

Download full-text PDF

Read full-text

Download citation

Copy link

▼


Article

PDF Available


Traffic Light Control Using Deep Policy-Gradient and Value-Function Based Reinforcement Learning

April 2017 · [IET Intelligent Transport Systems](#) 11(7)  
DOI:[10.1049/iet-its.2017.0153](#)

Authors:




**Sajad Mousavi**  
Harvard Medical School



**Michael Schukat**  
National University of Ireland, Galway



**Peter Corcoran**



**Enda Howley**  
National University of Ireland, Galway

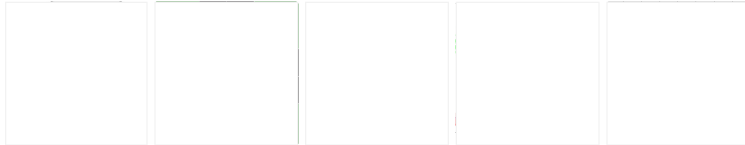
Citations (133)

References (53)

Figures (5)

Abstract and Figures

Recent advances in combining deep neural network architectures with reinforcement learning techniques have shown promising potential results in solving complex control problems with high dimensional state and action spaces. Inspired by these successes, in this paper, we build two kinds of reinforcement learning algorithms: deep policy-gradient and value-function based agents which can predict the best possible traffic signal for a traffic intersection. At each time step, these adaptive traffic light control agents receive a snapshot of the current state of a graphical traffic simulator and produce control signals. The policy-gradient based agent maps its observation directly to the control signal, however the value-function based agent first estimates values for all legal control signals. The agent then selects the optimal control action with the highest value. Our methods show promising results in a traffic network simulated in the SUMO traffic simulator, without suffering from instability issues during the training process.



Deep reinforcement... The intersection geometry for th... A comparison of performance of... Average Cumulative del... Average queue length of the...

Figures - uploaded by [Sajad Mousavi](#) Author content  
Content may be subject to copyright.

Discover the world's research

- 20+ million members
  - 135+ million publications
  - 700k+ research projects
- Join for free

Advertisement

[Download full-text PDF](#)[Read full-text](#)[Download citation](#)[Copy link](#)

Content uploaded by [Sajad Mousavi](#) Author content  
Content may be subject to copyright.

# Traffic Light Control Using Deep Policy-Gradient and Value-Function Based Reinforcement Learning

Seyed Sajad Mousavi  
Discipline of IT  
National University of Ireland,  
Galway  
s.mousavi1@nuigalway.ie

Michael Schukat  
Discipline of IT  
National University of Ireland,  
Galway  
michael.schukat@nuigalway.ie

Enda Howley  
Discipline of IT  
National University of Ireland,  
Galway  
enda.howley@nuigalway.ie

## ABSTRACT

Recent advances in combining deep neural network architectures with reinforcement learning techniques have shown promising potential results in solving complex control problems with high dimensional state and action spaces. Inspired by these successes, in this paper, we build two kinds of reinforcement learning algorithms: deep policy-gradient and value-function based agents which can predict the best possible traffic signal for a traffic intersection. At each time step, these adaptive traffic light control agents receive a snapshot of the current state of a graphical traffic simulator and produce control signals. The policy-gradient based agent maps its observation directly to the control signal, however the value-function based agent first estimates values for all legal control signals. The agent then selects the optimal control action with the highest value. Our methods show promising results in a traffic network simulated in the SUMO traffic simulator, without suffering from instability issues during the training process.

## CCS Concepts

•Theory of computation → Sequential decision making;

## Keywords

Traffic control, Reinforcement learning, Deep learning, Policy gradient method, Value-function method, Artificial neural networks

## 1. INTRODUCTION

With regard to fast growing population around the world, the urban population in the 21<sup>st</sup> century is expected to increase dramatically. Hence, it is imperative that urban infrastructure is managed effectively to contend with this growth. One of the most critical consideration when designing modern cities is developing smart traffic management systems. The main goal of a traffic management system is reducing traffic congestion which nowadays is one of the major issues of megacities. Efficient urban traffic management results in time and financial savings as well as reducing CO<sub>2</sub> emission into atmosphere. To address this issue, a lot of solutions have been proposed [23, 4, 1, 22]. They can be roughly classified into three groups. The first is pre-timed signal control, where a fixed time is determined for all green phases according to historical traffic demand, without considering possible fluctuations in traffic demand. The second

is vehicle-actuated signal control where, traffic information is used, provided by inductive loop at an equipped intersection to decide to control e.g. extending or terminating a green phase. adaptive signal control, where the signal timing is managed and updated automatically according to the current state of the intersection (i.e. traffic density, length of vehicles in each lane of the intersection, traffic flow fluctuation) [13]. In this study, we aim to propose two methods for traffic signal control by leveraging recent advances in machine learning and artificial intelligence.

Reinforcement learning [34] as a machine learning technique for traffic signal control problem has led to promising results [4, 30] and has shown a promising potential. It does not need to have a perfect knowledge of the environment in advance, for example traffic flow. Agents are able to gain knowledge and model the dynamic environment just by interacting with it. A reinforcement learning agent learns based on trial and error. It receives a scalar reward after taking each action in the environment. The obtained reward is based on how well the action and the agent's goal is to learn an optimal policy. The discounted cumulative reward is maximized through interaction with its environment. Aside from traffic signal control, reinforcement learning has been applied to a number of world problems such as cloud computing [12, 13].

Typically the complexity of using reinforcement learning in real world applications such as traffic signal control grows exponentially as state and action spaces increase. To deal with this problem, function approximation and hierarchical reinforcement learning approaches are used. Recently, deep learning has gained huge attention and has been successfully combined with reinforcement learning techniques to deal with complex optimization problems such as playing Atari 2600 games [27], Computer Go [33], etc., where the classical RL methods could not find optimal solutions. In this way, the current state of the environment is fed into a deep neural net (e.g. a convolutional neural network [20]) trained by reinforcement learning techniques to predict the next possible optimal action.

Inspired by the successes of combining reinforcement learning with deep learning paradigm and with regard to the complex nature of environment of traffic signal control, in this paper we aim to use the effectiveness of deep reinforcement learning to build adaptive signal control methods in order to optimize the traffic flow. In few previous studies have tried to apply deep reinforcement learning to traffic signal control.

arXiv:1704.08883v2 [cs.LG] 27 May 2017

[Download full-text PDF](#)[Read full-text](#)[Download citation](#)[Copy link](#)

learning in the traffic signal control problem [38, 14], in this research the state representation is different. Also, One of our methods uses policy gradient method which does not suffer from oscillations and instabilities during training process and can take full advantage of the available data of the environment to develop the optimal control policy.

We propose adaptive signal controllers by combination two reinforcement learning approaches (i.e. policy gradient and action-value function) and a deep convolution neural network, which perceive embedded camera observations in order to produce control signals in an isolated intersection. We conduct simulated experiments with our proposed methods in SUMO traffic simulator.

The rest of the paper is organized as follows. Section 2 provides related work in the area of traffic light control (TLC). Section 3 gives a brief review of reinforcement learning techniques which we have used in this research. Section 4 presents how to formulate the TLC problem as a reinforcement learning task and the proposed methods to solve the task. Then Section 5 provides simulation results and the performance of the proposed approaches. Finally Section 6 concludes the paper and give some directions for future research.

## 2. RELATED WORK

A lot of research has been done in academic and industry communities to build adaptive traffic signal control systems. In particular, significant research has been conducted employing reinforcement learning methods in the area of traffic light signal control [39, 2, 7]. These works have achieved promising results. However, their simulation testbeds have not been mature enough to be comparable with more realistic situations. Developing advance traffic simulation tools have made researchers develop novel state representation and reward functions for reinforcement learning algorithms, which could consider more aspects of complexity and reality of real-world traffic problems [13, 1, 8, 3]. All this these attempts viewed the traffic light control problem as a fully observable Markov decision process (MDP) and investigated whether Q-learning algorithm can be applied to it. However, Richter's study formulated the traffic problem as a partially observable MDP (POMDP) and applied policy gradient methods to guarantee local convergence under a partial observable environment [31].

By utilizing advances in deep learning and its application to different domains [10, 11], deep learning has gained attention in the area of traffic management systems. Previous research has used deep stacked autoencoders (SAE) neural networks to estimate Q-values, where each Q-value is corresponding to each available signal phase [21]. It considered measures of speed and queueing length as its state in each time step of learning process of its proposed method. Two recent studies by [38, 14] provided deep reinforcement learning agents that used deep Q-network [27] to map from given states to Q-values. Their state representations were a binary matrix of the positions of vehicles on the lanes of an intersection, and a combination of the presence matrix of vehicles, speed and the current traffic signal phase, respectively. However, we use raw visual input data of the traffic simulator snapshots as system states. Moreover, in addition to estimating Q-function, one of the proposed methods directly maps from the input state to a probability distribution over actions (i.e. signal phases) via deep policy gradient

method.

## 3. BACKGROUND

In this section, we will review Reinforcement (RL) approaches and briefly describe how RL real world problems where the number of state are extremely high so that the regular reinforcement techniques cannot deal with them.

### 3.1 Reinforcement Learning

A common reinforcement learning [34] setting Figure 1 where an RL agent interacts with an environment. The interaction is continued until reaching a terminal state or the agent meets a termination condition. In this paper, we consider problems that RL techniques are applied to, a Markov decision processes (MDPs). A MDP is a five-tuple  $\langle S, A, T, R, \gamma \rangle$  where  $S$  is the set of the state space of the environment,  $A$  is the set of the action space that the agent can use in order to interact with the environment,  $T$  is the transition function that gives the probability of moving between the environment states to the next states  $s_{t+1} \in S$  given the current state  $s_t \in S$  and action  $a_t \in A$ .  $R$  is the reward function and  $\gamma \in [0, 1]$  is known as the discount factor, which models the importance of the immediate rewards. At each time step  $t$ , the agent observes the state  $s_t \in S$  and, based on its observation, chooses an action  $a_t \in A$ . Taking the action, leads to the state transitions to the next states  $s_{t+1} \in S$  given the current state  $s_t \in S$  and action  $a_t \in A$ . Then, the agent receives a reward  $r_t$  which is determined by the reward function.

The goal of the learning agent in reinforcement learning framework is to learn an optimal policy  $\pi : S \rightarrow A$  which defines the probability of selecting action  $a_t$  given the state  $s_t$ , so that with following the underlying policy, the cumulative discounted reward over time is maximized. The discounted future reward,  $R_t$  at time  $t$  is defined as:

$$R_t = E[\sum_{k=0}^{\infty} \gamma^k r_{t+k}],$$

where the role of the discount factor  $\gamma$  is to weigh the worth of immediate and future rewards. In many problems, there are many states and actions so that it is impossible to apply classic reinforcement learning techniques, which consider tabular representations for states and action spaces. For example, in the traffic light optimization, that we interest in this paper, the state space is continuous. Hence, it is common to use function approximators [42] or decomposition and approximation techniques like Hierarchical Reinforcement Learning approaches [13, 15, 28] and advance HRL [16].

Different forms of function approximators are used with reinforcement learning techniques. For linear function approximation, a linear combination of state and action features  $f$  and learned weights  $w$  (i.e.  $\sum_i f_i w_i$ ) or a non-linear function approximation (e.g. neural network). Until recently, the majority of reinforcement learning has been applying linear function approximators. More recently, deep neural networks such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), stacked auto-encoders (SAE) also been commonly used as function approximators for reinforcement learning tasks [19, 26]. The interested reader is referred to [29] for a review of using deep neural networks for reinforcement learning.

[Download full-text PDF](#)[Read full-text](#)[Download citation](#)[Copy link](#)

**Figure 1: Deep reinforcement learning agent of traffic signal control.**

with reinforcement learning framework.

### 3.2 Deep learning and Deep Q-learning

Deep learning techniques are one of the best solutions to address high dimensional data and extract discriminative information from the data. Deep learning algorithms have the capability of automating feature extraction (the extraction of representations) from the data. The representation are learnt through the data which are fed directly into deep nets without using human knowledge (i.e. automated feature extraction). Deep learning models contain multiple layers of representations. Indeed, it is a stack of building blocks such as auto-encoders, Restricted Boltzmann Machines (RBMs) and convolutional layers. During training, the raw data is fed into a network consisting of multiple layers. The output of the each layer which is nonlinear feature transformations, is used as inputs to the next layers of the deep neural network. The output representation of the final layer can be used for constricting classifiers or those applications which can have the better efficiency and performance with abstract representation of the data in a hierarchical manner as inputs. A nonlinear transformation is applied at each layer on its input to try to learn and extract underlying explanatory factors. Consequently, this process learns a hierarchy of abstract representations.

One of the main advantages of deep neural networks is the capability of automating feature extraction from raw input data. A deep Q-learning Network (DQN) [26] uses this benefit of deep learning in order to represent the agent's observation as an abstract representation in learning an optimal control policy. The DQN method aggregates a deep neural network function approximator with Q-learning to learn action value function and as a result a policy  $\pi$ , the behaviour of the agent which tells the agent what action should be selected for each input state. Applying non-linear function approximators such as neural networks with model-free reinforcement learning algorithms in high-dimensional continuous state and action spaces, has some convergence problems [37]. The reasons for these issues are: 1) Consecutive states in reinforcement learning tasks have correlation. 2) The underlying policy of the agent is changing frequently, because of slight changes in Q-values. To cope with these problems, the DQN provides some solutions which improve the performance of the algorithm significantly. For the prob-

lem of correlated states, DQN uses the previous experience replay approach [24]. In this way, step, the DQN stores the agent's experience ( $s$ ) into a data set  $D$ , where  $s_t$ ,  $a_t$ , and  $r_t$  are then action and received reward, respectively at state at the next time step. To update the DQN utilizes stochastic minibatch updates with random sampling from the experience replay (previous observed transitions) at training time. strong correlations between consecutive samples approach to deal with aforementioned convergence which we also examine in this research, is the current methods. This approach has demonstrated convergence properties in some RL problems [35].

### 3.3 Policy Gradient Methods

A Policy Gradient (PG) method tries to optimize a parameterized policy function by gradient descent. Indeed, policy gradient methods are interested in policy space to learn policies directly, instead of learning state-value or action-value functions. Unlike other reinforcement learning algorithms, PG does not suffer from the convergence problems of estimating functions under nonlinear function approximation environments which might be partially observed. They can also deal with the complexity of continuous and action spaces better than purely value-based methods [35]. Policy gradient methods estimate policy using Monte Carlo estimates of the policy gradient. These methods are guaranteed to converge to a local optimum of their parametrized policy function. However, these methods result in high variance in their gradient estimates. Hence, in order to reduce the variance of the gradient estimators, some methods subtract a baseline function from policy gradients. The baseline function can be estimated in different manners [32, 40]. By inspiring these methods and successes of neural networks in feature abstractions, we use deep neural networks to learn an optimal traffic control policy directly in the control problem.

## 4. SYSTEM DESCRIPTION

In this section, we will formulate traffic light control problem as a reinforcement learning task by describing the actions and reward function. We then present a deep neural network and how to train the network.

### 4.1 State Representation

We represent the state of the system as an image or a snapshot of the current state of a graph (e.g. SUMO-GUI [18]) which is a vector of row of current view of the intersection at each station (as shown in Figure 1). This kind of representation is like putting a camera on an intersection with it to view the whole intersection. The state representation in the traffic light control literature usually represents the presence of a vehicle at the intersection as a Boolean-valued vector where a value 1 indicates the presence of a vehicle and a value 0 indicates the absence [38, 36], or a combination of the presence and another vector indicating the vehicle's speed at the intersection [14]. Regardless of these state representations that are using a prior knowledge provided, the

Download full-text PDF

Read full-text

Download citation

Copy link

assumptions which are not generalizable for the real world. For instance, they discretize a lane segment of an intersection into cells with a constant length  $c$  which is supposed to be the vehicle length to build the vehicle's speed and presence vectors. However, by feeding the state as an image to a convolutional neural network, the system can detect the location and presence of all vehicles with different lengths and as result the vehicles' queue on each lane. Furthermore, by stacking a history of consecutive observations as input, the convolutional layers of a deep network are able to estimate velocity and travel direction of vehicles. Hence, the system can implicitly benefit from these information as well.

## 4.2 Action Set

To control traffic signal phases, we define a set of possible actions  $A = \{\text{North/South Green (NSG), East/West Green (EWG)}\}$ . NSG allows vehicles to pass from North to South and vice versa, and also indicates the vehicles on East/West route should stop and not proceed through the intersection. EWG allows vehicles to pass from East to West and vice versa, and implies the vehicles on North/South route should stop and not proceed through the intersection. At each time step  $t$ , an agent regarding its strategy chooses an action  $a_t \in A$ . Depending the selected action, the vehicles on each lane are allowed to cross the intersection.

## 4.3 Reward Function

Typically an immediate reward  $r_t \in \mathbb{R}$  is a scalar value which the agent receives after taking the chosen action in the environment at each time step. We set the reward as the difference between the total cumulative delays of two consecutive actions, i.e.

$$r_t = D_{t-1} - D_t, \quad (2)$$

where  $D_t$  and  $D_{t-1}$  are the total cumulative delays in the current and previous time steps. The total cumulative delay at time  $t$ , is the summation of the cumulative delay of all the vehicles appeared from  $t = 0$  to current time step  $t$  in the system. The positive reward values imply the taken actions led to decrease the total cumulative delay and the negative rewards imply an increase in the delay. With regard to the reward values, the agent may decide to change its policy in certain states of the system in the future.

## 4.4 Agent's Policy

The agent chooses the actions based on a policy  $\pi$ . In the policy-based algorithm, the policy is defined as a mapping from the input state to a probability distribution over actions  $A$ . We use the deep neural network as the function approximator and refer its parameters  $\theta$  as policy parameters. The policy distribution  $\pi(a_t|s_t; \theta)$  is learned by performing gradient descent on the policy parameters. In the value-function based algorithm, the deep neural network is utilized to estimate the action-value function. The action-value function maps the input state to action values, which each represents the future reward that can be achieved for the given state and action. The optimal policy can then be extracted by performing a greedy approach to select the best possible action.

## 4.5 Objective Function and System Training

There are many measures such as maximizing throughput, minimizing and balancing queue length, minimizing the

**Algorithm 1** Deep Value-Function based reinforcement agent of traffic signal control with replay

```

1: Initialize parameters,  $\theta$  with random value
2: Initialize replay memory  $M$  with capacity
3: for each simulation do
4:   initialize  $s$  with current view of the intersection
5:   repeat # each step in the simulation
6:     choose action  $a$  according to  $\epsilon$ -greedy
7:     take action  $a$ , observe reward  $r$  and next state  $s'$ 
8:     store transition  $(s, a, r, s')$  in  $M$ 
9:      $s \leftarrow s'$ 
10:     $b \leftarrow$  sample random minibatch of transitions from the replay memory,  $M$ 
11:    for each transition  $(s_j, a_j, r_j, s'_j)$  in  $b$  do
12:      if  $s'_j$  is terminal then
13:         $y_j \leftarrow r_j$ 
14:      else
15:         $y_j = r_j + \gamma \max_{a'} Q(s'_j, a'; \theta^-_{i-1})$ 
16:      end if
17:      update parameters  $\theta$  according to  $y_j$ 
18:    end for
19:  until  $s$  is terminal
20: end for
```

delay, etc. in the traffic signal management research, the agent aims to maximize the total cumulative delay, which empirically has to maximize throughput and to reduce queue details discussed in Section 5.3).

The objective of agent is to maximize the expected cumulative discounted reward. We aim to maximize under the probability distribution  $\pi(a_t|s_t; \theta)$ :

$$J(\theta) = E_{\pi_\theta} \left[ \sum_{t=0}^T \gamma^t r_t \right] = E_{\pi_\theta} [R].$$

We divide the system training based two RL Value-function based and Policy-based. in **value-based approach**, the value function,  $Q_\pi(s, a)$  follows:

$$Q_\pi(s, a) = E_\pi [r_t + \gamma \max_{a'} Q(s', a') | s]$$

Where it is implicit that  $s, s' \in S$  and  $a \in A$ . The function can be parameterized,  $Q(s, a; \theta)$  with vector  $\theta$ . Typically, the gradient-descent method to learn parameters,  $\theta$  by trying to minimize the loss function of mean-squared error in  $Q$  value

$$J(\theta) = E_\pi [(r + \gamma \max_{a'} Q(s', a'; \theta) - Q(s, a; \theta))^2]$$

Where  $r + \gamma \max_{a'} Q(s', a'; \theta)$  is the target value. In the DQN algorithm, a target Q-network is used to solve the instability problem of the policy. The network with the target Q-network to obtain consistent targets by keeping the weight parameters ( $\theta^-$ ) and learning target fixed and updating them periodically. The target value of the DQN is represented as follows

$$y_i = r + \gamma \max_{a'} Q(s', a'; \theta^-_{i-1})$$



Download full-text PDF

Read full-text

Download citation

Copy link

**Figure 2: The intersection geometry for the traffic simulation.**

Where  $\theta^-$  is parameters of the target network. The stochastic gradient descent method is used in order to optimize equation (5). The parameters of the deep Q-learning algorithm are updated as follows:

$$\theta_i \leftarrow \theta_{i-1} + \alpha(y_i - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i) \quad (7)$$

Where  $y_i$  is the target value for iteration  $i$  and  $\alpha$  is a scalar learning rate. Algorithm 4.4 presents the pseudo-code for the training algorithm.

In **policy-based approach**, The gradient of the objective function represented in equation (3) is given by:

$$\nabla_{\theta} J = \sum_{t=0}^T E_{\pi_{\theta}} [\nabla_{\theta} \log(a_t | s_t; \theta) R_t]. \quad (8)$$

This equation (8) is standard learning rule of the REINFORCE algorithm [41]. It updates the policy parameters  $\theta$  in the direction  $\nabla_{\theta} \log(a_t | s_t; \theta)$  so that the probability of action  $a_t$  at state  $s_t$  is increased if it has led to high cumulative reward, however it is decreased if the action has result in a low reward. The gradient estimate in equation 2 results to have high variance. It is common to reduce the variance by subtracting a baseline function  $b_t(s_t)$  from the return  $R_t$ , without changing expectation. Commonly an estimate of the state value function is used as the baseline,  $b_t(s_t) = V^{\pi_{\theta_v}}(s_t)$ . Thus, the adjusted gradient is  $\nabla_{\theta} \log(a_t | s_t; \theta) (R_t - b_t(s_t))$ . The value  $R_t - b_t$  is known as the *advantage function*.

With regard to the advantage actor-critic method [25], computing a single update is done by selecting actions using the underlying policy for up to  $M$  steps or till a terminal state is met. In this way, the agent obtains up to  $M$  rewards from the environment at each update point and updates the policy parameters after every  $n \leq M$  steps regarding  $n$ -step returns. The vector parameters  $\theta$  is updated through the stochastic gradient descent method:

$$\theta \leftarrow \theta + \alpha \sum_t \nabla_{\theta} \log(a_t | s_t; \theta) A(s_t, a_t; \theta, \theta_v), \quad (9)$$

where  $A(s_t, a_t; \theta, \theta_v)$  is an estimate of the *advantage function* corresponding  $\sum_{i=0}^{n-1} \gamma^i r_{t+i} + \gamma^n V(s_{t+n}; \theta) - V(s_t; \theta_v)$ , where  $n$  might have different values with respect to the state, up to  $M$ . this process is an actor-critic algorithm, the policy

$\pi(a_t | s_t; \theta)$  refers to the actor and the estimate value function  $V^{\pi_{\theta_v}}(s_t)$  implies to the critic [rithm 4.5 shows the pseudo-code for the traini

**Algorithm 2** Deep Policy-Gradient based r learning agent of traffic signal control

```

1: Initialize parameters,  $\theta, \theta_v$  with random va
2: Initialize step counter  $t \leftarrow 0$ 
3: for each simulation do
4:   initialize  $s$  with current view of the inter
5:    $t_{start} = t$ 
6:   repeat
7:     perform action  $a$  according to policy  $\pi$ 
8:     observe reward  $r$  and next state  $s'$ 
9:      $t \leftarrow t + 1$ 
10:  until  $s$  is terminal or  $t - t_{start} == M$  (
11:  if  $s$  is terminal then
12:     $R = 0$ 
13:  else
14:     $R = V(s; \theta_v)$ 
15:  end if
16:  for  $i \in \{t - 1, \dots, t_{start}\}$  do
17:     $n \leftarrow M$  times
18:     $R \leftarrow r_i + \gamma R$ 
19:     $\theta \leftarrow \theta + \alpha \nabla_{\theta} \log(a_i | s_i; \theta) (R - V(s_i; \theta_v))$ 
20:     $\theta_v \leftarrow \theta_v + \frac{\partial (R - V(s_i; \theta_v))}{\partial \theta_v}$ 
21:  end for
```

## 5. EXPERIMENT AND RESULTS

In this section, we present the simulation where our experiments have been done. We t the details of the deep neural network utilis hyper-parameters to represent the agent's poli

### 5.1 Experiment Setup

We have used the Simulation of Urban MOBi [18] tool to simulate traffic in all experiment a well-known open source traffic simulator wl useful Application Programming Interfaces ( Graphical User Interface (GUI) view to mod networks as well as some possibilities to handle t ticular, we utilised SUMO-GUI v0.28.0. as it a snapshots of each step of the simulation. The geometry used in this study is shown in Figure 4 incoming lanes to the intersection and four o from the intersection. To generate traffic dema ferent directions (i.e. north-to-south and wes vice versa) to the road network, a uniform pr tribution with the probability 0.1 was used.

### 5.2 System Architecture and Hyper

We took the snapshots from the SUMO-GUI basic pre-processing. The snapshots are conver green-blue (RGB) representation to gray-scale them to  $128 \times 128$  frames. To enable our syst rize a history of the past observations, we sta four frames of the history and provided them t as input. So, the input to the network was a 1 image. We applied approximately the same as the Deep Q-Network (DQN) algorithm introdu

Table 1: Comparison of performance of the proposed methods against the SNN method

Evaluation Metric ( $\mu, \sigma; n = 100$ )	Policy-based	Value-function based	SNN
Queue (Vehicles)	(1.79, 0.073)	(1.74, 0.10)	(5.55, 0.73)
Cumulative Delay (s)	(11.25, 0.39)	(11.01, 0.69)	(41.40, 7.31)
Reward (r)	(6.14, 0.29)	(7.14, 0.44)	(1.73, 0.62)

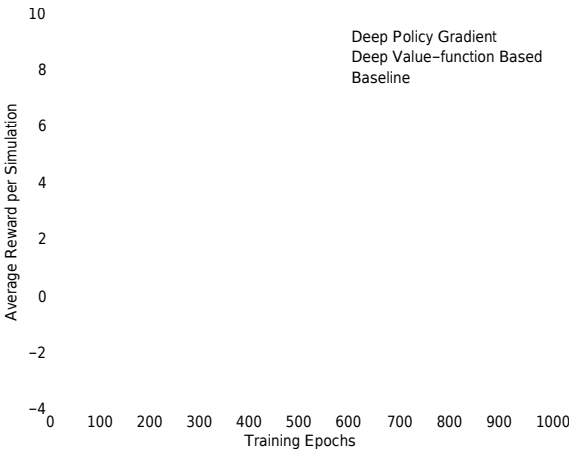


Figure 3: A comparison of performance of the average reward received during the evaluation time for the proposed method and the baseline.

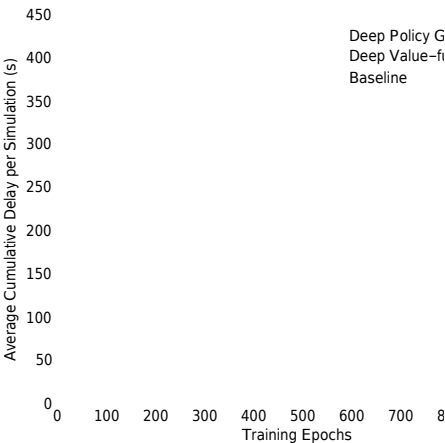


Figure 4: Average Cumulative delay per ing the suggested model and the baseline evaluation time.

et al. [26, 27]. The network consists a stack of two convolutional layers with filters  $16 \times 8$  and  $32 \times 4$  and strides 4 and 2, respectively. The final hidden layer is fully-connected with 256 hidden node. All three hidden layers are followed by a rectifier nonlinearity. The main difference with the network architecture of the DQN method is the last layer, where the last layer of DQN is a fully-connected linear layer with a number of output neurons (i.e. Q-values  $Q(a, s)$ ) corresponding to each action in a given Atari 2600 game, while in policy-based model the last layer represents two set of outputs, a softmax output resulting in a probability distribution over the actions A (i.e. the policy  $\pi(a, s)$ ), and a single linear output node resulting in the estimate of the state value function  $V(s)$ . For value-function model we used the architecture, the same as the DQN. The output layer is corresponding to action values. In all of our experiments, the discount factor was set to  $\gamma = 0.99$  and all weights of the network were updated by the Adam optimizer [17] with a learning rate  $\alpha = 0.00001$  and with mini batches of size  $M$  (up to 32), the maximum number of steps that the agent can take to follow its policy and afterwards needs to update it. The network was trained for about 1050 epoch, approximately 2 million time steps. Each epoch is corresponded 10 episodes and each episode was a complete SUMO-GUI simulation. The learned policies by the agent was evaluated every 10 episodes by running SUMO-GUI for 5 episodes and averaging the resulting rewards, total cumulative delay and queue length.

To evaluate our proposed method we also built a shallow neural network (SNN) with one hidden layer. The hidden layer has 64 hidden nodes followed by a rectifier nonlinearity. The output layer is a fully-connected linear layer with

a number of output neurons corresponding to signal phase in the intersection. Two vectors input state of the network. The first represent ber of queued vehicles at the lanes of the inte North, South, East and West) and the second c to the current traffic signal phase of the inters is trained with the same hyper-parameters and method (i.e. the gradient decent algorithm) as methods.

5.3 Results and Discussion

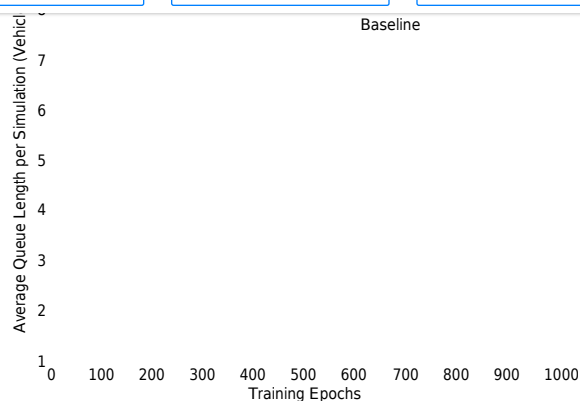
To evaluate the performance of the proposed compared them against a baseline traffic cont troller that gives an equal fixed time to each intersection. We ran SUMO-GUI simulator for model using the configuration setting explaine 5.2 and compared the average reward, average lative delay and average queue length achieved line. Figure 3 shows the received average rew: agent follows a certain policy. As shown in Figu posed method performs significantly better thar and results more reward magnitudes by doing This gradually increasing reward reflects the ity to learn an optimal control policy in a st: Unlike using deep reinforcement learning for es Q-values in traffic light optimisation problem | posed agent doesn't suffer stability issues. In o the learned policy by the agent, two of the n performance metrics in the traffic signal contro implemented: the cumulative delay and queue ures 4 and 5 illustrate the performance comp learing agent regarding average cumulative de

Download full-text PDF

Read full-text

Download citation

Copy link



**Figure 5: Average queue length of the intersection using the proposed model and the baseline during the evaluation time.**

average queue length metrics, respectively, to the baseline, while the agent is following the learning policy over time. The plots clearly show the agent is able to find a policy resulting in minimizing queue length and total cumulative delay. Moreover, these graphs reveal that by using the reward function for reducing cumulative delay, the intersection queue length is reduced as well as the total cumulative delay of all vehicles.

We also compared the proposed methods with the SNN, which is a shallow neural network with one hidden layer. Table 1 reports a comparison of the proposed models and the SNN model in terms of the average and standard deviation ( $\mu$ ,  $\sigma$ ) of average queue length, average cumulative delay time and the received average reward metrics. The results on Table 1 are calculated from the last 100 training epochs of each method. Comparing the metrics shown in Table 1, demonstrates that the proposed models significantly outperform the SNN method. Based on the data in Table 1 we can induce 67% and 72% reductions in the average cumulative delay and queue length for the policy gradient method and 68% and 73% reductions for value-function-based method compared to the SNN. Furthermore, we can see that the proposed methods have received average rewards superior to the SNN. Considering these results, it is obvious that the policy gradient and value-function agents could learn the control policies better than the SNN approach.

## 6. CONCLUSION

In this paper, we applied deep reinforcement learning algorithms with focusing on both policy and value-function based methods to traffic signal control problem in order to find optimal control policies of signalling, just by using raw visual input data of the traffic simulator snapshots. Our approaches have led to promising results and showed they could find more stable control policies compared to previous work of using deep reinforcement learning in traffic light optimization. In our work, we developed and tested the proposed methods in a small application, extending the work for more complex traffic simulations, for instance consider-

coordination problem between agents would be for future research.

## Acknowledgments

We would like to thank FotoNation company to work with their GPU cluster.

## REFERENCES

- [1] M. Abdoos, N. Mozayani, and A. L. Bazzaz. Multi-agent system for traffic signals control. *Engineering Applications of Artificial Intelligence*, 26(5):1575–1587, 2013.
- [2] B. Abdulhai, R. Pringle, and G. J. Karak. Reinforcement learning for true adaptive traffic control. *Journal of Transportation Engineering*, 129(3):278–285, 2003.
- [3] I. Arel, C. Liu, T. Urbanik, and A. Kohls. Reinforcement learning-based multi-agent network traffic signal control. *IET Intelligent Transport Systems*, 4(2):128–135, 2010.
- [4] P. Balaji, X. German, and D. Srinivasan. Traffic signal control using reinforcement learning. *IET Intelligent Transport Systems*, 4(3):171–176, 2010.
- [5] A. G. Barto and S. Mahadevan. Recent advances in hierarchical reinforcement learning. *Discr Dynamic Systems*, 13(4):341–379, 2003.
- [6] J. Baxter, P. L. Bartlett, and L. Weaver. Learning with infinite-horizon, policy-gradient estimation. *Journal of Artificial Intelligence Research*, 2001.
- [7] E. Brockfeld, R. Barlovic, A. Schadschneider, and M. Schreckenberg. Optimizing traffic light cellular automaton model for city traffic. *Review E*, 64(5):056132, 2001.
- [8] Y. K. Chin, N. Bolong, A. Kiring, S. S. Y. K. T. K. Teo. Q-learning based traffic optimization management of signal timing plan. *Intern Journal of Simulation, Systems, Science & Technology*, 12(3):29–35, 2011.
- [9] T. Degris, P. M. Pilarski, and R. S. Sutton. Reinforcement learning with continuous action space. In *American Control Conference 2012*, pages 2177–2182. IEEE, 2012.
- [10] L. Deng. A tutorial survey of architecture algorithms, and applications for deep learning. *APSIPA Transactions on Signal and Information Processing*, 3:e2, 2014.
- [11] M. Duggan, J. Duggan, E. Howley, and E. Barrett. Autonomous network aware vm migration cloud data centres. In *Cloud and Autonomous Computing (ICCAC), 2016 International Conference on*, pages 24–32. IEEE, 2016.
- [12] M. Duggan, K. Flesk, J. Duggan, E. Howley, and E. Barrett. A reinforcement learning approach to dynamic selection of virtual machines in cloud data centres. In *Sixth International Conference on Innovating Computing Technology. IEEE*, 2016.
- [13] S. El-Tantawy, B. Abdulhai, and H. Abdel-Aty. Multiagent reinforcement learning for intersection traffic control.



Download full-text PDF

Read full-text

Download citation

Copy link

Citations (133)

References (53)

... The single objective RL is shown in Fig. 4. The learning agent optimizes the parameter that is used to design the reward function. In traffic signal control, the parameter can be the traffic feature of road intersection like queue length [12], waiting time [11,38], delay [14,39], etc. Here delay is . ...

... The state is designed in many ways in the literature. Some approaches take single traffic feature information for state [22,48] and some consider multiple features [14,38,39,49]. To give the agent more information about the environment, in CRA the state is designed with three traffic features information. ...

... Delay as a reward function is designed using Equation 6. This reward function is used in multiple studies [14,39] since it is also a crucial parameter for minimize. 4. HRA: This approach is proposed in [19]. ...

Adaptive traffic signal control system using composite reward architecture based deep reinforcement learning

[Article](#) [Full-text available](#)

Jan 2021 · [IET INTELL TRANSP SY](#)

 Abu Rafe Md. Jamil ·  Kishan Kumar Ganguly ·  Naushin Nower

[View](#)

... En effet, le peu d'usagers autorisé à sortir serait contraint d'attendre inutilement à certains feux de signalisation, bien que la chaussée soit complètement dégagée. Les méthodes de résolution en ligne doivent donc prendre des décisions à la volée, en tenant compte de la situation présente et en émettant parfois des hypothèses sur le futur (par exemple en analysant un flux de données de capteurs ou de vidéo-surveillance au cours de l'optimisation [114]). ...

... D'autre part, les méthodes en ligne considèrent l'optimisation dynamique des réglages en tenant compte de l'information sur le trafic en temps réel (par exemple à l'aide d'un flux de vidéo-surveillance). Bien que ces méthodes en ligne s'avèrent efficaces à petite échelle, elles sont généralement difficiles à transposer à l'échelle d'une ville entière [160,114,54]. En outre, la grande majorité des feux de circulation est encore paramétrée par des réglages fixes [9,55,121,132,137,106]. Aussi, les contributions présentées dans la section 2.4.1 ainsi qu'aux chapitres 3, 4 et 5 se concentrent essentiellement sur des problèmes d'optimisation hors ligne pour la recherche du réglage optimal

Download

des feux de signalisation préprogrammés, selon la vraisemblance du trafic urbain dans les villes étudiées. ...

Download citation

Copy link

### Méta-modélisation, simulation et optimisation de flux urbains

Thesis

Nov 2020

Florian Leprêtre

[View](#) [Show abstract](#)

... Recent works Gao et al. (2017); Wan and Hwang (2018); Mousavi et al. (2017) use neural networks as function approximators to avoid the dimensionality and computing limitations of table based methods in large state-action spaces, showing DRL TSC can be more efficient than some earlier methods. The first two use discrete cell encoding vectors to represent the system, which are passed to a Convolutional Neural Network (CNN), whereas the second directly uses pixels in the same manner. ...

### Reinforcement Learning for Traffic Signal Control: Comparison with Commercial Systems

Preprint

Full-text available

Apr 2021

[View](#) [Show abstract](#)

... During the decision process, the policy that the agent takes combines both exploitation of already learned policies and exploration of new policies that never met before. Studies using similar RL frameworks to manipulate ATSC are not rare in the last two decades and have provided beneficial references for research [14,15,16,17,18]. For instance, a single-agent model-free Q-learning algorithm was developed for optimizing signal timing in a single intersection [14]. In this study, the authors used queue length as the state representation and accumulative delay between two action cycles as the reward. ...

### Network-wide traffic signal control optimization using a multi-agent deep reinforcement learning

Preprint

Full-text available

Apr 2021

[View](#) [Show abstract](#)

... A widely used class of algorithms in the literature are value-based methods [41,43,45]. These algorithms try to extract the near optimal policy based on the value function, which is defined in Eq. 6. ...

### Availability-aware and energy-aware dynamic SFC placement using reinforcement learning

Article

Full-text available

Apr 2021 · J Supercomput

● Guto Leoni · ● Judith Kelner · ● Patricia Takako Endo

[View](#) [Show abstract](#)

... Recent advances in artificial intelligence (AI) and machine learning have made image-based modeling and analysis (e.g., classification, real time prediction, and image segmentation) even more successful in different applications [23,24,25,26]. Also, with the advent of

Download

nanotechnology semiconductors, a new generation of Tensor Processing Units (TPUs) and Graphical Processing Units (GPUs) can provide an extraordinary computation capability for data-driven methods [27]. ...

Download citation

Copy link

### Aerial imagery pile burn detection using deep learning: The FLAME dataset

[Article](#) [Full-text available](#)
Mar 2021 · [COMPUT NETW](#)

● Alireza Shamsoshoara · ● Fatemeh Afghah · ● Abolfazl Razi · ● Erik Blasch

[View](#) [Show abstract](#)

... In 2015, deep reinforced learning was firstly introduced to traffic signal control optimization in [25] and further refined in 2016 by Van der Pol et al. [26], while considering the coordination of multiple intersections in a small network. In 2017, a traffic signal control policy has been trained by deep policy gradient and applied to a large traffic network by assuming multiple intersections could be controlled with the same agent [27], [28]. The result showed promising potential for policy-based reinforcement learning for traffic signal control. ...

### Boosted Genetic Algorithm using Machine Learning for traffic control optimization

[Preprint](#)

Mar 2021

● Tuo Mao · ● Mihaita Adriana Simona · Fang Chen · Hai L. Vu

[View](#) [Show abstract](#)

... In [29] builds two kinds of reinforcement learning algorithms, namely, deep policy-gradient (PG) and value-function-based agents, that can predict the best possible traffic signal for traffic intersections. The adaptive traffic light control agent receives a snapshot of the current graphical traffic simulator and generates a control signal. ...

### Recent development of smart traffic lights

[Article](#) [Full-text available](#)

Mar 2021

A'isya Nur Aulia Yusuf · Ajib Setyo Arifin · ● F.Y. Zulkifli

[View](#) [Show abstract](#)

... The former aims to find suitable fixed-time signal plans based on historical traffic demand, while the latter dynamically adjusts the signal state according to the traffic information detected in real time. Although traffic responsive methods are technically sound, their performance depends heavily on real-time sensor systems [7] and they are generally difficult to apply to the whole city owing to the high operational cost [8,9]. Besides, the majority of traffic lights in real-world work under fixed signal timing plans and the traffic flows tend to repeat similar patterns like morning and evening peaks. ...

### Surrogate-assisted cooperative signal optimization for large-scale traffic networks

[Preprint](#)

Mar 2021

● Yongsheng Liang · ● Zhigang Ren · Lin Wang · Wenhao Du

[View](#) [Show abstract](#)

Download full-text PDF

Read full-text

Download citation

Copy link

... To solve this problem, a combination of RL with deep neural networks has shown promising results. The authors of the study performed in [36] used value function-based agents and deep policygradient to predict the optimal signal at intersections. A snapshot of the current state is received to these adaptive traffic light controls to produce signals at each time step. ...

Traffic Flow Management of Autonomous Vehicles Using Deep Reinforcement Learning and Smart Rerouting

Article Full-text available

Mar 2021  
Anum Mushtaq · Irfan Ul Haq · Muhammad Usman · Omair Shafiq  
View Show abstract

Show more

Recommendations Discover more

Project

Predictive Communication for UAV Networks  
Abolfazl Razi · Fatemeh Afghah · Jonathan Ashdown · [...] · Kurt Turck  
View project

Project

Intensive Care systems  
Joerg Kampmann · Michael Schukat · E. Schwarzer · [...] · G. Lau  
View project

Project

Digital Certificate-based Port Knocking for Connected Embedded Systems  
Basim Mahbooba · Michael Schukat  
  
This research aims to reinforce existing port knocking methods with a digital certificate for alternative authentication among IoT devices. Such concepts will be complementary to other cryptographi ... [more]  
View project

Conference Paper Full-text available

Traffic Light Control Using Deep Reinforcement Learning Agent  
April 2017  
Sajad Mousavi · Michael Schukat · Enda Howley

Recent advances in combination deep neural network architectures and reinforcement learning techniques have shown promising

[Download full-text PDF](#)[Read full-text](#)[Download citation](#)[Copy link](#)[View full-text](#)Thesis [Full-text available](#)

Researching Advanced Deep Learning Methodologies in Combination with Reinforcement Learning Techniqu...

December 2018

● Sajad Mousavi · ● Michael Schukat · ● Enda Howley

Artificial intelligence (AI) field concerns to build autonomous agents that learn to do tasks successfully in complex and uncertain environments. AI provides powerful techniques which are used to solve many real-world problems ranging from computer science, industry, games, music to hospitals and medicine. What makes it applicable in various domains is a machine learning approach, which is the ... [\[Show full abstract\]](#)

[View full-text](#)

Article

Deep Reinforcement Learning for Traffic Light Control in Vehicular Networks

March 2018

Xiaoyuan Liang · Xunsheng Du · ● Guiling Wang · Zhu Han

Existing inefficient traffic light control causes numerous problems, such as long delay and waste of energy. To improve efficiency, taking real-time traffic information as an input and dynamically adjusting the traffic light duration accordingly is a must. In terms of how to dynamically adjust traffic signals' duration, existing works either split the traffic signal into equal duration or extract ... [\[Show full abstract\]](#)

[Read more](#)Conference Paper [Full-text available](#)

Deep Learning Methodologies in Combination with Reinforcement Learning Techniques

April 2016

● Sajad Mousavi · ● Michael Schukat · ● Enda Howley

Before a reinforcement learning agent (software or hardware) can choose an action, it must have a good representation of the environment in which the agent is to be learned. Hence, perception is one of the key problems that must be solved before the agent can decide to select an optimal action to take. Learning good representations of high-dimensional state or action spaces is a major challenge ... [\[Show full abstract\]](#)

[View full-text](#)Article [Full-text available](#)

Using a Deep Reinforcement Learning Agent for Traffic Signal Control

November 2016

● Wade Genders · ● Saiedeh N Razavi

Ensuring transportation systems are efficient is a priority for modern society. Technological advances have made it possible for transportation systems to collect large volumes of varied data on an unprecedented scale. We propose a traffic signal control system which takes advantage of this new, high quality data, with minimal abstraction compared to other proposed systems. We apply modern deep ... [\[Show full abstract\]](#)

[View full-text](#)



Download full-text PDF

Read full-text

Download citation

Copy link



Company	Support	Business solutions
<a href="#">About us</a>	<a href="#">Help Center</a>	<a href="#">Advertising</a>
<a href="#">News</a>		<a href="#">Recruiting</a>
<a href="#">Careers</a>		