traffic scenario design and how our RL model works. Section 4 shows all our experiments  
40 and results. Finally in section 5, we conclude our work.

## 41 LITERATURE REVIEW

### 42 Classical Methods

43 There exist a set of classical transportation methods for traffic signal control problems. They are  
44 based on the knowledge of signal planning and phase changing rules, which can be intuitively

1 understood and explained. Generally, these classical methods work as baselines comparing to RL-  
2 based methods. Most of the classical methods use equations related to groundtruth of physics, such  
3 as relationship between time, length and speed.

4 GreenWave (5) is one of the most classical methods focusing on unidirectional traffic signal.  
5 It optimizes an offset for signals in each intersection by reducing the number of stops of vehicles  
6 when traveling in one certain direction. However, it requires a pre-defined signal plan applied to all  
7 intersections. A similar method named Maxband (6) has nearly the same setting with GreenWave.  
8 The difference occurs in the objective, as Maxband considers vehicles traveling along two opposite  
9 directions.

10 Actuated control is a more flexible way than static signal control, and is a more controllable  
11 way than RL-based signal control. There are rules for changing signal, by detecting whether there  
12 exists a sufficient time gap between current vehicle and its successive vehicle. A maximum and  
13 a minimum time limit is needed for phases, so that there will be balance for each direction in the  
14 intersection. This control method is commonly used in Germany.

15 Queue length of an intersection can be considered as "pressure" of the intersection. Max-  
16 pressure control is used to balance the pressure of each intersection and avoid over-saturation  
17 where the pressure is formally defined as the difference of queue length between incoming lanes  
18 and outgoing lanes (7).

## 19 Reinforcement Learning Methods

20 Because of more data resources and more computational power, RL-based methods have been  
21 applied widely in traffic signal control field. Since the traffic scenario is depicted via MDP, as  
22 discussed above, we will consider the traffic lights as agents, and their phases as the states (1, 2).  
23 Whether the traffic lights change the phase and which phase they change to will be the design of  
24 action. Usually, the reward of the algorithm is related to the total delay of vehicles.

25 However, there are much more possible settings of state, action and reward of RL problem.  
26 The length of waiting queue in front of traffic light (1), the level of congestion happened in a lane  
27 (8), waiting time of vehicles (9) could also be defined as the reward of agents. Although the goal of  
28 Reinforcement Learning problem is to maximize the expected return, according to different state  
29 setting, the reward function can be different. A reward function can also be a linear combination of  
30 several factors, so that we will know more about how the influence of factors to the traffic signals.  
31 From different perspectives, RL-based algorithms can be categorised differently. We can choose  
32 the most suitable algorithm by considering whether we need to learn the state transition function,  
33 whether we learn the policy parameter, and whether we learn by using a table for state-action pairs.

34 Game theory is a useful tool that can be applied into Reinforcement Learning methods, es-  
35 pecially multi-agent RL model. By using the knowledge of game theory, signals or traffic lights are  
36 considered as players in a game, and they choose their best actions according to their payoffs (i.e.  
37 objectives in traffic signal control setting) (10). Finding an equilibrium in a game theory is similar  
38 to the communication between agents in Reinforcement Learning and find the optimal solution.  
39 Moreover, different game theory models can be considered to fit different traffic environments and  
40 different settings of agents.

41 RL-based method can be also used together with a classical one. There is an example of  
42 doing so which is PressLight (11), using the idea of max-pressure (7). By this way, the knowledge  
43 in transportation field can be well used and make reinforcement learning methods more intuitive,  
44 as well as outperform both approaches.

## 1 Special Traffic Scenarios

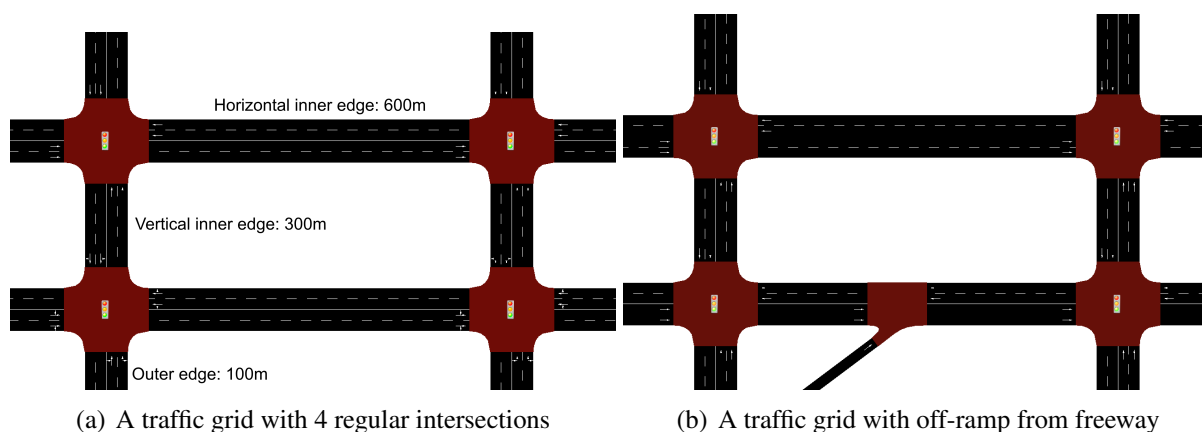
2 Comparing to the road network that only consists simple ordinary 4-way intersections (i.e. cross-roads) (1), special traffic scenarios are less common appeared, both in our daily life and in research area. This means special traffic scenarios are worth to be researched into. T-junction is one of the 5 simplest intersections, but it contains more possible phase design than common crossroad (12). A 6 combination of two T-junctions with different directions is called LR-type staggered intersection 7 (13). Congestions and even accidents tend to happen more in those non-regular traffic area.

8 Bottlenecks and intersections near freeways are another type of special traffic scenarios. 9 Vehicles in both scenarios will face different width of roads before and after crossing an intersec- 10 tions. Moreover, due to different speed limit in different type of roads, traffic collisions appear 11 more than other type of traffic scenarios. Existing research shows that if we apply traffic signal 12 control, even basic control system, there will be a reduce of around 20% of collisions. (14)

13 Not only special traffic networks affect, some traffic events matter in traffic signal control. 14 (15) and (16) mention there may be uncertain supply happened as traffic disruptions on links such 15 as incidents, freeway crash and lane closure. These events can be treated as part of special traffic 16 scenarios and will affect the control of signal and the traffic flow of the total traffic networks.

## 17 PROBLEM STATEMENT

### 18 Scenario Design



**FIGURE 1** : Scenario design that (b) simply adds an off-ramp from freeway in the traffic grid of (a). If the traffic flow in the ramp is 0, these two environments are exactly the same for RL model.

19 We design a simple traffic grid that influenced by a freeway off-ramp shown in Figure 20 1(b). This grid is a generalization from regular traffic grid, which contains four regular crossroad 21 intersections, as displayed in Figure 1(a). The scenario simulates how vehicles exiting from the 22 freeway will merge the city freeway traffic flow. Original road setting is 300 meters with 4 lanes 23 for every vertical inner edges, and 600 meters with 4 lanes for every horizontal inner edges. For 24 outer edges we only show 100 meters for simplicity. For each regular intersection, we allow traffic 25 inflow to come from two outer edges of the intersection with 25% probability to add new vehicle 26 every timestep. Moreover, we can set random number of vehicles in the grid when initializing the 27 scenario. If the vehicle flow in off-ramp is 0, meaning that no vehicles added in the ramp, then 28 scenarios with or without the ramp are the same for RL models. In each regular intersection, we

1 set a traffic light with specified phases, without setting the time limit of phases. The traffic lights  
2 are what we treat as agents in our RL formulation. There is no traffic light in the area that ramp is  
3 connected with the grid in the very basic setting.

4       Some variables in this scenarios can be changed and then used to compare how the environ-  
5 ment change will affect the traffic, which serve as a reference when government design their real  
6 traffic environment. Changeable variables include the inflow rate of vehicles in outer edges and  
7 ramp, the length of ramp, the position of ramp (i.e. how close it is to the right bottom intersection).  
8 Other variables such as the length of inner edges, the speed limit of vehicles or the number of lanes  
9 will also influence, but with smaller effect.

## 10 **Environment Design**

11 Working as a reinforcement learning problem, we need to specify the environment setting of this  
12 Markov Decision Process. Key elements in the environment includes agent, state, action and  
13 reward are designed as:

14       • **State:** Our state is defined as vehicle information for each intersection. The state com-  
15 ponent includes which intersection is selected, the speed of each vehicle, the distance  
16 between each vehicle and selected intersection and the edge where each vehicle is on.  
17 This state design is the same as the work of Wu et al. (17).

18       • **Action:** Because we design the action phase to make sure that vehicles can only pass the  
19 intersections horizontally or vertically, the action space  $\mathbf{a}$  only consists of a list of float  
20 variables ranging from 0 to 1, specifying whether a traffic light is supposed to switch or  
21 not. The actions are sent to the traffic light in the grid from left to right and then top to  
22 bottom.

23       • **Reward:** Reward is calculated by the negative per vehicle delay minus a penalty for  
24 switching traffic lights.

25       • **Policy agent:** Since it is an environment with multiple agents (lights) and observations,  
26 we can consider both "single-agent" and "multi-agent" policy environment. Here we use  
27 the term "agent" again for policy trainer as it is defined in the RLlib (18). This agent  
28 definition is only used in algorithm discussion. If these is an super-agent controlling  
29 the observation of four traffic lights, we treat this environment as single-agent policy  
30 environment. Otherwise, in multi-agent environment, agents will model each light's own  
31 goal and action.

32       • **Termination:** An epoch is terminated if the expected timesteps are reached.

## 33 **METHODOLOGY**

34 Our research focuses more on the special transportation scenario and its comparison. Conse-  
35 quently, we will implement widely-applied methods from RLlib (18, 19). In this section, we  
36 would firstly introduce the vanilla policy gradient algorithm. Then we discuss the Asynchronous  
37 Actor-Critic Agents (A3C) and Proximal Policy Optimization (PPO) algorithms, which belongs  
38 to Policy Gradient method family but more efficient. These three algorithms are the basic models  
39 discussed in the experiment part.

## 1 Policy Gradient

2 The goal of reinforcement learning is to find an optimal strategy for agents to achieve the optimal  
3 rewards. As an important member, policy-gradient algorithms do not suffer from issues that value-  
4 based RL methods might have, such as complexity problem from the continuous states and actions  
5 (20). The basic idea of policy gradient is that, by sampling trajectories (sequences of state, action  
6 and reward tuples) from the environment from current policy function, we can gather the rewards  
7 with different trajectories where high-reward paths receive more weights.

8 Suppose we have a differentiable policy function  $\pi_\theta(\mathbf{a}|\mathbf{s})$ , where  $\mathbf{a}$  is the action,  $\mathbf{s}$  is the  
9 state, and  $\theta$  represents the weight parameters to be learned. Suppose the reward function is  $r(\mathbf{s}, \mathbf{a})$ .  
10 The optimal policy can be obtained by

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [r(\tau)] = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

11 where  $\tau$  denotes the sampled trajectory of  $(\mathbf{s}_t, \mathbf{a}_t)$  tuples from policy  $\pi_\theta$ . The optimization  
12 process can be approximated by the gradient descent method. Since we are updating the weight  
13 parameters of policy  $\pi_\theta$ , the method is referred as Policy Gradient. The updating rule can be  
14 approximated by (21):

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \left( \sum_t \nabla_{\theta} \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[ \left( \sum_t \nabla_{\theta} \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \right] \end{aligned} \quad (1)$$

15 The vanilla Policy Gradient algorithm is summarized in Algorithm 1. Within each epoch,  
16 step 2 samples state-action-reward trajectories and step 3 and 4 calculate policy updating where  
17  $J(\theta)$  stands for the objective function we evaluate. The process can assign larger weights on  
18 high-reward samples to make them easier to be sampled. Thus, these well performed samples will  
19 further improve the policy updating.

---

### Algorithm 1 Vanilla Policy Gradient

---

Initialize: A differentiable policy function  $\pi_\theta(\mathbf{a}|\mathbf{s})$ , learning rate  $\alpha$

1: **for** each epoch **do**

2: Sample trajectory  $\mathbf{s}_{i,t}, \mathbf{a}_{i,t}, r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}), \mathbf{s}_{i,t+1}$  from  $\pi_\theta(\mathbf{s}|\mathbf{a})$  in the simulator

3: Calculate  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[ \left( \sum_t \nabla_{\theta} \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \right]$

4:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

5: **end for**

---

## 20 Actor-Critic Agents

21 Even though policy gradient methods enforce good sampling, the sampled rewards in practice  
22 usually have large variances, inducing large variance of the estimated gradients. Therefore, the  
23 vanilla policy gradient method is unstable and inefficient.

24 A3C add asynchronous mechanism where each agents interacts with its respective environ-  
25 ments asynchronously, learning interaction with the others. While the core idea is still Actor-Critic,

1 which solves the variance problems by adding a value function, as "critic", to the vanilla policy  
 2 gradient algorithm. It evaluates the value function of the current policy by observing the states and  
 3 rewards feedback from the environment. The policy  $\pi_\theta$  is the "actor" that generates the trajectory  
 4 in Algorithm 1. Before updating the policy, the value function approximates the expected rewards  
 5 using methods like neural network using the generated samples. The value function is used to as-  
 6 sist in improving the policy's parameter  $\theta$  updating by comparing the average value of the current  
 7 state.

---

**Algorithm 2** Actor-critic
 

---

Initialize: A differentiable policy function  $\pi_\theta(\mathbf{a}|\mathbf{s})$ , learning rate  $\alpha$

1: **for** each epoch **do**

2:   Sample trajectory  $\mathbf{s}_{i,t}, \mathbf{a}_{i,t}, r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}), \mathbf{s}_{i,t+1}$  from  $\pi_\theta(\mathbf{s}|\mathbf{a})$  in the simulator

3:   fit  $\hat{V}_\phi^\pi(\mathbf{s})$  to sampled rewards

4:   Evaluate the advantage as  $\hat{A}^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) = r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1}) - \hat{V}_\phi^\pi(\mathbf{s}_{i,t})$

5:   Calculate  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N [(\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})) \hat{A}^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})]$

6:    $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

7: **end for**

---

8           As shown in Algorithm 2, step 3 updates the value function using the generated states-  
 9 actions-reward tuples and then used to calculate the advantage  $\hat{A}^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$  in step 4. Advantage  
 10 estimates the improvement of the corresponding action achieves compared to the average value  
 11 function we have estimated. Thus, by multiplying the gradient of the policy function with the  
 12 advantage in step 5, we can update the parameters so that high-reward actions will be selected  
 13 more likely.

### 14 Proximal Policy Optimization

15 In practice, vanilla policy gradient is not scalable for large scale problems. To solve the prob-  
 16 lem of hard convergence in vanilla policy gradient, Schulman et. al propose PPO algorithm  
 17 with adaptive Kullback–Leibler divergence (KL-divergence) (22). KL-divergence is defined as  
 18  $\bar{D}_{KL}(P||Q) \approx -\sum_i P(i) \ln \frac{Q(i)}{P(i)}$ , which penalizes the objective if the new policy  $P$  is different from  
 19 the old one  $Q$ . Details can be found in Algorithm 3.

20           It can be found that the key is step 4 where KL-divergence constraints ensure that the new  
 21 policy and the old one would not differ greatly. In this way, the convergence of the algorithm can  
 22 be guaranteed.

## 23 EXPERIMENTAL RESULTS

24 For experiments, we customized scenario network and both single and multi-agent environments  
 25 based on Berkeley Flow project (17). Simulation is conducted in SUMO (23) with 3000 timesteps  
 26 in each epoch. Each case is run 5 times under different random seeds.

### 27 Reinforcement Learning Experiments

28 We have run three basic algorithms with RLlib package (19) for both single-agent and multi-agent  
 29 policy environment including Asynchronous Actor-Critic Agents (A3C), Policy Gradient (PG), and  
 30 Proximal Policy Optimization (PPO). As mentioned, here the term "single-agent" means there is a



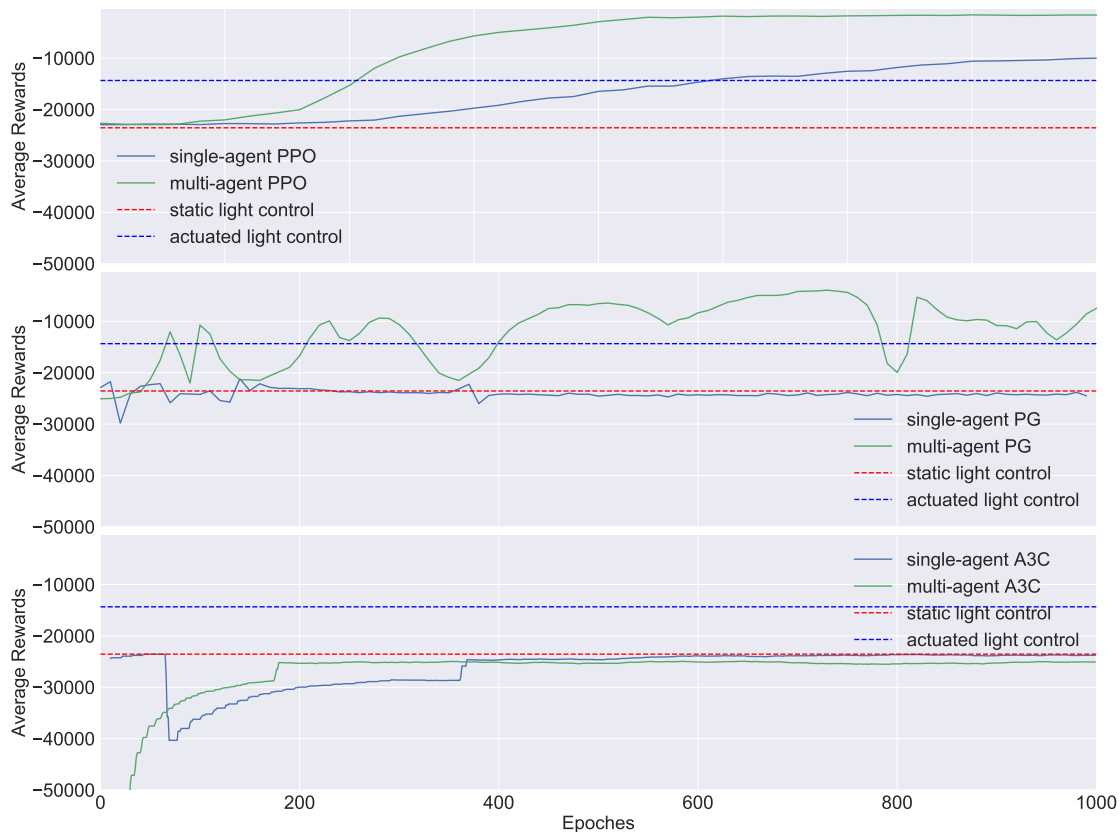
**Algorithm 3** PPO

---

Initialize: policy parameter  $\theta_0$ , KL penalty  $\beta_0$ , target KL-divergence  $\delta$

- 1: **for**  $k=0,1,2,\dots$  **do**
- 2:   Sample trajectory  $\mathbf{s}_{i,t}, \mathbf{a}_{i,t}, r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}), \mathbf{s}_{i,t+1}$  from  $\pi_{\theta_k}(\mathbf{s}|\mathbf{a})$  in the simulator
- 3:   Evaluate the advantage as  $\hat{A}^{\pi_{\theta_k}}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) = r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t})$
- 4:   Compute policy update  $\theta_{k+1} = \arg \min_{\theta} J(\theta_k) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$  by taking K steps of minibatch stochastic gradient descent.
- 5:   **if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$  **then**
- 6:      $\beta_{k+1} = 2\beta_k$
- 7:   **else if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta/1.5$  **then**
- 8:      $\beta_{k+1} = \beta_k/2$
- 9:   **end if**
- 10: **end for**

---



**FIGURE 2 :** Comparison of multi-agent and single-agent in three algorithms.

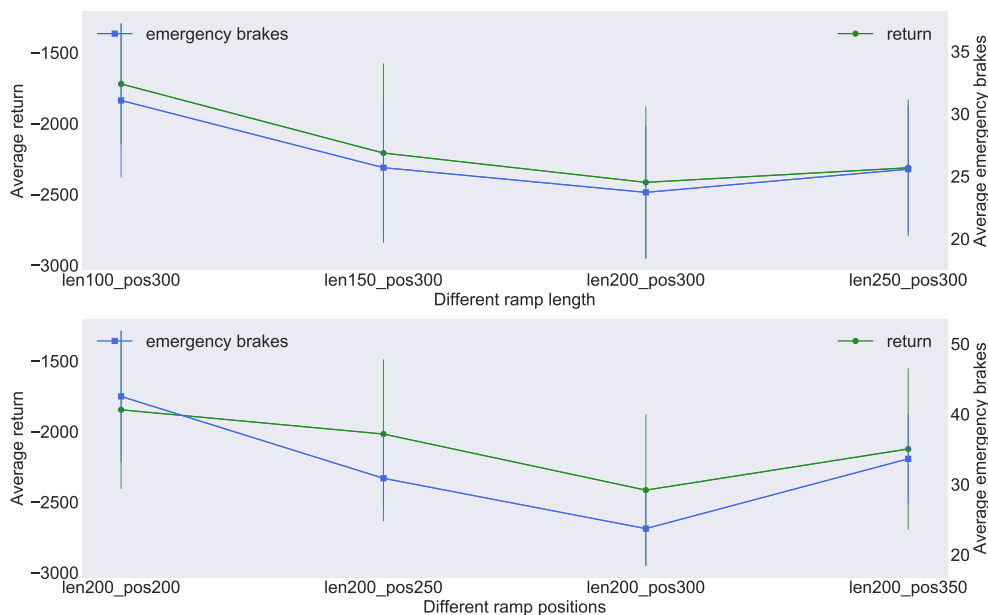
- 1 super agent controls the policy of 4 lights, instead of meaning we only use one traffic light's
- 2 observation. Contrarily, "multi-agent" means we train separately policies for different observations
- 3 of the lights. Besides the RL control, There are two non-RL baseline methods comparing to three
- 4 RL algorithms mentioned before. The first one is static light phase length. We set all green light
- 5 phases with a length of 31 seconds and all yellow phases with 3 seconds. The second baseline

1 method is actuated traffic signals which is more commonly used in Germany. It will prolong traffic  
 2 phases whenever a continuous stream of traffic is detected. and it switches to the next phase after  
 3 detecting a sufficient time gap between successive vehicles.

4 Figure 2 show the mean rewards of PPO, PG and A3C algorithms respectively. From the  
 5 result, we observe that A3C is not suitable for this traffic scenario. It even performs worse than  
 6 static light control. And in A3C case, there is no much difference between single and multi-  
 7 agent algorithms. Multi-agent PG works better than its single version, and it may be used in this  
 8 traffic signal control problem, since it at least performs better than actuated signal control. But  
 9 the problem is that it will not converge. Actually, PPO works quite well, especially in multi-agent  
 10 environment. After convergence, multi-agent PPO performs 10 times better than the single one.

11 Because our task is to discuss the influence from infrastructural settings and the traffic flow  
 12 (i.e. the independent variables), we only need to pick up the best performance algorithm as the  
 13 base model. It can be found that multi-agent PPO suits this problem most, thus for the following  
 14 infrastructure experiments, we keep on using this model setting as default, and evaluate on how  
 15 infrastructure settings like ramp length and position affect.

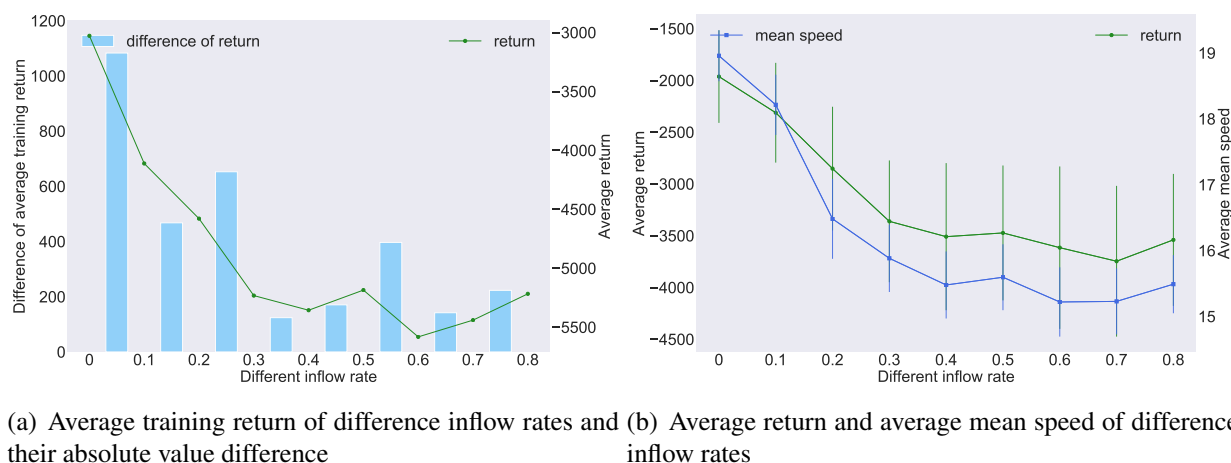
## 16 Infrastructure Experiments



**FIGURE 3** : Average return and number of emergency brakes under different ramp infrastructure settings

17 We then run experiments with different network settings to discuss how infrastructure of  
 18 special scenarios will influence the RL control. The length of the ramp varies from 100 meters to  
 19 250 meters, and the position of the ramp to the right bottom intersection varies from 200 meters  
 20 to 350 meters, both with 50 meters as intervals. In addition to the returns, we would like to apply  
 21 reinforcement learning on traffic signal control to bring efficiency along with safety. Herein, we  
 22 use the number of emergency brakes happened in each epoch as safety index. Figure 3 shows  
 23 the results of average return and the number of emergency brakes under different infrastructure

1 settings after 100 epochs training with 5 different random seeds.  
 2 We observe that shortest ramp length gives best training reward, while there may not be  
 3 large difference when ramp is too long (e.g. more than 200 meters). Slightly closer to the one of  
 4 the intersections will have positive impact on traffic efficiency. Because the vehicles from ramp  
 5 will turn right, the return of ramp closer to the left intersection is better than the return of the one  
 6 closer to the right intersection. This result can be used as a simulation reference for infrastructure  
 7 building. For safety issue, we find out that the average return and the number of emergency brakes  
 8 have similar trend. However, larger the return is better, while less the emergency brakes is better.  
 9 The setting that gives best training reward, which is ramp with length 100 and position 300, does  
 10 not guarantee the safest control. The reason could be that balanced infrastructure setting, such as  
 11 position 300 which denotes the ramp in the middle of the inner edge, will influence more on the  
 12 flow by separating the road equally, but will make less emergency brakes also by separating the  
 13 road equally.



(a) Average training return of difference inflow rates and (b) Average return and average mean speed of difference inflow rates

**FIGURE 4** : Discussion of the impacts of inflow rate on training rewards.

14 We still need to answer why we use special scenarios as our research targets. The experi-  
 15 ments on different inflow rates of ramp edge reflect our main motivation working on special traffic  
 16 scenarios. From figure 4(a) and 4(b), we can clearly find that there is a large gap between inflow  
 17 probability from 0 to 0.1. If we set the probability as 0, it means there is no inflow from ramp, thus  
 18 this scenario becomes an ordinary traffic grid with only four crossroads. Once we add 0.1 prob-  
 19 ability of inflow to the ramp, the rewards will decrease a lot, comparing to other 0.1 probability  
 20 change (e.g. from 0.1 to 0.2). When the probability of inflow is more than 0.4, there is no much  
 21 difference between the return, as well as the average mean speed, also seen from Figure 4(b). We  
 22 have run several simulations and then noticed that when the inflow probability rate is greater than  
 23 0.4, there will be congestion on the ramp. In this case, the traffic on freeway that generate the ramp  
 24 will be affected more than the traffic grid that is connected to the ramp. Thus, the return we get  
 25 that reflect the traffic situation for only grid area will not follow a decrease trend and change a lot  
 26 when the inflow probability is greater than 0.4.

## 27 CONCLUSIONS

28 In this project, we manage to build a special traffic scenario to see how special scenario like adding  
 29 an off freeway ramp to regular intersections will affect RL training performance from both training

1 returns and safety. By applying multi-agent PPO we have found that adding a freeway off-ramp  
2 will greatly influence the model performance, where difference is much higher than increasing  
3 the same amount of inflow rate. Also, we further explore the infrastructure influence, and run  
4 experiments on different ramp length and position combination. We find that shorter ramp length  
5 will promote the training reward while ramp position in the middle of the road will undermine the  
6 training performance. Therefore, the infrastructure designs have unignorable impacts in RL-based  
7 TSC problem. Thus more attention should be paid upon the infrastructure design methods for  
8 RL-based signal controllers.

9 Several interesting questions stem from our paper, which we plan to study in the future.  
10 For safety issue, we would like to add safety index directly into reinforcement learning reward,  
11 instead of just consider the number of emergency brakes. More special scenarios should be built  
12 and analyzed. We expect to find the most compatible algorithm and infrastructure for each special  
13 scenario. Autonomous vehicle is another direction. How to treat vehicles and traffic lights as  
14 agents together in one reinforcement learning control system is a challenge.

## 15 **ACKNOWLEDGEMENT**

16 This research is supported by the Natural Sciences and Engineering Research Council (NSERC)  
17 of Canada, Fonds de recherche du Quebec - Nature et technologies (FRQNT), and the Canada  
18 Foundation for Innovation (CFI). The authors declare no competing financial interests with respect  
19 to this work.

## 20 **AUTHOR CONTRIBUTIONS**

21 D.Z. and Z.M. designed and performed the research; D.Z., Z.M. and L.S. wrote the paper.

## 22 **REFERENCES**

- 23 [1] Wei, H., N. Xu, H. Zhang, G. Zheng, X. Zang, C. Chen, W. Zhang, Y. Zhu, K. Xu, and Z. Li,  
24 CoLight: Learning Network-Level Cooperation for Traffic Signal Control. In *Proceedings of*  
25 *the 28th ACM International Conference on Information and Knowledge Management*, Asso-  
26 ciation for Computing Machinery, New York, NY, USA, 2019, CIKM '19, p. 1913–1922.
- 27 [2] Wei, H., G. Zheng, H. Yao, and Z. Li, Intellilight: A reinforcement learning approach for  
28 intelligent traffic light control. In *Proceedings of the 24th ACM SIGKDD International Con-*  
29 *ference on Knowledge Discovery & Data Mining*, 2018, pp. 2496–2505.
- 30 [3] Lim, K., J. H. Kim, E. Shin, and D.-G. Kim, A signal control model integrating arterial  
31 intersections and freeway off-ramps. *KSCE Journal of Civil Engineering*, Vol. 15, No. 2,  
32 2011, pp. 385–394.
- 33 [4] Federal Highway Administration, U. D. o. T., *Ramp Management and Control Handbook*,  
34 2006.
- 35 [5] Roess, R., E. Prassas, and W. McShane, *Traffic Engineering*. Traffic Engineering, Pear-  
36 son/Prentice Hall, 2004.
- 37 [6] Little, J. D. C., M. D. Kelson, and N. H. Gartner, MAXBAND : a versatile program for setting  
38 signals on arteries and triangular networks, 1981.

- 1 [7] Varaiya, P., The Max-Pressure Controller for Arbitrary Networks of Signalized Intersections.  
2 *Advances in Dynamic Network Modeling in Complex Transportation Systems*, 2013, pp. 27–  
3 66.
- 4 [8] Bakker, B., S. Whiteson, L. Kester, and F. Groen, *Traffic Light Control by Multiagent Rein-*  
5 *forcement Learning Systems*, Vol. 281, pp. 475–510, 2010.
- 6 [9] Wei, H., G. Zheng, H. Yao, and Z. Li, IntelliLight: A Reinforcement Learning Approach for  
7 Intelligent Traffic Light Control, 2018, pp. 2496–2505.
- 8 [10] Elhenawy, M., A. A. Elbery, A. A. Hassan, and H. A. Rakha, An Intersection Game-Theory-  
9 Based Traffic Control Algorithm in a Connected Vehicle Environment. In *2015 IEEE 18th*  
10 *International Conference on Intelligent Transportation Systems*, 2015, pp. 343–347.
- 11 [11] Wei, H., C. Chen, G. Zheng, K. Wu, V. Gayah, K. Xu, and Z. Li, PressLight: Learning  
12 Max Pressure Control to Coordinate Traffic Signals in Arterial Network. In *Proceedings of*  
13 *the 25th ACM SIGKDD International Conference on Knowledge Discovery Data Mining*,  
14 Association for Computing Machinery, New York, NY, USA, 2019, KDD '19, p. 1290–1298.
- 15 [12] Borg, D. L. and K. Scerria, Constrained Dynamic Control of Traffic Junctions. *Procedia*  
16 *Computer Science*, Vol. 32, 2014, pp. 293 – 300, the 5th International Conference on Ambi-  
17 ent Systems, Networks and Technologies (ANT-2014), the 4th International Conference on  
18 Sustainable Energy Information Technology (SEIT-2014).
- 19 [13] Cai, Z., M. Xiong, D. Ma, and D. Wang, Traffic design and signal timing of staggered in-  
20 tersections based on a sorting strategy. *Advances in Mechanical Engineering*, Vol. 8, No. 4,  
21 2016, p. 1687814016641292.
- 22 [14] Cambridge Systematics, I., *Twin Cities Ramp Meter Evaluation*, ????
- 23 [15] Zhu, F. and S. V. Ukkusuri, Accounting for dynamic speed limit control in a stochastic traffic  
24 environment: A reinforcement learning approach. *Transportation Research Part C: Emerging*  
25 *Technologies*, Vol. 41, 2014, pp. 30 – 47.
- 26 [16] Aslani, M., M. S. Mesgari, and M. Wiering, Adaptive traffic signal control with actor-critic  
27 methods in a real-world traffic network with different traffic disruption events. *Transportation*  
28 *Research Part C: Emerging Technologies*, Vol. 85, 2017, pp. 732 – 752.
- 29 [17] Wu, C., A. Kreidieh, K. Parvate, E. Vinitzky, and A. M. Bayen, *Flow: A Modular Learning*  
30 *Framework for Autonomy in Traffic*, 2017.
- 31 [18] Liang, E., R. Liaw, R. Nishihara, P. Moritz, R. Fox, J. Gonzalez, K. Goldberg, and I. Sto-  
32 ica, Ray rllib: A composable and scalable reinforcement learning library. *arXiv preprint*  
33 *arXiv:1712.09381*, 2017.
- 34 [19] Liang, E., R. Liaw, R. Nishihara, P. Moritz, R. Fox, J. Gonzalez, K. Goldberg, and I. Sto-  
35 ica, Ray RLLib: A Composable and Scalable Reinforcement Learning Library. *CoRR*, Vol.  
36 abs/1712.09381, 2017.

- 1 [20] Mao, C., Y. Liu, and Z.-J. M. Shen, Dispatch of autonomous vehicles for taxi services: A deep  
2 reinforcement learning approach. *Transportation Research Part C: Emerging Technologies*,  
3 Vol. 115, 2020, p. 102626.
- 4 [21] Sutton, R. S. and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- 5 [22] Schulman, J., F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal policy optimiza-  
6 tion algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 7 [23] Lopez, P. A., M. Behrisch, L. Bieker-Walz, J. Erdmann, Y.-P. Flötteröd, R. Hilbrich,  
8 L. Lücken, J. Rummel, P. Wagner, and E. Wießner, Microscopic Traffic Simulation using  
9 SUMO. In *The 21st IEEE International Conference on Intelligent Transportation Systems*,  
10 IEEE, 2018.