# Ensemble Reinforcement Learning for Large-Scale Traffic Signal Control

### Zhenzhe Ying
443407384@qq.com
HangZhou, ZheJiang

### Zhuoer Xu
xuzhuoer.rex@gmail.com
HangZhou, ZheJiang

### Hui Li
lihuiknight@gmail.com
HangZhou, ZheJiang

### Haotian Wang
1031821435@qq.com
HangZhou, ZheJiang

### Shiwen Cui
278044960@qq.com
HangZhou, ZheJiang

## ABSTRACT

Traffic signal control is an important problem that affects people's daily life in commuting. However, it has been quite a challenge to design an effective policy for large-scale traffic signal control because of both computation and memory issues. In this paper, we formulate the problem as a Markov Decision Process (MDP). Specifically, we model each valid traffic light intersection as an agent and all the intersections (agents) need to make a *joint policy* coordinately so that the traffic system can serve more vehicles with less *delay*. However, learning the joint policy for the large-scale traffic signal control is intractable because its action space is exponentially large over the number of intersections. To address this issue, we employ deep reinforcement learning to approximate the *joint policy* for the traffic signal control based on a carefully designed MDP, where the reward is a self-defined *delay* index and the state is designed based on some heuristic methods. To make our policy generalize well on the unknown road networks, we learn multiple policies on some randomly generated road networks and archive the final policy via ensemble learning. We conduct our method on multiple partially observed traffic flows as well as a real-world road network: 1004 traffic lights in Nanchang, one of the largest cities in China. The experimental results show that our policy has a strong generalization and won second place in KDD CUP 2021 on city brain challenge from 1156 teams [1].

## 1 INTRODUCTION

Traffic congestion is one of the biggest issues in the city traffic, however, it's unclear the substantive cause. Is it because that the

---

[1] The details about the City Brain challenge (KDD CUP 2021) can be found from http://www.yunqiacademy.org/poster.
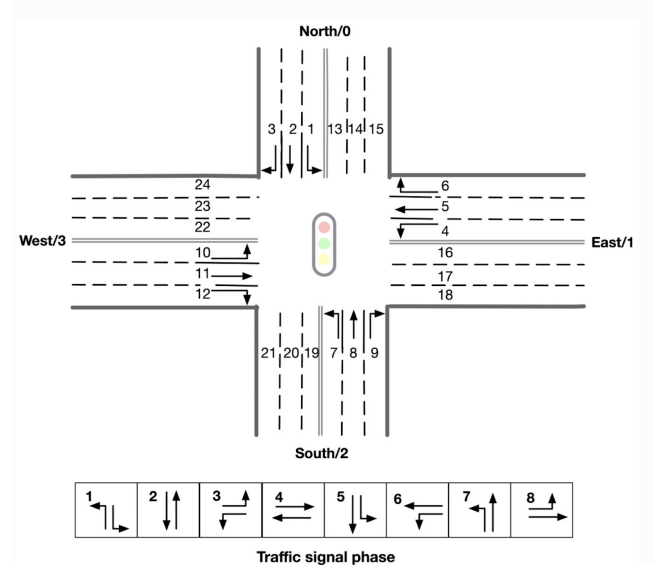
---

**Figure 1: Illustration of a a four-leg intersection and its traffic signal phases (8 phases). When modeling the traffic signal control as a Markov Decision Process (MDP), we consider the intersection as an agent and traffic signal phase as the actions.**

number of vehicles has exceeded the capacity of the city or that we fail to utilize the road network at its maximum capacity? For example, Tokyo and New York City rank similarly by traffic congestion index. However, Tokyo has 43% more registered vehicles than New York City while Tokyo only has 15% more signalized intersections and 32% more road length than New York City [4]. Why is Tokyo able to serve more vehicles than New York City? Is New York City operating the traffic at its maximum capacity? To answer these questions, we are invited to design a policy to coordinate the traffic signals and find the maximum number of vehicles that can be served with a city-scale road network in the City Brain challenge from KDD CUP 2021.

The paper is organized as follows. In section 2, we formally introduce the problem, evaluation metric and challenges in traffic signal control, then formulate the problem as a Markov Decision Process (MDP) and present the proposed method in section 3. We

show the experimental results in section 4 and give the conclusion and future work in section 5.

## 2 PROBLEM AND CHALLENGES

**Problem.** Figure 1 presents a four-leg intersection. In each period of the time step, only one of eight types of signal phases can be selected and a corresponding pair of non-conflict traffic movements will be served. For example, phase-1 gives right-of-way for left-turn traffic from northern and southern approaches. In the real-world road network, there are typically a large number of intersections. The intelligent traffic signal system needs to make a decision from a learned or other manually specific joint policy to coordinate the movements of vehicles at each intersection. We are required to design a joint policy for the traffic system to coordinate the traffic signals and make the traffic system serve more vehicles while maintaining a lower acceptable delay.

**Evaluation Metrics.** The *total number of served vehicles* refers to the total number of vehicles entering the network. The trip delay index is computed as actual travel time divided by travel time at free-flow speed. For an uncompleted trip, the free-flow speed is used to estimate the travel time of the rest of the trip. The delay index is computed as average trip delay index overall vehicles served, *i.e.*, $D = \frac{1}{N} \sum_{i=1}^{M} D_i$, where $M$ refers to the number of vehicles. $D_i$ refers to the delay of $i$-th vehicle, which is defined by

$$D_i = \frac{TT_i + TT_i^r}{TT_i^f} \tag{1}$$

where $TT_i$ refers to travel time of vehicle $i$, $TT_i^r$ refers to the rest of trip travel time which is estimated with free-flow speed, and $TT_i^f$ refers to full trip travel time at free-flow speed.

**Challenges.** Many researchers have addressed this problem via optimization techniques, such as reinforcement learning (RL) [5, 8, 12], however, this problem is still far from solved because of the following challenges:

- There is a lack of efficient reinforcement learning environment which can simulate the movements for hundreds of thousands of vehicles in the large-scale road network [10]. For example, the environment provided by City Brain KDD CUP 2021 is not that efficient. Specifically, the provided environment costs about 3 minutes for one episode (about 300 steps) when running on a standard 8-cores machine. It's difficult for an RL agent to play millions of episodes, which is typically necessary for an RL agent to learn its policy well from scratch [8, 9].
- The action space is exponentially large over a large number of valid intersections (agents). For example, each agent has 8 actions in our scenario. The road map with 1000 intersections will have $8^{1000}$ joint actions, which makes it intractable to learn the optimal policy within the limited time and computation resources.
- Reinforcement learning typically leads to an overfitting policy on the specific environment. However, our policy is evaluated in an unknown environment with limited submissions. It requires us to learn a well-generalized policy.

In the next section, we provide a practical approach to handle the above challenges.

## 3 METHOD

In this section, we formulate the traffic signal control as a Markov Decision Process (MDP) and employ ensemble reinforcement learning to learn the approximate joint policy for large-scale traffic signal control.
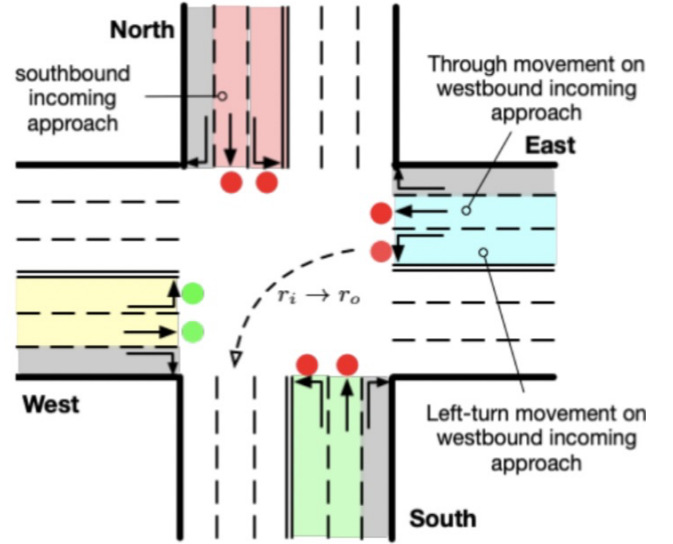
### 3.1 Environment



**Figure 2: Traffic movement for the specific signal phase.**

The environment simulates the vehicle movements and provides the observation for all vehicles, roads and $N$ intersections (agents) [2]. A road consists of a set of lanes. There are two kinds of lanes in our setting: incoming lanes and outgoing lanes (also known as approaching/entering lane and receiving/exiting lane). There are some important parameters in the lanes, including speed limit, length, latitude, longitude and *etc*. An intersection may have a traffic signal or not. For the traffic signal control problem, we only need to consider the intersections that have signals. A movement signal has two states: *green* denotes the corresponding movement is allowed and *red* denotes the movement is prohibited. For a traffic signal, there are at most 8 phases (numbered from 1 to 8). Each phase allows a pair of non-conflict traffic movements for one intersection. Figure 2 illustrates one of the traffic movements for the specific signal phase. The *action* is defined as the traffic signal phase for each intersection to be selected in the next 10 seconds. If an agent is switched to a different phase, there will be a 5 seconds period of *all-red* at the beginning of the next phase, which means all vehicles could not pass this intersection.

After analyzing the environment carefully, we find there are five challenges for this task:

- The traffic flow exceeds the load of road net. It is impossible for the traffic system to deliver all vehicles within the given

---

[2]The source code of environment for round 3 on KDD CUP is available in https://github.com/CityBrainChallenge/KDDCup2021-CityBrainChallenge-starter-kit

period of time. It's necessary to design an intelligent policy to deliver high-value vehicles first.

- There is uncertainty in the environment. The same policy could generate two different delay values. The perturbation is within a certain range [0, 0.001].
- The environment is not efficient to play millions of episodes. We need to develop multi-processing tasks within the limited memory and computation resources.
- It's difficult to design a suitable observation (state) because each intersection may have 0 to 1000 vehicles. Also, the designed state should represent not only vehicles in the current intersection but also its neighbor intersections.
- Each intersection has 8 actions. A road network with $N$ intersections contains $8^N$ different actions. The action space of joint policy will be exponentially large over the number of intersections. It is intractable to develop hand-coded solutions and compute the optimal solutions [2].
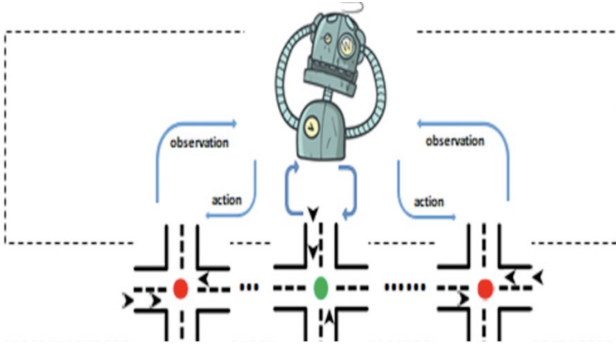
## 3.2 Formulation



**Figure 3: Learning reinforcement learning policy from the formulated MDP.**

The environment simulates the vehicle movements and provides an initial state $s^0 \in \{S_i\}_{i=1}^N$ for all vehicles, roads and $N$ intersections (agents) as illustrated by figure 3. After observing $t$-th state $s^t$, all the agents take actions $a^t \in \{\mathcal{A}_i\}_{i=1}^N$ according to the specific policy $\pi(a^t|s^t)$ and received their rewards $r^t \in \{r_i\}_{i=1}^N$. The environment transforms its state from $s^t$ to $s^{t+1} \in \{S_i\}_{i=1}^N$. Therefore, the whole process will be a standard Markov Decision Process (MDP), which is characterized by a tuple $(\{S_i\}_{i=1}^N, \{\mathcal{A}_i\}_{i=1}^N, \{r_i\}_{i=1}^N, \gamma)$, where $N$ is the number of agents, $S_i$ is the local state of the agent $i$ and $\mathcal{A}_i$ denotes the set of action available to agent $i$, $\gamma$ denotes the decay factor. Typically, $0 \le \gamma \le 1$. Our goal is to learn a joint policy $\pi$ that maximizes the following optimization problem

$$\eta(\pi) = \mathbb{E}_\pi \left[ \sum_{t=0}^\infty \gamma^t \sum_{i=1}^N r_i^t \right] \qquad (2)$$

To make the above optimization problem tractable, we need to define the two important factors in the MDP: the state $\{S_i\}_{i=1}^N$ and the reward $\{r_i\}_{i=1}^N$, which are important in our scenario. In this paper, we employ Deep Q-Network (DQN) [6] to learn the policy $\pi$ via an end-to-end fashion. We introduce the details of the two factors in the next sections.

## 3.3 State with Feature Engineering

We derived more than 200 features in our model. Because of the limitation of space, we only introduce some important features as following.

*Vehicle information.* (1) The number of vehicles that are closer or further to the intersections. (2) According to our observation, it's interesting to notice that most of the roads are longer than 200 units, which indicates that the vehicle positions are very important. We segment the road into several pieces by the movement time of the vehicles moving from their positions to the intersection in free-flow speed, which makes it easy for the agent to determine when to switch to the other phase.

*Road information.* The length of roads. Typically, fixing other environmental variables, the shorter lane will be much busy. When a heavy traffic jam happens in a lane, the other lane in the same road is not available and vehicles are not able to enter the intersection. So we put the road length into features, and specify a larger weight on the vehicle within the shorter road.

*Route prediction.* The route of vehicle matters a lot in the competition while the metric depends on $TT_i^f$ as defined in equation (1). However, we can not observe the whole route in the environment. we predict the whole route for each vehicle according to its history route. This trick reduces the delay performance by about 0.01.

## 3.4 Reward Function

We design a new reward according to the evaluation metric in equation 1, which is defined by

$$D_i^t = \sum_{j=0}^M \left( \frac{TT_j + TT_j^r}{1000} \right) \qquad (3)$$

$$r_i^t = \sum_{j=0}^2 \left( \frac{D_i^{t+j}}{(j * 0.6 + 1)} \right) \qquad (4)$$

where $i$ refers to agent id, $t$ refers to the step number in each episode.

Note that, in equation 4, we sum the rewards of sequential three steps to make the model benefit from long-term reward.

## 3.5 Training

As we have already considered the interactions of neighbouring agents in the state (see section 3.3), we employ independent deep Q-Learning (DQN) to optimize the problem in equation (2). However the environment is not efficient to play millions of episodes when the number of vehicles is larger than 10000, we develop a multi-processing learning task. Note that, DQN is an off-policy algorithm and makes it suitable for the multi-processing task [7]. Specifically, in the competition, we train each model with 4 ∼ 8 process by 24 ∼ 48 hours. Some other articles [3, 11] also provide the task-specific rewards, however, they don't meet the requirement in our scenario.
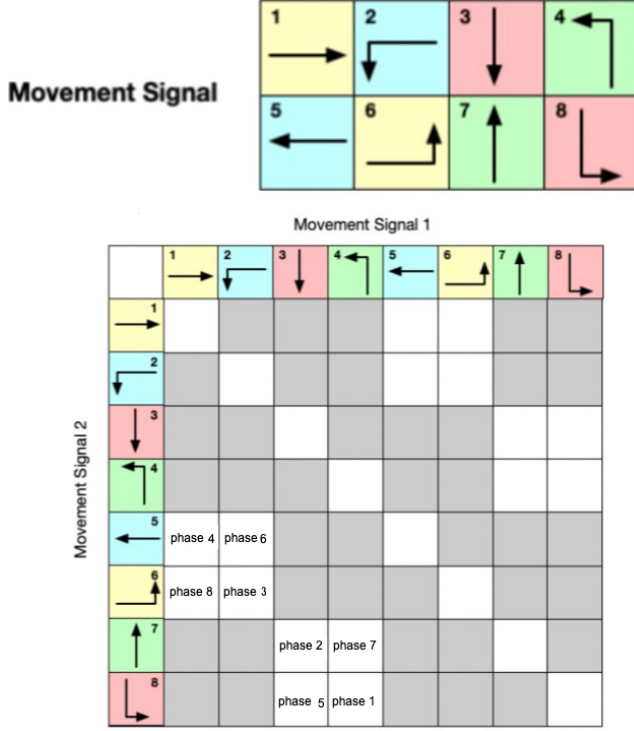
**Figure 4: Signal and Phase.**

## 3.6 Ensemble Learning

It's interesting to point out that the RL model trained on one speci-
fied traffic flow typically overfits the flow and performs worse on
the private flow provided by the leaderboard. To make our policy
generalize well, we use different features combinations to train 4
models. Because the neural network has a big uncertainty in re-
gression task [1], we average their predictions. Then we develop
a rule-based system by evaluating the value of action with differ-
ent trials. The final action is an ensemble based on the weighting
average of rule-based system score ranking and neural networks'
average score.

## 4 EXPERIMENT RESULTS

Since the leaderboard has changed a lot before the deadline, we
do not have historical experimental data. Table 1 shows the top 10
team submissions.

## 5 CONCLUSION AND FUTURE WORK

In this paper, we propose a practical ensemble reinforcement-learning
method to handle the problem of multi-intersection traffic signal
control, especially for large-scale networks. Our proposed method
has shown its strong performance and generalization ability. In the
future, one can use multi-agent reinforcement learning to capture
the neighbors' information for each intersection.

**Table 1: Experimental Results of different teams on Round 3. The larger of total served vehicles will be better. For the same total served vehicles, the lower delay will be better. GoodGoodStudy refers to our method.**

| Team Id | Total Served Vehicle | Delay |
|---|---|---|
| IntelligentLight | 360637 | 1.4027135308052847 |
| **GoodGoodStudy** | **358895** | **1.4052081109433103** |
| 4PQC_team | 358888 | 1.4025206513613906 |
| SmartLight | 357229 | 1.4027102867089065 |
| BOE_IOT_AIBD | 356463 | 1.4013649198707516 |
| SUMO | 335599 | 1.401792220782304 |
| IF_BigData | 335376 | 1.4020449312353032 |
| alphabeta | 318720 | 1.4034693515888321 |
| bingo | 317954 | 1.4010931499550239 |
| D | 316941 | 1.4021008708625378 |

## REFERENCES

[1] Filipe Aires, Catherine Prigent, and William B Rossow. 2004. Neural network uncertainty assessment using Bayesian statistics: A remote sensing application. *Neural computation* 16, 11 (2004), 2415–2458.

[2] Bram Bakker, Shimon Whiteson, Leon Kester, and Frans CA Groen. 2010. Traf-fic light control by multiagent reinforcement learning systems. In *Interactive Collaborative Information Systems.* Springer, 475–510.

[3] Chacha Chen, Hua Wei, Nan Xu, Guanjie Zheng, Ming Yang, Yuanhao Xiong, Kai Xu, and Zhenhui Li. 2020. Toward a thousand lights: Decentralized deep reinforcement learning for large-scale traffic signal control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 3414–3421.

[4] KDD CUP. 2021. City Brain Challenge. http://www.yunqiacademy.org/poster

[5] Mengyu Guo, Pin Wang, Ching-Yao Chan, and Sid Askary. 2019. A reinforcement learning approach for intelligent traffic signal control at urban intersections. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC).* IEEE, 4242–4247.

[6] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *nature* 518, 7540 (2015), 529–533.

[7] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc G Bellemare. 2016. Safe and efficient off-policy reinforcement learning. *arXiv preprint arXiv:1606.02647* (2016).

[8] Chao Qu, Hui Li, Chang Liu, Junwu Xiong, James Zhang, Wei Chu, Weiqiang Wang, Yuan Qi, and Le Song. 2020. Intention propagation for multi-agent rein-forcement learning. *arXiv preprint arXiv:2004.08883* (2020).

[9] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, and Julian Schrittwieser et al. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 7587 (2016).

[10] Zheng Tang, Milind Naphade, Ming-Yu Liu, Xiaodong Yang, Stan Birchfield, Shuo Wang, Ratnesh Kumar, David Anastasiu, and Jenq-Neng Hwang. 2019. Cityflow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 8797–8806.

[11] Hua Wei, Chacha Chen, Guanjie Zheng, Kan Wu, Vikash Gayah, Kai Xu, and Zhenhui Li. 2019. Presslight: Learning max pressure control to coordinate traffic signals in arterial network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 1290–1298.

[12] Hua Wei, Nan Xu, Huichu Zhang, Guanjie Zheng, Xinshi Zang, Chacha Chen, Weinan Zhang, Yanmin Zhu, Kai Xu, and Zhenhui Li. 2019. Colight: Learning network-level cooperation for traffic signal control. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management.* 1913–1922.