

# Continuous residual reinforcement learning for traffic signal control optimization

Mohammad Aslani, Stefan Seipel, and Marco Wiering

**Abstract:** Traffic signal control can be naturally regarded as a reinforcement learning problem. Unfortunately, it is one of the most difficult classes of reinforcement learning problems owing to its large state space. A straightforward approach to address this challenge is to control traffic signals based on continuous reinforcement learning. Although they have been successful in traffic signal control, they may become unstable and fail to converge to near-optimal solutions. We develop adaptive traffic signal controllers based on continuous residual reinforcement learning (CRL-TSC) that is more stable. The effect of three feature functions is empirically investigated in a microscopic traffic simulation. Furthermore, the effects of departing streets, more actions, and the use of the spatial distribution of the vehicles on the performance of CRL-TSCs are assessed. The results show that the best setup of the CRL-TSC leads to saving average travel time by 15% in comparison to an optimized fixed-time controller.

**Key words:** continuous state reinforcement learning, adaptive traffic signal control, microscopic traffic simulation.

**Résumé :** Le contrôle de feux de circulation peut être naturellement considéré comme un problème d'apprentissage par renforcement. Malheureusement, c'est une des catégories de problèmes d'apprentissage par renforcement les plus difficiles à aborder en raison de son vaste espace d'états. Une méthode directe en vue de relever ce défi est de contrôler les feux de circulation en se basant sur l'apprentissage par renforcement continu. Bien que les contrôleurs de feux de circulation aient eu du succès, ils peuvent devenir instables et ne pas converger aux solutions presque optimales. Nous développons des contrôleurs de feux de circulation adaptatifs en nous basant sur l'apprentissage par renforcement résiduel continu (« CRL-TSC ») qui est plus stable. L'effet de trois fonctions caractéristiques est examiné de façon empirique en ayant recours à une simulation de trafic microniveau. En outre, les effets des rues de partance, de plus d'actions et de l'utilisation de la distribution spatiale des véhicules sur la performance des CRL-TSC sont évalués. Les résultats montrent que la meilleure configuration du CRL-TSC donne lieu à une économie en temps de déplacement moyen de 15% en comparaison avec un contrôleur à temps fixe optimisé. [Traduit par la Rédaction]

**Mots-clés :** état continu d'apprentissage par renforcement, contrôle de feux de circulation adaptatif, simulation du trafic microniveau.

## 1. Introduction

The high trend of population growth in cities and consequently a high level of accumulation and concentration of economic and social activities in urban areas lead to a growing demand for transportation (Bhatta 2010; Rodrigue et al. 2017). Such an increase in transportation demand renders the current transportation infrastructures incapable of successfully handling transportation needs.

Heavy traffic congestion and long vehicle queues on signalized approaches are common phenomena that are currently observable every day in large cities. Traffic congestion usually arises from the excess of demand in comparison to the available capacity of the streets. One of the effective solutions for alleviating traffic congestion is to embed intelligent transportation system (ITS) technologies into the traffic infrastructures to make them work more efficiently (Chowdhury and Sadek 2003). One of the cornerstones of ITS that attracted attention of a lot of researchers and practitioners is developing adaptive traffic signals.

Adaptive traffic signal control is a real-time traffic management strategy in which traffic signal timing changes, or adapts, according to the actual traffic demand. It uses the observed information

to immediately adapt to traffic demand (Aslani et al. 2017; El-Tantawy et al. 2013; Aslani et al. 2018). In this context, multi-agent systems (MAS) have become very popular in traffic control owing to the analogies of their characteristics (e.g., distribution, intelligence, and autonomy) with the traffic control nature (Barcelos de Oliveira and Camponogara 2010). In this study, which uses MAS to control traffic, there are two types of agents: traffic signal agents (active agents) that have learning capabilities and vehicle agents (passive agents) that have behaviors of acceleration, deceleration, and overtaking.

Due to the complexity and uncertainty arising in traffic environments, it is difficult to resolve the problem with preprogrammed multi-agent behaviors. Therefore, a learning mechanism is required by which the agent can gain the necessary knowledge while making a decision and interacting intelligently with an uncertain environment. Within such a context, reinforcement learning serves as a promising approach for training agents such that agents never see examples of the correct behavior, but instead receive positive or negative rewards for the actions they try (Kaelbling et al. 1996; Sutton and Barto 1998). It allows agents to

Received 11 July 2017. Accepted 26 February 2018.

**M. Aslani.** Department of Industrial Development, IT and Land Management, University of Gävle, Gävle, Sweden.

**S. Seipel.** Department of Industrial Development, IT and Land Management, University of Gävle, Gävle, Sweden; Division of Visual Information and Interaction, Department of Information Technology, Uppsala University, Uppsala, Sweden.

**M. Wiering.** Institute of Artificial Intelligence and Cognitive Engineering, University of Groningen, Groningen, the Netherlands.

**Corresponding author:** Mohammad Aslani (email: [moh.aslani@gmail.com](mailto:moh.aslani@gmail.com)).

Copyright remains with the author(s) or their institution(s). Permission for reuse (free in most cases) can be obtained from [RightsLink](https://www.nrcresearchpress.com/cjce).

automatically determine the ideal behavior to maximize their performance (Schwartz 2014; van Otterlo and Wiering 2012). Numerous reinforcement learning algorithms have been developed in the literature; however, the temporal difference learning methods (Sutton 1988) are the most relevant to the traffic signal control problem.

Conventional reinforcement learning methods need to first discretize the state space and then apply a suitable algorithm for a discrete stochastic system. This discretization has some drawbacks. A coarse discretization results in a jerky output and poor performance, while a fine discretization, which may lead to better performance, necessitates not only a large memory storage but also many learning trials. To eliminate these difficulties, we develop adaptive traffic signal controllers founded on continuous reinforcement learning. Continuous reinforcement learning rests on the concept of generalization through function approximators (Sutton and Barto 1998). The idea behind it is that the agent requires no direct experiencing of all states since the values of state-actions are estimated from the values of similar state-actions (Szepesvári 2010). In continuous reinforcement learning, the original state space is mapped onto a feature space through a feature function. The performance of the continuous reinforcement learning methods is highly dependent on the suitability of the selected feature function. With this end in mind, three different feature functions, namely tile coding, triangular-shaped functions (TSFs), and radial basis functions (RBFs) are compared. The combination of reinforcement learning with function approximation may become unstable and fail to converge to a near-optimal solution. To overcome this challenge, the design of adaptive traffic signal controllers is done through residual algorithms (Baird 1999).

The rest of the paper is organized as follows: Section 2 reviews related work. Section 3 describes the principles of reinforcement learning. The proposed adaptive traffic signal controller based on continuous residual reinforcement learning is presented in Section 4. Section 5 presents the experimental setup and results. Section 6 provides a discussion of our findings and Section 7 concludes the paper.

## 2. Related work

Adaptive traffic signal control optimizes the traffic signal scheduling parameters based on current traffic conditions to achieve a set of specific goals. Different methods have been proposed in traffic engineering to adaptively control traffic signals, e.g., SCOOT (Hunt et al. 1981), SCAT (Sims and Dobinson 1980), OPAC (Gartner 1983), PRODYNE (Henry et al. 1983), and RHODES (Head et al. 1992). In recent decades, different methods in artificial intelligence have attracted the interest of experts in traffic signal control. Neural networks (Bishop 1995; Samarasinghe 2016), fuzzy inference systems (Mamdani 1974; Takagi and Sugeno 1985; Zadeh 1965), and reinforcement learning (Sutton and Barto 1998) are three machine learning approaches drawn upon for developing adaptive traffic signal controllers.

In the research done by Srinivasan et al. (2006), two traffic signal control systems based on neural networks were developed. The first system was developed using the integration of simultaneous perturbation stochastic approximation (SPSA) and a fuzzy neural network. In this method, SPSA was used in modeling the weight update process of a five-layer fuzzy neural network. The hybrid neural network based on a multi-agent system, as the second system, was developed to solve the online distributed control problem. Each agent includes a five-layer fuzzy neural network for online decision-making. The learning process contains three steps, namely reinforcement learning, adjustment of learning rates and weights, and adjustments of fuzzy relations. A microscopic traffic simulation of the central business district of Singapore was developed to be used as a test-bed for assessing the

performance of the systems. The results demonstrated that the second system outperforms the first one in more complex scenarios with multiple traffic peaks.

Chiou and Huang (2013) employed the integration of a fuzzy inference system and a stepwise genetic algorithm to develop adaptive traffic signal controllers. Since fuzzy inference systems are not able to learn as such and requires that the knowledge base is derived from experts' knowledge, the stepwise genetic algorithm was deployed to tune both the form of fuzzy membership functions and fuzzy rules. Also, traffic flows and queue lengths were selected as the input variables and the extension of green time was chosen as the control variable. The control system was tested in an isolated intersection and a 1x3 traffic network. Through the experimental results, they conclude that the proposed system is efficient and robust.

In Bi et al. (2014), a distributed traffic signal control system founded on type-2 fuzzy logic control was adopted. A differential evolution method was deployed to tune the knowledge base. The proposed method was benchmarked against type-1 fuzzy logic control and fixed-time methods on a grid-type network composed of 11 intersections. The results revealed that the proposed method has better performance.

In this research, reinforcement learning is employed to develop adaptive traffic signal controllers. Within such a context, Wiering (2000) employed modeled-based reinforcement learning to minimize the waiting time. Vehicles have the ability to communicate with traffic signals. The average waiting time estimated by vehicles is transmitted to the traffic signal located at the next intersection. Then, a traffic signal selects a phase with the maximum sum of the vehicles' gains. The gain is defined as the difference between the vehicle's waiting time when the light is red and when it is green. Results indicated that the proposed method reduces waiting time by 22% in comparison to a fixed-time controller.

In Steingröver et al. (2005), the authors extended Wiering's approach using extra information from neighboring intersections. Although adding new information of the congestion on the next lanes allows the agents to learn how to handle traffic when the departing links are congested, it makes the state-space bigger and decreases the convergence speed.

Salkham et al. (2008) proposed a traffic signal control system using collaborative reinforcement learning in which each signalized intersection learns the suitable phase timing by collaborating with neighboring controllers. A pair of the phase number and its status (busy/not busy) was considered as the state space. Also, the action of each controller was the phase duration. The performance of the proposed method was evaluated in a real-world simulated traffic environment of downtown Dublin. It was benchmarked against a fixed-time system and a SAT-like algorithm (Richter 2006) that emulates SCATS' behavior of saturation balancing. The experiments showed that the proposed system significantly outperforms other methods in terms of average waiting time.

Medina et al. (2010) used reinforcement learning to adaptively control traffic signals. Each traffic signal controller senses the number of vehicles approaching its intersection. The state space is augmented by the numbers of vehicles stopped on departing lanes approaching adjacent intersections. The results revealed that the proposed adaptive traffic signal controllers led to lower delays as well as a more balanced distribution of the delay among all vehicles in comparison to a fixed-time method.

Medina and Benekohal (2012) employed the Q-learning algorithm and an approximate dynamic programming algorithm to develop traffic signal control strategies. At each intersection, the learning controller takes into account not only the local state but also the congestion state of neighboring intersections. A real-world traffic simulation was carried out in VISSIM to test the efficiency of the proposed systems. The proposed systems were

compared with TRANSYT7F and the results indicated that they led to 13% lower average delay.

In El-Tantawy et al. (2013), a coordinated traffic signal control scheme based on multi-agent reinforcement learning was developed. In this scheme, each agent that controls one intersection coordinates its actions with neighboring intersections. The state space contains the index of the current green phase, elapsed time of the current phase, and maximum queue lengths associated with each phase. The action of each agent is to extend the current phase or to switch to another phase. The performance of the proposed scheme was evaluated in a simulated network of 59 intersections in downtown Toronto. The results revealed that their method outperforms the current control scheme of the study area by 26% regarding average travel time.

Employing discrete state reinforcement learning for traffic signal control, which is naturally continuous, may bring about a low convergence speed and poor performance. The more reasonable solution is to employ continuous reinforcement learning that has the ability to perform accurately on unseen data. Within such a context, Prashanth and Bhatnagar (2011) enabled the traffic signal controller to handle large state space by the combination of Q-learning and a function approximator. Queue lengths and elapsed time for the red signal are variables forming the state space. The objective of the controller is to minimize the queue lengths by considering fairness among different approaching links such that no lane has a long red time. The results showed that Q-learning with function approximation significantly outperforms the fixed-time controller.

In Abdoos et al. (2014), the authors proposed a hierarchical multi-agent reinforcement learning architecture to provide different levels of control for a traffic network. The intelligent agents are divided into two groups: local agents that are responsible for controlling each intersection and global agents that adjust actions of the local agents. Local agents and global agents adapt to prevailing traffic conditions through standard Q-learning and continuous Q-learning, respectively. There are nine local agents (3×3 junction grid) and three global agents that each supervises three local agents. The results revealed that the proposed method outperforms standard Q-learning in terms of delay time.

In this research, we develop adaptive traffic signal controllers based on continuous residual reinforcement learning (CRL-TSC) that is more stable and the performance of the best CRL-TSC is compared with fixed-time, standard Q-learning, and actor-critic controllers. Also, the effect of three feature functions, namely tile coding, radial basis functions (RBFs), and triangular-shaped functions (TSFs) are empirically investigated. Moreover, the impacts of considering departing links and the spatial distribution of vehicles in the state space, as well as more actions in the action space are evaluated. Departing link variables provide CRL-TSCs with the ability to handle the spillback phenomenon and indirect cooperation with neighboring intersections. The spatial distribution causes the distance of the vehicles to the associated intersection to be, to some extent, regarded. Investigating the effect of more actions determines if increasing the flexibility in the action space improves the performance.

### 3. Reinforcement learning

Reinforcement learning originally stems from the study of animal intelligence (Thorndike 1911) and has been developed as a major branch of machine learning for solving sequential decision-making problems. Reinforcement learning is an approach to learn an optimal policy of an agent by interacting with its surrounding environment such that it maximizes some numerical value that represents a long-term objective.

In reinforcement learning, the decision-maker is called an intelligent autonomous agent and everything except the agent is referred to as the environment (Sutton and Barto 1998). In many

cases, the environment has the Markovian property with respect to the agent's perception. The Markovian property means that the result of an action does not depend on all previous actions and visited states (history), but only depends on the current state. A fundamental formalism for reinforcement learning, especially in stochastic domains is called a Markov decision process (MDP) (van Otterlo and Wiering 2012). In fact, the basic elements of the reinforcement learning problem can be formalized by an MDP. MDP consists of four elements:  $S$ ,  $A$ ,  $R_{ss'}^a$ , and  $P_{ss'}^a$ , where  $S$  is the set of states,  $A$  is the set of agent's actions,  $P_{ss'}^a$  is the probability of going from state  $s$  to  $s'$  after taking action  $a$ , and  $R_{ss'}^a$  is the average reward for the transition from state  $s$  to  $s'$  by taking action  $a$ . The decision-making function of the agent that specifies which action should be taken in each state is called the policy  $\pi(s, a)$ . In other words, the policy is a mapping from states to actions and indicates the probability of selecting action  $a$  in state  $s$ . In this research, we employ Boltzmann policy to balance between exploration and exploitation.

Another element of reinforcement learning is the use of action values. While a reward function shows how good an action is in an immediate term, an action value specifies how good a particular action is in a long-term sense. The action value that shows the value of performing an action in a state and thereafter following the policy  $\pi$  is calculated by eq. (1).

$$(1) \quad Q^\pi(s, a) = E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\}, 0 \leq \gamma \leq 1$$

where  $Q^\pi(s, a)$  is the state-action value, that corresponds to the expected return when starting in the state  $s$  and taking action  $a$  and following the policy  $\pi$  thereafter;  $r_{t+k+1}$  is the reward obtained when arriving into states  $s_{t+1}$ ,  $s_{t+2}$  etc.;  $\gamma$  is the discount factor that represents the difference in importance between future rewards and instant rewards. The objective of reinforcement learning is to find a policy that maximizes the action values. The state space of the traffic environment is very large and continuous, and this makes conventional reinforcement learning inefficient. In this research, continuous reinforcement learning is employed to tackle this challenge.

### 4. Continuous residual reinforcement learning traffic signal controller

The continuous residual reinforcement learning traffic signal controller (CRL-TSC) is an autonomous learner that iteratively interacts with the traffic environment to find the optimal or near-optimal signal timing plan. CRL-TSC tunes the parameters of the traffic signal controller in response to the changing traffic conditions. At the beginning of each phase, each CRL-TSC senses the current traffic state ( $s_t$ ) of its local intersection. The traffic state is represented by a vector whose elements are the number of vehicles on each approaching street. Through this representation, the traffic load which is easy to be measured through existing sensors is encoded in the definition of the environment state. Also, this state definition allows us to manage vehicles with many passengers (e.g., buses). It should be noted that each dimension of the state space is normalized between 0 and 1.

After sensing the current traffic state, CRL-TSC selects a green time duration (action), i.e., a value from [20, 30, 40, 50, 60, 70, 80, 90] seconds ( $a_t$ ). Once a green time is selected, CRL-TSC waits to the end of the current phase duration which is the summation of the green and yellow interval. Then, CRL-TSC receives a reward signal ( $r_{t+1}$ ). The reward signal provided to CRL-TSC is defined as the negative total number of the vehicles waiting on all the streets leading to the associated junction. Using this reward function causes assigning longer green time to the streets with heavier traffic congestion and higher input traffic flows. In fact, if the selected green time leads to passing (eliminating) more vehicles



from the streets with a high input traffic flow, it receives a greater reward. To put it simply, it considers both traffic congestions and the input traffic flow of the approaching streets.

Once the immediate reward signal ( $r_{t+1}$ ) is received, CRL-TSC senses the new traffic condition ( $s_{t+1}$ ) and selects another green time duration ( $a_{t+1}$ ) for the next traffic light configuration.

Through  $r_{t+1}$ ,  $s_{t+1}$ , and  $a_{t+1}$  the value of ( $s_t$ ,  $a_t$ ) is estimated. The generalization concept is drawn upon for estimating the value of each state-action pair. It enables traffic signals to perform successfully on unseen states. In fact, there is a natural metric on the state space in such a way that close states require similar behaviors. Thus, CRL-TSCs are able to contend with states never exactly experienced before and they can learn efficiently by generalizing from previously (similar, close) experienced states. Indeed, CRL-TSCs require no direct experience of all different environment states and can obtain the value of a state from that of other similar states.

The value of each state-action pair to be approximated at time  $t$  under policy  $\pi$ ,  $Q^\pi(s_t, a_t)$ , is represented as a linear function  $\hat{Q}^\pi(s_t, a_t) = \theta^T \phi(s_t, a_t)$  where  $s_t$  is the state at time  $t$ ,  $\theta$  is a set of scalar weights, and  $\phi$  is a feature function that encodes the similarity relationship of the state-actions with that of their values (Sutton and Barto 1998). Both  $\theta$  and  $\phi$  are  $(n \times k)$ -dimensional vectors where  $n$  is the total number of features and  $k$  is the number of actions,  $\phi(s_t, a_t)$  is defined according to eq. (2).

$$(2) \quad \begin{aligned} \phi^T(s_t, a_t) &= [\varphi_1(s_t) \cdot b_1, \dots, \varphi_n(s_t) \cdot b_1, \varphi_1(s_t) \cdot b_2, \dots, \varphi_n(s_t) \cdot b_k] \\ b_i &= \begin{cases} 1, & \text{if } a_t = a_i \\ 0, & \text{if } a_t \neq a_i \end{cases} \end{aligned}$$

Choosing the right feature function type is critical for successful learning. Among different feature functions employed in linear function approximators, tile coding, RBF, and TSF are the most exploited techniques. Tile coding generalizes the state space into partitions called tilings (Albus 1975). Each tiling consists of a set of non-overlapping grid cells called a tile. The membership value of the triggering state to different tiles is either 0 or 1 (eq. (3)). There is always just one feature active in each tiling layer. Let  $N$  represent the dimension of the state space and  $m_j$  is the number of tiles on  $j$ th dimension. The tile coding features are created as follows:

$$(3) \quad \varphi_i(s_t) = \begin{cases} 0 & \text{if } s_t \notin \text{tile}_i \\ 1 & \text{if } s_t \in \text{tile}_i \end{cases} \quad 1 < i < n, \quad n = m_1 \times m_2 \times m_3 \dots \times m_N$$

Their reliance on a set of binary features makes tile coding one of the most explored feature functions. In this research, each state variable is partitioned into a finite set of tiles and then the tiling is created by combining the tiles in each state variable in a vector. Each tiling has the same number of tiles in each dimension.

RBF provides a continuous representation of states instead of a binary representation. In fact, RBF builds a more complex representation using a distance metric leading to non-binary features. The activation of each RBF feature continuously decays away from the center of the RBF. The output of the  $i$ th RBF centered around  $s_t$  is calculated according to eq. (4).

$$(4) \quad \varphi_i(s_t) = e^{-\sum_{j=1}^N \frac{(s_t^j - \mu_{ij})^2}{2\sigma_{ij}^2}} \quad 1 < i < n, \quad n = m_1 \times m_2 \times m_3 \dots \times m_N$$

where  $\sigma_{ij}$  and  $\mu_{ij}$  are the standard deviation and center of RBF $_i$  on the  $j$ th dimension and  $s_t^j$  is the  $j$ th dimension of the state at time  $t$ . Figure 1 shows the parameters of RBFs and how RBFs are located on each dimension. As it is clear, the distance between the centers

Fig. 1. Layout of RBFs and their parameters on the  $j$ th dimension.

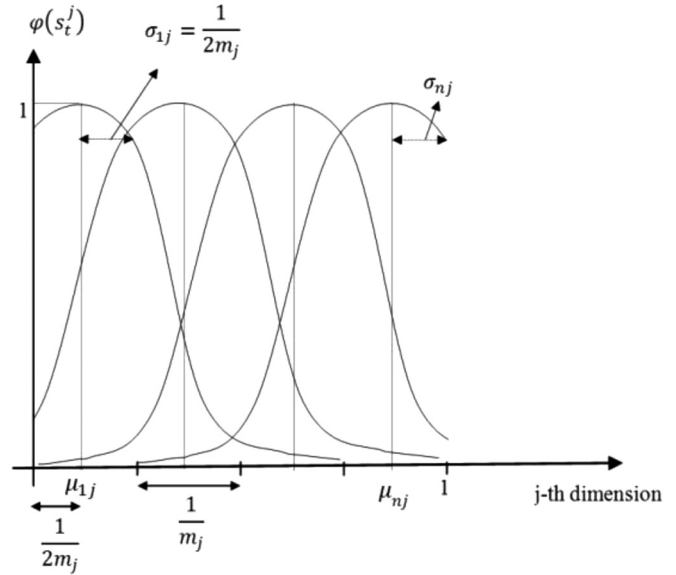
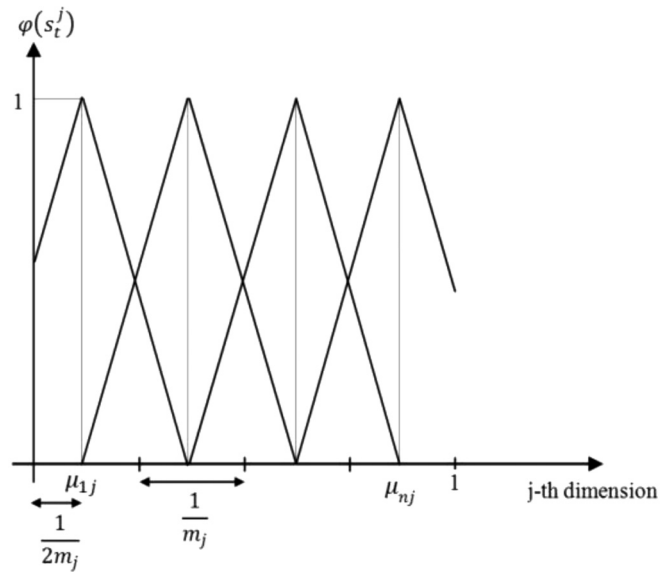


Fig. 2. Layout of TSFs and their parameters on the  $j$ th dimension.



of two consecutive RBFs and the standard deviation of each RBF on the  $j$ th dimension are  $\frac{1}{m_j}$  and  $\frac{1}{2m_j}$  respectively, where  $m_j$  is the number of RBFs on the  $j$ th dimension. In fact, the centers of RBFs are distributed in the state space as a fixed uniform grid.

TSF is a function whose figure takes the shape of a triangle. Equation (5) is used to calculate the degree of membership to different TSFs in each state variable.

$$(5) \quad \varphi_i(s_t) = \prod_{j=1}^N (1 - |s_t^j - \mu_{ij}| \cdot m_j) \quad 1 < i < n, \quad n = m_1 \times m_2 \times m_3 \dots \times m_N$$

The output of a TSF is zero when state  $s$  is far from the center ( $\mu$ ). Figure 2 indicates the arrangement of TSFs on the  $j$ th dimension. It is evident that maximally two features per dimension become active in each state.

The number of tiles, RBFs, and TSFs, as well as the locations of RBFs and TSFs centers, greatly affect the accuracy and validity of

Fig. 3. CRL-TSC.

```

Initialize  $\theta, \phi, \alpha, \gamma$ , and  $\omega$ 
A is the action set
 $t \leftarrow 0$ 
loop
   $s_t, a_t \leftarrow$  initial state and action of the episode
  repeat
    Set the current phase duration to  $a_t + \text{yellow time}$ 
    Wait until the end of the phase
    Observe the number of vehicles on each approaching street
    Calculate reward  $r_{t+1}$ 
     $\delta_{t+1} \leftarrow r_{t+1} - \theta^T \phi(s_t, a_t)$ 
    Estimate the state  $s_{t+1}$ 

    for all  $a_i \in A$  do
       $Q(s_{t+1}, a_i) \leftarrow \theta^T \phi(s_{t+1}, a_i)$ 
       $\rho \leftarrow \rho + \exp(\omega \cdot Q(s_{t+1}, a_i))$ 
    end for
    Uniformly draw a number  $P \in [0, 1]$ 
     $d \leftarrow 0$ 
    for all  $a_i \in A$  do //Boltzmann
       $Pr(s_{t+1}, a_i) \leftarrow \frac{\exp(\omega \cdot Q(s_{t+1}, a_i))}{\rho}$ 
      if  $P \leq Pr(s_{t+1}, a_i) + d$ 
         $a_{t+1} \leftarrow a_i$ 
        break
      else
         $d \leftarrow d + pr(s_{t+1}, a_i)$ 
      end if
    end for

     $Q(s_{t+1}, a_{t+1}) \leftarrow \theta^T \phi(s_{t+1}, a_{t+1})$ 
     $\delta_{t+1} \leftarrow \delta_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ 
     $\Delta\theta \leftarrow \alpha \delta_{t+1} \phi(s_t, a_t)$ 
     $\Delta\theta' \leftarrow \alpha \delta_{t+1} (\phi(s_t, a_t) - \gamma \phi(s_{t+1}, a_{t+1}))$ 
     $\beta \leftarrow \frac{\Delta\theta \cdot \Delta\theta'}{\Delta\theta \cdot \Delta\theta' - \Delta\theta' \cdot \Delta\theta'}$ 
     $\Delta\theta'' \leftarrow (1 - \beta) \Delta\theta + \beta \Delta\theta'$ 
     $\theta_{t+1} \leftarrow \theta_t + \Delta\theta''$ 
     $t \leftarrow t + 1$ 
  until  $s_t$  in terminal
end loop

```

the learning performance of CRL-TSCs. Also, the poorly placed tiles, RBFs, and TSFs can prevent CRL-TSCs to correctly estimate the value function even in some simple domains. Therefore, in our experiments, we generate different feature functions of tile coding, RBF, and TSF with a different density of features on each state variable. In fact, we evaluate and compare the performance of CRL-TSCs by considering different numbers of tiles, RBFs, and TSFs on each dimension (see Section 5.1).

In all feature functions, after the values of features are calculated, they are normalized so that the total sum of them becomes 1. The scalar weights vector ( $\theta$ ) is updated such that the following cost function is minimized (eq. (6)).

$$(6) \quad C = E[(r_{t+1} + \gamma \theta^T \cdot \phi(s_{t+1}, a_{t+1}) - \theta^T \cdot \phi(s_t, a_t))^2]$$

A good strategy for minimizing eq. (6) is to try to minimize it on the observed examples. Stochastic gradient descent is able to do this by tuning the scalar weight vector after each example observed. Therefore, by using stochastic gradient descent,  $\theta$  is updated according to eq. (7).

$$(7) \quad \Delta\theta = \alpha(r_{t+1} + \gamma \theta_t^T \cdot \phi(s_{t+1}, a_{t+1}) - \theta_t^T \cdot \phi(s_t, a_t)) \cdot \nabla_{\theta} Q(s_t, a_t)$$

$$\theta_{t+1} = \theta_t + \Delta\theta$$

In this equation,  $\alpha$  is the learning rate and  $\nabla_{\theta} Q(s_t, a_t) = \phi(s_t, a_t)$ . Although this method is a very simple and fast way for updating  $\theta$ , it is not guaranteed to converge due to the fact that changing the value of one state usually changes the values of other states such as that of the state  $s_{t+1}$ . Consequently, the estimated value ( $\hat{Q}(s_t, a_t)$ ) may move away from its target value. To tackle this problem, the scalar weight vector ( $\theta$ ) can also be updated according to eq. (8).

$$(8) \quad \Delta\theta' = \alpha(r_{t+1} + \gamma \theta_t^T \cdot \phi(s_{t+1}, a_{t+1}) - \theta_t^T \cdot \phi(s_t, a_t)) \cdot (\nabla_{\theta} Q(s_t, a_t) - \gamma \nabla_{\theta} Q(s_{t+1}, a_{t+1}))$$

$$\theta_{t+1} = \theta_t + \Delta\theta'$$

where  $\nabla_{\theta} Q(s_{t+1}, a_{t+1}) = \phi(s_{t+1}, a_{t+1})$ . This residual learning method considers the states  $s_t$  and  $s_{t+1}$  to improve stability and convergence properties. However, this method does not always learn as quickly as eq. (7). In fact, eq. (7) is fast but unstable; eq. (8) can be stable but slow in terms of convergence. Therefore, the best solution is to combine the two methods to gain the advantages of them (fast and stable learning). This can be achieved by using a weighted average of two gradient vectors (eq. (9)).

$$(9) \quad \Delta\theta'' = (1 - \beta) \Delta\theta + \beta \Delta\theta'$$

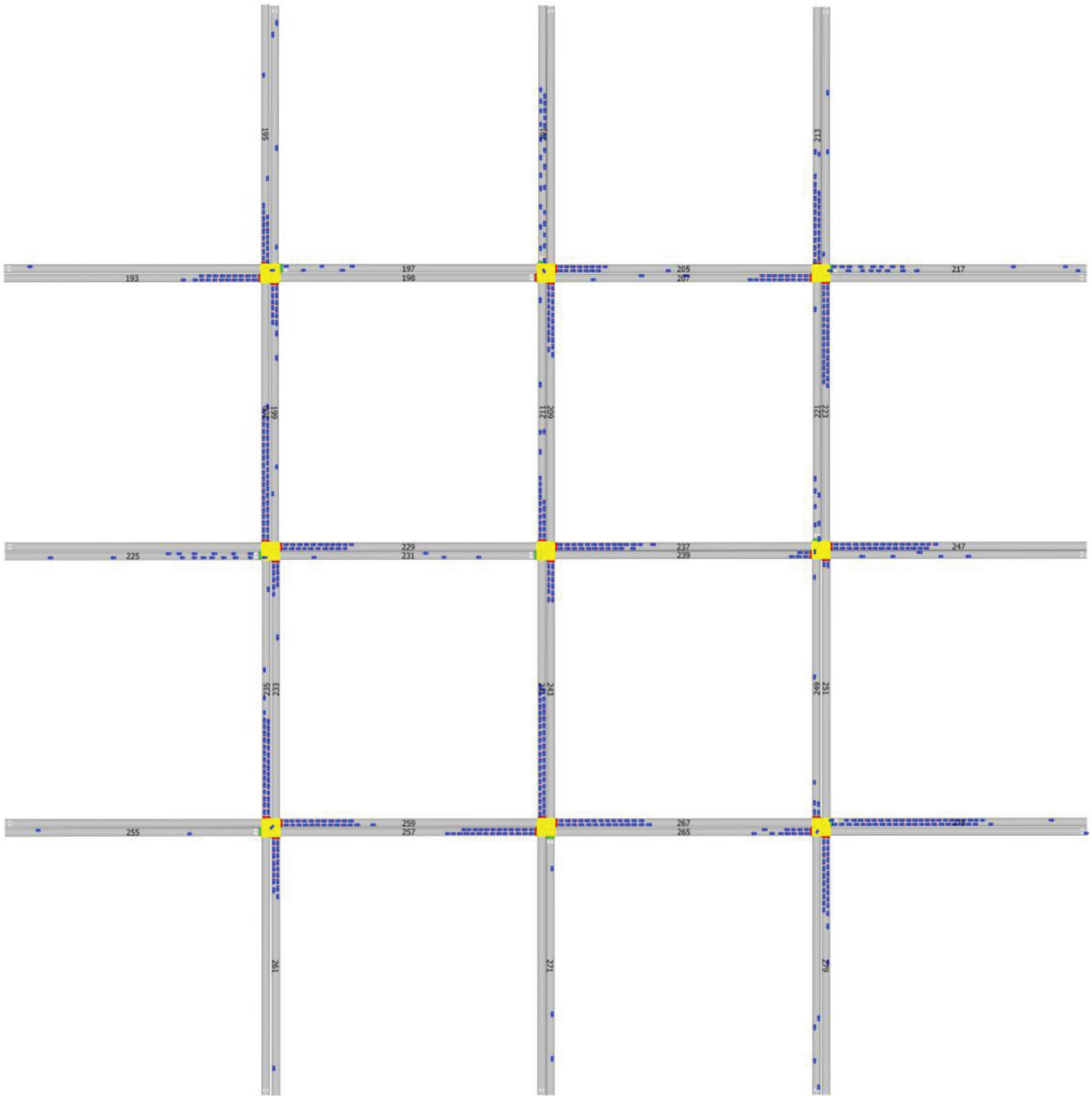
$$\Delta\theta'' = \alpha(r_{t+1} + \gamma \theta_t^T \cdot \phi(s_{t+1}, a_{t+1}) - \theta_t^T \cdot \phi(s_t, a_t)) \cdot (\nabla_{\theta} Q(s_t, a_t) - \beta \gamma \nabla_{\theta} Q(s_{t+1}, a_{t+1}))$$

$$\theta_{t+1} = \theta_t + \Delta\theta''$$

where  $\beta \in [0, 1]$  attenuates the effect of the successor state ( $s_{t+1}$ ). In this research,  $\beta$  is adapted during the learning process using eq. (10) (Baird 1999).

$$(10) \quad \beta = \frac{\Delta\theta \cdot \Delta\theta'}{\Delta\theta \cdot \Delta\theta' - \Delta\theta' \cdot \Delta\theta'}$$

Another point to be noted is the way how CRL-TSCs select the suitable actions in each traffic state. Basically, action selection should be based on the value of state-action pairs ( $Q(s_t, a_t)$ ). How-

**Fig. 4.** Microscopic traffic simulation. [Colour online.]

ever, owing to the fact that CRL-TSCs do not possess the correct values of each state-action pair at the beginning of the learning process, they need to explore different green time durations regardless of their values to achieve accurate estimations of state-action values. As time ( $t$ ) goes by and CRL-TSCs obtain better estimations, they should rely less on an exploration through random green time selection and begin to exploit their obtained knowledge of the traffic environment by choosing those green times that possess a fairly high value. In this research, to trade-off between exploration and exploitation, the Boltzmann exploration strategy is employed. The probability of selecting each available action is calculated according to eq. (11).

$$(11) \quad \Pr(s_t, a_i) = \frac{\exp(\omega \cdot Q(s_t, a_i))}{\sum_{j=1}^k \exp(\omega \cdot Q(s_t, a_j))}$$

where  $k$  is the number of actions,  $a_i$  is each action available in state  $s_t$ , and the parameter  $\omega$  controls the exploration rate. The higher the  $\omega$  value, the sharper the distribution becomes. For  $\omega \rightarrow \infty$ , it converges to the greedy policy. By using the Boltzmann policy, actions with high values are more likely to be selected than actions with a lower value. Figure 3 shows the algorithm of each CRL-TSC.

**Table 1.** Traffic network configuration.

Properties	Value
Number of intersections	9
Number of links	48
Average length of links	250
Number of lanes per links	2
Maximum speed	50 km/h
Number of input/output centroids	12
Arrival distribution	Gaussian
Simulation duration	700 h

**Table 2.** Parameters of vehicles.

Properties	Mean value	Standard deviation
Reaction time	1 s	0.0 s
Reaction time at stop	1.35 s	0.0 s
Length	4 m	0.5 m
Width	2 m	0.0 m
Maximum speed	100 km/h	10 km/h
Maximum acceleration	3 m/s <sup>2</sup>	0.2 m/s <sup>2</sup>
Maximum deceleration	6 m/s <sup>2</sup>	0.5 m/s <sup>2</sup>

## 5. Simulation experiments

### 5.1. Implementation

The microscopic traffic simulation in this research comprises a traffic network, vehicles, and CRT-TSCs. The employed 3×3 grid network for which one CRL-TSC controls one intersection is depicted in Fig. 4. All the streets are two-way (bi-directional) with two lanes on each side. The capacity of each lane is 40 vehicles. The length of each street is 250 m. Vehicles enter the network using a Gaussian distribution through 12 sources that lie on the borders of the network. In each intersection, it is assumed that among all the vehicles approaching the intersection, 33.3% of them go straight, 33.3% turn left, and 33.3% turn right. The traffic network configuration parameters are shown in Table 1. We have used this small grid to accomplish a careful analysis of the impact of different parameters on the performance. However, the proposed CRL-TSC can be easily used in larger traffic networks.

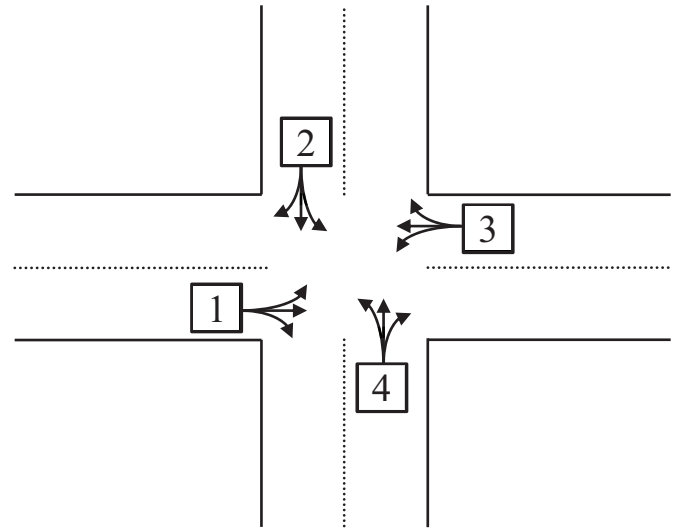
The movement of a vehicle depends on the external properties of the vehicle (vehicle type), such as length, width, maximum speed, and acceleration, as well as internal characteristics of the human driver including reaction time (s) and reaction time at stop (s) (Casas et al. 2010). The parameters incorporated in the vehicle movements are shown in Table 2.

The driving speed ( $v_d$ ) is determined by taking four factors into account: maximum speed ( $v_m$ ), section speed limit ( $v_s$ ), speed acceptance factor ( $f_s$ ), and speed of the following vehicle. The section speed limit is the maximum allowed speed of the vehicles passing a section. The maximum allowed speed of all the sections is 50 km/h. The speed acceptance factor shows the degree of accepting the speed limits of sections by a driver. The value of  $f_s$  for each vehicle is drawn from a Gaussian distribution function with the mean of 1.1 and standard deviation of 0.1. The driving speed is calculated by  $v_d = \text{minimum}(v_m, f_s \times v_s)$ . Also, the driving speed momentarily changes based on the speed of the vehicle ahead. If the vehicle in front has a lower speed, the follower should reduce its driving speed according to the car following model (Gipps 1981) or change its driving lane (Gipps 1986).

The traffic simulation is carried out for 300 h. Also, each one hour is referred to as an episode. Since all the intersections are 4-way crossroads, the system contains homogeneous CRL-TSCs. A signals group is assigned to each approaching street as shown in Fig. 5. Thus, each CRL-TSC controls a signalized intersection that has four phases each assigned to an approaching link.

The point that is important in the learning of CRL-TSCs is the value of the learning rate ( $\alpha$ ) and discount factor ( $\gamma$ ). The best found values for the learning rate for tile coding, RBF, and TSF are 0.1, 0.075, and 0.075, respectively. Also, the discount factor is set to 0.99. It should be noted that these values were obtained by trial and error. Moreover, the value of  $\omega$  (Boltzmann parameter) increases from 0.0 to 10.0 during the first 200 episodes (training period) and then it is kept constant at 10.0 over the last 100 episodes (test period) to evaluate the learning performance of the system.

Another point that is the key to success of the CRL-TSCs' effectiveness is the number of tiles, RBFs, and TSFs. Choosing the

**Fig. 5.** Order of the phases.

wrong number of tiles, RBFs, and TSFs can ruin the generalization property of CRL-TSCs. Hence, to show the impact of the set-up of tiles, RBFs, and TSFs on the learning performance and find the optimal ones, the performance of CRL-TSCs are evaluated based on different numbers of tiles (3, 5, 7, and 9), RBFs (3, 5, 7, and 9), and TSFs (3, 5, 7, and 9) on each dimension of the state space. Due to the space limitations in this paper, the results of investigating the optimal number of tiling layers in the tile coding approach are not presented and it is directly set to 3. In fact, 3 tilings worked best in the preliminary experiments.

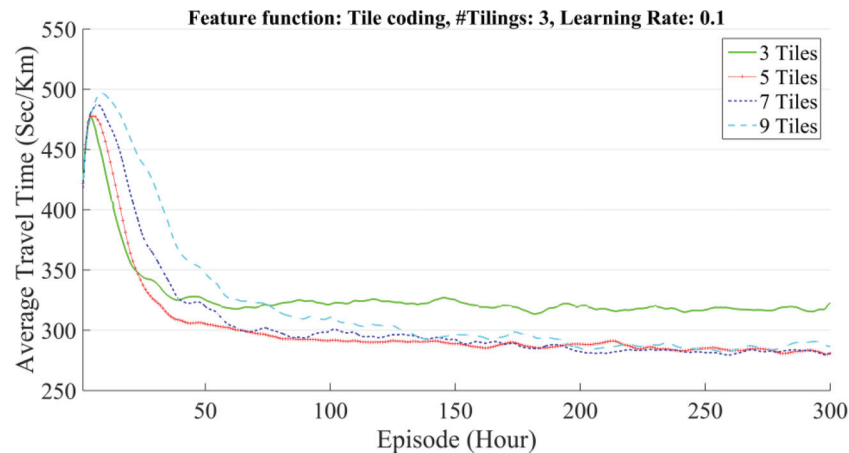
Since each CRL-TSC controls one intersection and each intersection has four approaching streets, the state space has a dimension of 5 (index of the current green phase as well as 4 approaching streets). The total number of features for each CRL-TSC based on the RBFs and TSFs approaches is  $n = m \times m \times m \times m \times \text{ph}$ , where  $m$  is the number of features on each dimension (e.g., 3, 5, 7, and 9) and  $\text{ph} = 4$  is the number of phases of the traffic signals. Also, the total number of features based on the tile coding approach is  $n = k \times (m \times m \times m \times m \times \text{ph})$ , where  $k = 3$  is the number of tiling layers.

The traffic simulation has been implemented using AIMSUN. Three indexes, average travel time (s/km), average stop time (s/km), and average stop numbers (#/veh/km) are used to assess the performance of CRL-TSCs. The average travel time is the average time that a vehicle requires to travel one kilometre. The average stop time is the average time that a vehicle stays at a standstill status in traveling one kilometre inside the traffic network. The average stop number is the average number of stops per vehicle per kilometre.

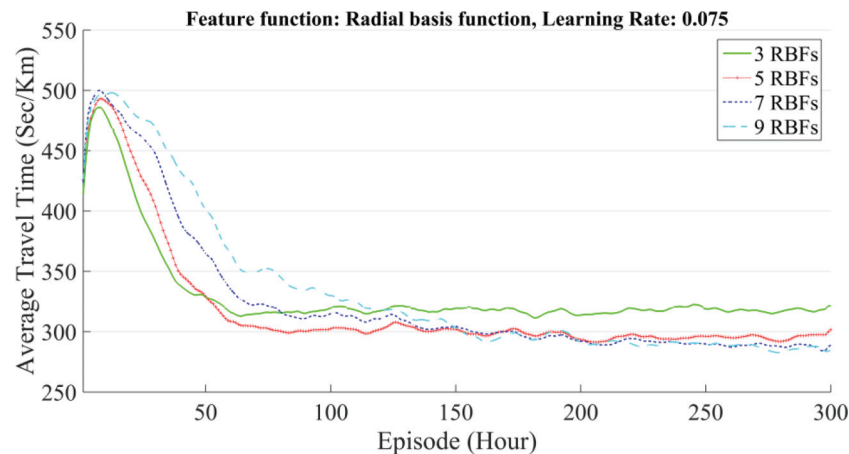
### 5.2. Results

Figure 6 presents the performance of CRL-TSCs with 3 tilings and different numbers of tiles (3, 5, 7, and 9) in terms of average

**Fig. 6.** The learning performance of CRL-TSCs for different numbers of tiles for tile coding. [Colour online.]



**Fig. 7.** The learning performance of CRL-TSCs for different numbers of RBFs. [Colour online.]



travel time. The learning curves are the average of repeated (five times) simulations. As it is clear, 3 tiles has the poorest performance due to the lack of segmentation in the state variables. In other words, the number of features is not high enough to provide the CRL-TSCs with fairly flexible generalization ability to adapt to different traffic states. Increasing the number of tiles from 3 to 5 results in considerable improvement in comparison to 5 to 9 tiles. Indeed, the considerable difference between the learning performance of 3 and 5 tiles indicates the pivotal role of the number of tiles and their set-up. It is evident that the average performance of 5, 7, and 9 tiles over the last 100 episodes are very close to each other.

The learning curves of CRL-TSCs for different RBFs are depicted in Fig. 7. Three RBFs is outperformed by others (5, 7, and 9 RBFs) owing to the insufficient number of features. Similar to tile coding, there is a significant difference between 3 RBFs and others, that proves the importance of the number of RBFs. Also, the performance of 7 RBFs is almost in line with 9 RBFs over the last 100 episodes. Comparing Fig. 7 with Fig. 6 reveals that tile coding slightly outperforms RBF in terms of average travel time.

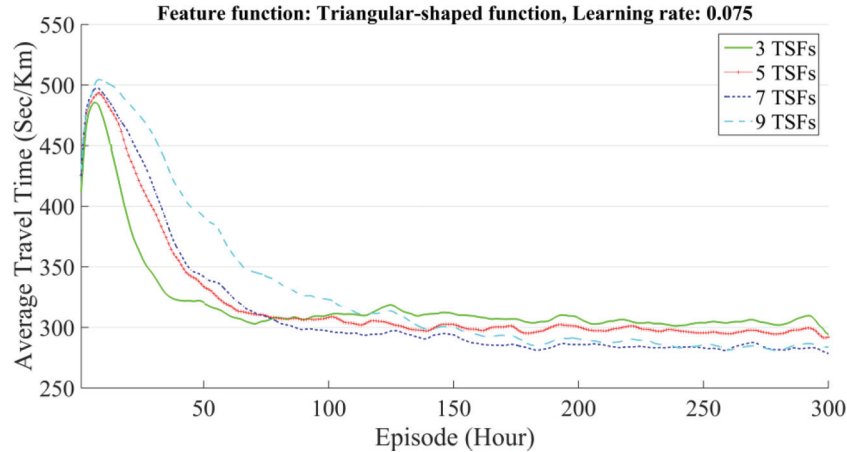
Figure 8 shows the performance of CRL-TSCs for different numbers of TSFs. Similar to tile coding and RBF, the convergence speed decreases as the number of TSFs increases. This is due to the increase in a number of scalar weights parameters. Comparing Figs. 6–8 indicates that 3 TSFs has better performance in comparison to 3 RBFs and 3 tiles.

Table 3 compares the average performance, as measured by the three performance indicators described earlier, of different CRL-TSCs over the last 100 episodes. The best feature functions are shown in boldface. It is evident that 7 features in all three feature functions lead to the best outcome. Also, increasing the number of features from 3 to 7 has improved the performance, but increasing from 7 to 9 could not make the results better. Therefore, increasing the number of tiles does not necessarily lead to an improvement in the performance of the system. Comparing the results reveals that the CRL-TSC with 3 tilings and 7 tiles is the best controller.

The next issue is the impact of the action space on the performance of CRL-TSCs. The small action space leads to a high learning speed, but at the same time might result in lower flexibility. In fact, it should be determined if increasing the flexibility in the action space results in a significant improvement in CRL-TSC performance. With this end in view, the performance of CRL-TSC with a higher number of actions is evaluated. The new action space has 16 actions with steps of 5 s, i.e., [20, 25, 30, ..., 90] s. This allows us to determine whether the shorter green times are crucial in traffic signal control. Figure 9 indicates the performance of CRL-TSC with the new action space and compares it with the initial controller. As it is clear, increasing the number of actions from 8 to 16 does not significantly affect the final performance of CRL-TSCs. Therefore, using the new action space is unreasonable as it just increases the state-action space size.

Another issue is the effect of augmenting the state space with the departing links connected to the junction. Considering de-



**Fig. 8.** The learning performance of CRL-TSCs for different numbers of TSFs. [Colour online.]**Table 3.** Average performance of CRL-TSCs for different numbers of tiles, RBFs, and TSFs over the last 100 episodes.

Feature function	Avg. travel time (s/km)	Avg. stop time (s/km)	Avg. stop numbers (#/veh/km)
3 tiles	318±5	231±5	4.52±0.05
5 tiles	285±4	199±4	4.07±0.05
7 tiles	<b>282±4</b>	<b>196±4</b>	<b>4.02±0.05</b>
9 tiles	286±5	200±5	4.02±0.05
3 RBFs	318±6	231±6	4.40±0.07
5 RBFs	295±4	209±4	4.17±0.05
7 RBFs	<b>289±5</b>	<b>203±5</b>	<b>4.08±0.06</b>
9 RBFs	289±6	203±6	4.08±0.06
3 TSFs	304±5	218±5	4.22±0.05
5 TSFs	297±4	211±4	4.21±0.05
7 TSFs	<b>284±4</b>	<b>198±4</b>	<b>4.04±0.05</b>
9 TSFs	286±5	200±5	4.05±0.05

parting links in the state space provides the CRL-TSCs with the ability to handle the spillback phenomenon and indirect cooperation with the neighboring traffic signals. On the other hand, it leads to a substantial increase in the state space size. The new state space of each agent has a dimension of  $5+D$ . The first five elements are the index of the current green time and the number of vehicles waiting on each of the four approaching streets. The remaining components indicate the number of vehicles on each departing street. Figure 10 shows the performance of CRL-TSC with the augmented state space and compares it with the initial one. It is evident that considering the departing streets improves the performance of CRL-TSCs.

The average performance of CRL-TSCs with the new state space and action space over the last 100 episodes is indicated in Table 4. Increasing the number of actions could not significantly affect the learning performance, but adding the departing streets could improve the performance by 3.2%.

## 6. Discussion

Discrete reinforcement learning methods have been widely used for adaptive traffic signal control (Abdoos et al. 2011, 2013). To validate the performance of the best CRL-TSCs, we benchmark its performance against standard Q-learning (Watkins and Dayan 1992) and actor-critic (Konda and Borkar 1999). Q-learning is categorized as an off-policy method that learns the values of state-actions on the basis of the optimal policy independent of the policy being followed. During the learning process, the agent stores a particular Q-value for each state-action pair. Equation (12)

shows how state-action values are updated (Sutton and Barto 1998).

$$(12) \quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

where  $0 < \alpha \leq 1$  is the learning rate,  $\gamma$  is the discount factor, and  $s_t$  is the discretized state. The optimal values of the discount factor and learning rate are 0.99 and 0.01. Unlike Q-learning, actor-critic is an on-policy method for which state-action values are updated based on the policy being followed. It has a separate memory structure for estimating state-values (the critic) and a mapping from states to a preference value for each action (the actor) (van Otterlo and Wiering 2012). The values of states are updated according to the following equation:

$$(13) \quad V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

In the actor, the values of all state-action pairs are updated by eq. (14).

$$(14) \quad P(s_t, a_t) \leftarrow P(s_t, a_t) + \beta[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

where  $0 < \beta \leq 1$  is the learning rate of the actor and  $P(s_t, a_t)$  indicates the tendency to select action  $a_t$  in state  $s_t$ . The learning rates and discount factor of the actor-critic method are respectively set to 0.01 and 0.99.

For employing discrete reinforcement learning, it is necessary to discretize the state space. To do so, each state variable is first discretized into a finite set of regions. Then, the union of all state variables covers the whole state space. In this way, the total number of states is the product of the number of regions for each discretized state variable. In this research, the first four state variables (approaching streets) are discretized into 6 regions and the last state variable (index of the current green phase) is discretized into 4 regions. Thus, the total number of states is  $n = 6 \times 6 \times 6 \times 6 \times 4 = 5184$ .

**The Boltzmann method is used for balancing between exploration and exploitation.** Similar to CRL-TSC the Boltzmann parameter gradually increases from 0.0 to 10.0 over the first 200 episodes and then it is kept constant at 10.0 over the last 100 episodes. Figure 11 compares the performance of the best CRL-TSC with standard Q-learning and actor-critic in terms of average travel time.

As depicted in Fig. 11, the best CRL-TSC outperforms the other reinforcement learning methods. Table 5 compares the average performance of these controllers based on average travel time,

Fig. 9. The learning performance of CRL-TSC for two different action spaces. [Colour online.]

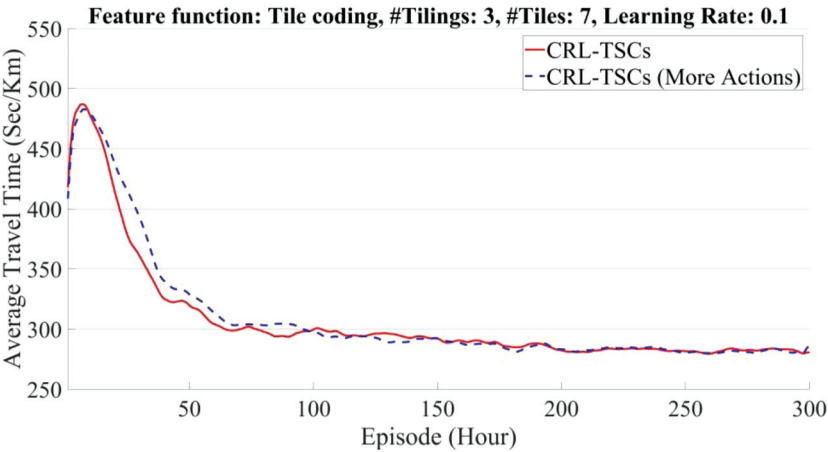


Fig. 10. The learning performance of CRL-TSC for two different state spaces. [Colour online.]

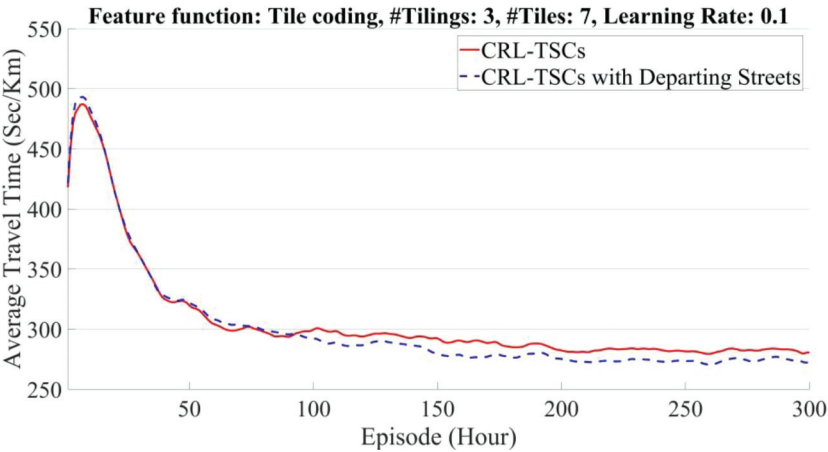


Table 4. Comparison of different CRL-TSCs with different state and action spaces.

Controller	Avg. travel time (s/km)	Avg. stop time (s/km)	Avg. stop numbers (#/veh/km)
1-CRL-TSC	282±4	196±4	4.02±0.04
2-CRL-TSC with more actions	282±4	196±4	3.97±0.05
3-CRL-TSC with departing links	273±4	191±4	3.89±0.05
% Improvement controller 1 vs. 2	0	0	1.2
% Improvement controller 1 vs. 3	3.1	2.6	3.2

stop time, and stop numbers over the last 100 episodes. It is clear that the best CRL-TSC saves average travel time by 30% and 20% in comparison to actor-critic and Q-learning.

Increasing the air pollution and fuel consumption are two destructive consequences of inefficient traffic control systems. In this paper, we evaluate the impact of the best CRL-TSC on decreasing air pollution and fuel consumption. With this end in mind, fuel consumption and emission rates (HC, CO, and NO<sub>x</sub>) are incorporated into the microscopic traffic simulation.

Vehicle specific power (VSP), a technique featuring engine power, is utilized to model the vehicle fuel consumption rate. The VSP shows how the vehicle operating conditions affect the fuel consumptions. The VSP for typical light-duty vehicles depends on the speed, acceleration (deceleration), and roadway slope on the basis of second-by-second cycles (Jiménez-Palacios 1999). According to the study conducted by the Air Quality Control Company of

Tehran (AQCC 2014), the exponential function fits the relationship well between the VSP and the fuel consumption rate (eq. (15)).

(15) 
$$FC\left(\frac{L}{s}\right) = a \times e^{b \times VSP} + c \times VSP + d$$

In eq. (15); *a*, *b*, *c*, and *d* are the constant coefficients. The constant coefficients have been calibrated by the Air Quality Control Company of Tehran (AQCC 2014).

Emission rates of vehicles depend on different parameters such as vehicle speed, vehicle mileage, engine temperature and vehicle load. However, only the first parameter (i.e., vehicle speed) is considered for modeling emission rates in the present research. A 6th order polynomial function is employed for modeling the emission rate (eq. (16)) (Boulter et al. 2009; TRL 1999).

(16) 
$$\begin{aligned} \text{Emission rate}\left(\frac{g}{km}\right) \\ = \frac{A + Bv + Cv^2 + Dv^3 + Ev^4 + Fv^5 + Gv^6}{v} \end{aligned}$$

where A–G are coefficients and *v* is the speed of the vehicle in km/h. The coefficients in eq. (16) have been calibrated for different air pollutants (CO, HC, and NO<sub>x</sub>) by AQCC (2014). The total fuel consumption and traffic-generated air pollution for the best CRL-TSC, actor-critic, and Q-learning are presented in Fig. 12. It is clear

Fig. 11. Comparing the performance of the best CRL-TSC with standard Q-learning and actor-critic. [Colour online.]

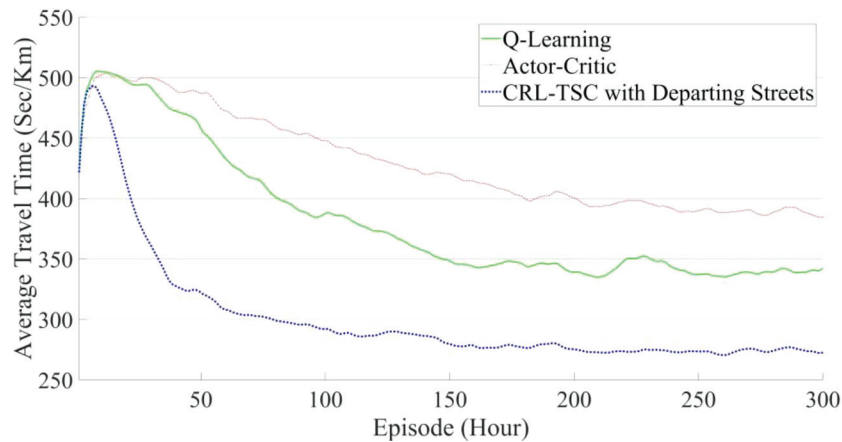
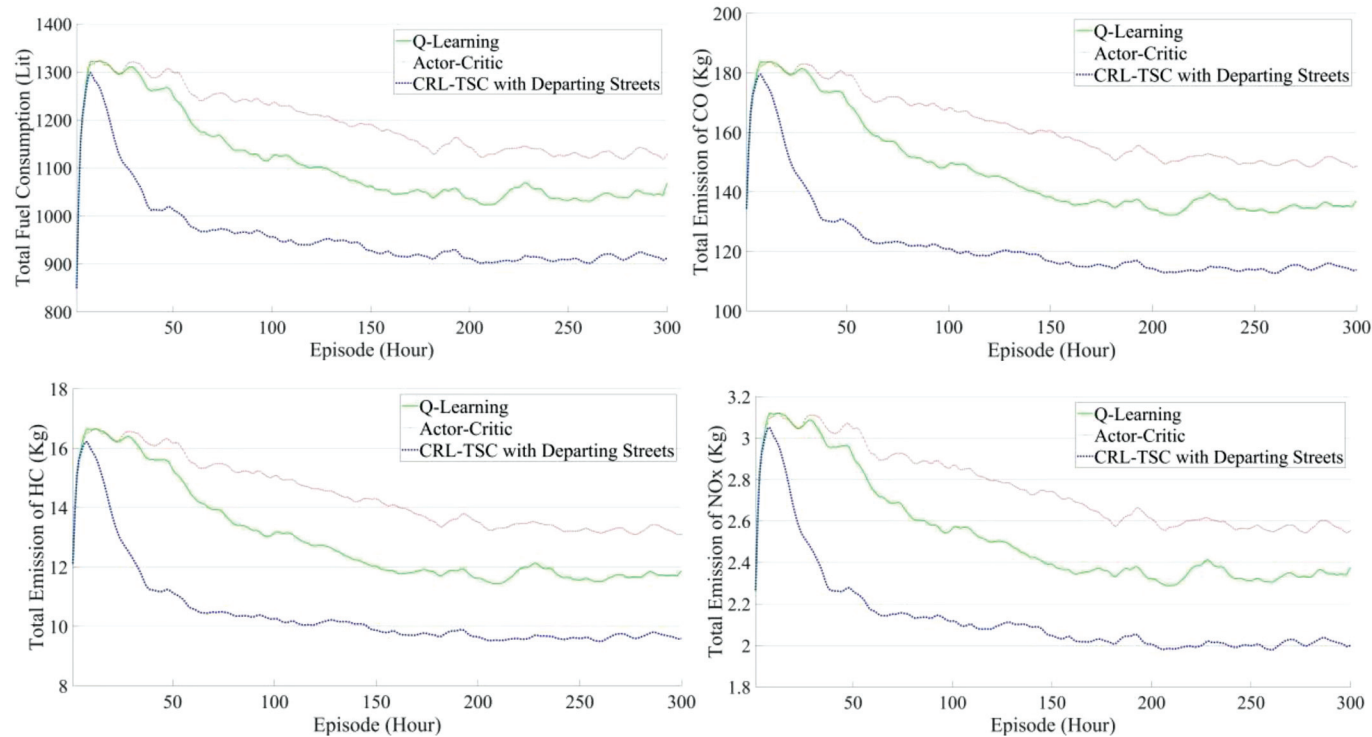


Table 5. Comparison of the best CRL-TSC, standard Q-learning, and actor-critic.

Controller	Avg. travel time (s/km)	Avg. stop time (s/km)	Avg. stop numbers (#/veh/km)
Actor-critic	392±7	304±7	4.44±0.06
Q-learning	340±8	254±7	4.26±0.07
CRL-TSC with departing links	273±4	191±4	3.89±0.05
% improvements CRL-TSC vs. actor-critic	30.4	37.2	12.39
% improvements CRL-TSC vs. Q-learning	19.70	24.8	8.7

Fig. 12. Learning curves of the best CRL-TSC, discrete state Q-learning, and actor-critic in terms of fuel consumption and emissions of air pollutants. [Colour online.]



that the CRL-TSC has the best performance with regard to air pollution and fuel consumption.

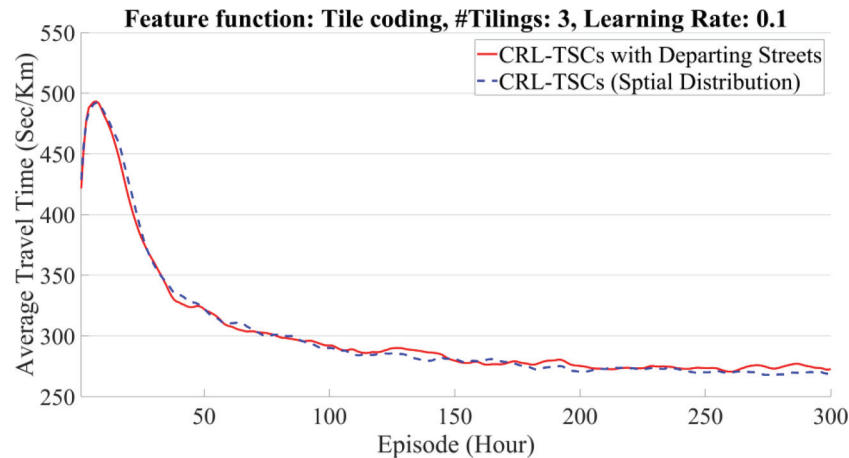
Also, Table 6 compares the average performance of these controllers in terms of fuel consumption and air pollution over the last 100 episodes. It is clear that CRL-TSC decreases the total fuel consumption by 20% and the total emission of NO<sub>x</sub> by 22% relative

to actor-critic as well as fuel consumption by 13% and emission of NO<sub>x</sub> by 14% in comparison with Q-learning.

So far CRL-TSCs do not consider the spatial distribution of vehicles along the approaching streets. To provide them with this ability, each street is split into two parts and the number of vehicles per street-segment is used as state variables. In fact, the vehi-

**Table 6.** Comparison between the best CRL-TSC with Q-learning and actor-critic.

Controller	Total fuel (L)	Total CO (kg)	Total HC (kg)	Total NO <sub>x</sub> (kg)
Actor-critic	1131±22	150±3.1	13.3±0.29	2.58±0.05
Q-learning	1042±24	135±3.3	11.7±0.31	2.34±0.06
CRL-TSC with departing links	910±18	114±2.4	9.6±0.21	2.00±0.04
% improvements CRL-TSC vs. actor-critic	19.5	24	27.8	22.5
% improvements CRL-TSC vs. Q-learning	12.7	15.5	17.9	14.5

**Fig. 13.** The learning performance of CRL-TSC for two different state spaces. [Colour online.]**Table 7.** Comparison of the best CRL-TSC and an optimized fixed-time controller.

Controller	Avg. travel time (s/km)	Avg. stop time (s/km)	Avg. fuel (L)
Best CRL-TSC	273	191	910
Fixed-time	323	237	968
% improvement best CRL-TSC vs. fixed-time	15	19	6

cles that are more than 100 m away from the intersections are considered in separate variables from the closer ones. The (near)-optimal number of tilings and tiles are found to be 3. Figure 13 shows the learning performance of the new CRL-TSC and compares it with the former design. It is clear that considering the spatial distribution can slightly improve the performance, although the improvement is not significant.

The best CRL-TSC is benchmarked against an optimized fixed-time controller to verify the performance of the proposed method. In the fixed-time controller, the timing of signals is always fixed no matter how the traffic loads change. As shown in Table 7, the best CRL-TSC results in a lower average travel time, stop time, and fuel consumption. The most notable improvement is average stop time.

## 7. Conclusion

This paper described a continuous state reinforcement learning traffic signal controller (CRL-TSC). CRL-TSC has the ability of generalization which enables it to perform accurately on unseen states. The tile coding, RBFs, and TSFs approaches as different function approximation methods were applied in CRL-TSC to tackle the challenges arising from the continuous state space in the traffic network.

A small traffic network (3×3 grid network) was employed for the thorough evaluation of CRL-TSC and the influence of the parameters. CRL-TSCs were designed in such a way that they can be deployed by basic existing traffic infrastructures in developing countries such as Iran. Owing to the distributed feature of CRL-

TSCs, they are extendable to a traffic network with many intersections. Results show that function approximation is crucial in the performance of CRL-TSCs. The excessive rise and reduction in the number of features may ruin the generalization property of CRL-TSCs. Also, the results indicate that the best CRL-TSC is feasible and effective, compared with standard Q-learning and actor-critic and it greatly reduces the average travel time for vehicles as well as fuel consumption and emissions.

## References

- Abdoos, M., Mozayani, N., and Bazzan, A.L.C. 2011. Traffic light control in non-stationary environments based on multi agent Q-learning. In *Proceedings of the 14th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, Washington, DC, pp. 1580–1585. doi:10.1109/ITSC.2011.6083114.
- Abdoos, M., Mozayani, N., and Bazzan, A.L.C. 2013. Holonic multi-agent system for traffic signals control. *Engineering Applications of Artificial Intelligence*, 26(5–6): 1575–1587. doi:10.1016/j.engappai.2013.01.007.
- Abdoos, M., Mozayani, N., and Bazzan, A.L.C. 2014. Hierarchical control of traffic signals using Q-learning with tile coding. *Applied Intelligence*, 40(2): 201–213. doi:10.1007/s10489-013-0455-3.
- Albus, J.S. 1975. A new approach to manipulator control: the Cerebellar Model Articulation Controller (CMAC). *Journal of Dynamic Systems, Measurement, and Control*, 97: 220–227. doi:10.1115/1.3426922.
- AQCC. 2014. The coefficient of emissions in the warm state for gasoline light duty vehicles of Iran. Air Quality Control Company of Tehran, Municipality Tehran.
- Aslani, M., Mesgari, M.S., and Wiering, M. 2017. Adaptive traffic signal control with actor-critic methods in a real-world traffic network with different traffic disruption events. *Transportation Research Part C: Emerging Technologies*, 85: 732–752. doi:10.1016/j.trc.2017.09.020.
- Aslani, M., Mesgari, M.S., Seipel, S., and Wiering, M. 2018. Developing adaptive traffic signal control by actor-critic and direct exploration methods. *Proceedings of the Institution of Civil Engineers - Transport*. doi:10.1680/jtran.17.00085.
- Baird, L.C. III. 1999. Reinforcement learning through gradient descent. *School of Computer Science, Carnegie Mellon University Pittsburgh, Pittsburgh, PA*.
- Barcelos de Oliveira, L., and Camponogara, E. 2010. Multi-agent model predictive control of signaling split in urban traffic networks. *Transportation Research Part C: Emerging Technologies*, 18(1): 120–139. doi:10.1016/j.trc.2009.04.022.
- Bhatta, B. 2010. Analysis of urban growth and sprawl from remote sensing data. Springer, Verlag Berlin Heidelberg.
- Bi, Y., Srinivasan, D., Lu, X., Sun, Z., and Zeng, W. 2014. Type-2 fuzzy multi-intersection traffic signal control with differential evolution optimization. *Expert Systems with Applications*, 41(16): 7338–7349. doi:10.1016/j.eswa.2014.06.022.



- Bishop, C.M. 1995. *Neural networks for pattern recognition*. Clarendon Press, Oxford.
- Boulter, P.C., Barlow, T.J., and McCrae, I.S. 2009. Emission factors 2009: Report 3 – exhaust emission factors for road vehicles in the United Kingdom. TRL, UK.
- Casas, J., Ferrer, J.L., Garcia, D., Perarnau, J., and Torday, A. 2010. Traffic simulation with Aimsun. *In* *Fundamentals of Traffic Simulation*. Edited by J. Barceló. International Series in Operations Research & Management Science. Springer New York, New York, NY, pp. 173–232.
- Chiou, Y.-C., and Huang, Y.-F. 2013. Stepwise genetic fuzzy logic signal control under mixed traffic conditions. *Journal of Advanced Transportation*, **47**(1): 43–60. doi:[10.1002/atr.1205](https://doi.org/10.1002/atr.1205).
- Chowdhury, M.A., and Sadek, A.W. 2003. *Fundamentals of intelligent transportation systems planning*. Artech House, Norwood, MA.
- El-Tantawy, S., Abdulhai, B., and Abdelgawad, H. 2013. Multiagent Reinforcement Learning for Integrated Network of Adaptive Traffic Signal Controllers (MARLIN-ATSC): methodology and large-scale application on downtown Toronto. *IEEE Transactions on Intelligent Transportation Systems*, **14**(3): 1140–1150. doi:[10.1109/TITS.2013.2255286](https://doi.org/10.1109/TITS.2013.2255286).
- Gartner, N.H. 1983. OPAC: A demand-responsive strategy for traffic signal control. *Transportation Research Record: Journal of the Transportation Research Board*, **906**: 75–81.
- Gipps, P.G. 1981. A behavioural car-following model for computer simulation. *Transportation Research Part B: Methodological*, **15**(2): 105–111. doi:[10.1016/0191-2615\(81\)90037-0](https://doi.org/10.1016/0191-2615(81)90037-0).
- Gipps, P.G. 1986. A model for the structure of lane-changing decisions. *Transportation Research Part B: Methodological*, **20**(5): 403–414. doi:[10.1016/0191-2615\(86\)90012-3](https://doi.org/10.1016/0191-2615(86)90012-3).
- Head, K.L., Mirchandani, P.B., and Sheppard, D. 1992. Hierarchical framework for real-time traffic control. *Transportation Research Record*, **1360**: 82–88.
- Henry, J.J., Farges, J.L., and Tufal, J. 1983. The PROLYN real-time traffic algorithm. *IFAC Proceedings Volumes*, **16**(4): 305–310. doi:[10.1016/S1474-6670\(17\)62577-1](https://doi.org/10.1016/S1474-6670(17)62577-1).
- Hunt, P.B., Robertson, D.I., Bretherton, R.D., and Winton, R.I. 1981. SCOOT – a traffic responsive method of coordinating signals, TRL, Crowthorne, UK.
- Jiménez-Palacios, J.L. 1999. Understanding and quantifying motor vehicle emissions with vehicle specific power and tunable infrared laser differential absorption spectrometer remote sensing. Department of Mechanical Engineering, Massachusetts Institute of Technology, MA.
- Kaelbling, L.P., Littman, M.L., and Moore, A.W. 1996. Reinforcement learning: a survey. *Journal of Artificial Intelligence Research*, **4**: 237–285.
- Konda, V.R., and Borkar, V.S. 1999. Actor-critic like learning algorithms for Markov decision processes. *SIAM Journal on Control and Optimization*, **38**(1): 94–123. doi:[10.1137/S036301299731669X](https://doi.org/10.1137/S036301299731669X).
- Mamdani, E.H. 1974. Application of fuzzy algorithms for control of simple dynamic plant. *Proceedings of the Institution of Electrical Engineer*, **121**(12): 1585–1588. doi:[10.1049/piee.1974.0328](https://doi.org/10.1049/piee.1974.0328).
- Medina, J.C., and Benekohal, R.F. 2012. Q-learning and approximate dynamic programming for traffic control: a case study for an oversaturated network. *In* *Proceedings of the 91st Annual Meeting of the Transportation Research Board*, Washington, DC.
- Medina, J.C., Hajbabaie, A., and Benekohal, R.F. 2010. Arterial traffic control using reinforcement learning agents and information from adjacent inter-sections in the state and reward structure. *In* *Proceedings of the 13th International IEEE Conference on Intelligent Transportation Systems (ITSC)*. IEEE, Funchal, pp. 525–530.
- Prashanth, L., and Bhatnagar, S. 2011. Reinforcement learning with function approximation for traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, **12**(2): 412–421. doi:[10.1109/TITS.2010.2091408](https://doi.org/10.1109/TITS.2010.2091408).
- Richter, S. 2006. Learning road traffic control: towards practical traffic control using policy gradients. Albert-Ludwigs-Universität Freiburg, Fakultät für Angewandte Wissenschaften, Institut für Informatik, Germany.
- Rodrigue, J.-P., Comtois, C., and Slack, B. 2017. *The Geography of transport system*. Routledge, New York.
- Salkham, A., Cunningham, R., Garg, A., and Cahill, V. 2008. A collaborative reinforcement learning approach to urban traffic control optimization. *In* *Web Intelligence and Intelligent Agent Technology*, 2008. Proceedings of the WI-IAT '08. IEEE/WIC/ACM International Conference, 9–12 December 2008, IEEE, pp. 560–566.
- Samarasinghe, S. 2016. *Neural networks for applied sciences and engineering: from fundamentals to complex pattern recognition*. Auerbach Publications, New York.
- Schwartz, H.M. 2014. *Multi-agent machine learning: a reinforcement approach*. Wiley, NJ.
- Sims, A.G., and Dobinson, K.W. 1980. The Sydney coordinated adaptive traffic (SCAT) system philosophy and benefits. *IEEE Transactions on Vehicular Technology*, **29**(2): 130–137. doi:[10.1109/T-VT.1980.23833](https://doi.org/10.1109/T-VT.1980.23833).
- Srinivasan, D., Choy, M.C., and Cheu, R.L. 2006. Neural networks for real-time traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, **7**(3): 261–272. doi:[10.1109/TITS.2006.874716](https://doi.org/10.1109/TITS.2006.874716).
- Steingröver, M., Schouten, R., Peelen, S., Nijhuis, E., and Bakker, B. 2005. Reinforcement learning of traffic light controllers adapting to traffic congestion. *In* *Proceedings of the Belgium-Netherlands Artificial Intelligence Conference, BNAIC'05*, Brussels, pp. 216–223.
- Sutton, R.S. 1988. Learning to predict by the methods of temporal differences. *Machine Learning*, **3**(1): 9–44. doi:[10.1023/A:1022633531479](https://doi.org/10.1023/A:1022633531479).
- Sutton, R.S., and Barto, A.G. 1998. *Reinforcement learning: an introduction*. MIT Press, Cambridge, MA.
- Szepesvári, C. 2010. *Algorithms for reinforcement learning*. Morgan & Claypool Publishers.
- Takagi, T., and Sugeno, M. 1985. Fuzzy identification of systems and its applications to modeling and control. *IEEE Transactions on Systems, Man, and Cybernetics*, **SMC-15**(1): 116–132. doi:[10.1109/TSMC.1985.6313399](https://doi.org/10.1109/TSMC.1985.6313399).
- Thorndike, E.L. 1911. *Animal intelligence; experimental studies*. The Macmillan Company, New York.
- TRL. 1999. *Methodology for calculating transport emissions and energy consumption*. Transport Research Laboratory, Crowthorne, UK.
- van Otterlo, M., and Wiering, M. 2012. Reinforcement learning and Markov decision processes. *In* *Reinforcement Learning: State-of-the-Art*. Edited by M. Wiering and M. van Otterlo. Springer, Berlin, Heidelberg, pp. 3–42.
- Watkins, C., and Dayan, P. 1992. Q-learning. *Machine Learning*, **8**(3): 279–292. doi:[10.1007/BF00992698](https://doi.org/10.1007/BF00992698).
- Wiering, M. 2000. *Multi-agent reinforcement learning for traffic light control*. Stanford, CA, pp. 1151–1158.
- Zadeh, L.A. 1965. Fuzzy sets. *Information and Control*, **8**(3): 338–353. doi:[10.1016/S0019-9958\(65\)90241-X](https://doi.org/10.1016/S0019-9958(65)90241-X).