

# *A reinforcement learning-based assignment scheme for EVs to charging stations*

Conference or Workshop Item

Accepted Version

Aljaidi, M., Aslam, N., Chen, X. ORCID: <https://orcid.org/0000-0001-9267-355X>, Kaiwartya, O., Al-Gumaei, Y. A. and Khalid, M. (2022) A reinforcement learning-based assignment scheme for EVs to charging stations. In: 2022 IEEE 95th Vehicular Technology Conference: (VTC2022-Spring), 19-22 June 2022, Helsinki, Finland. doi: <https://doi.org/10.1109/VTC2022-Spring54318.2022.9860535> Available at <https://centaur.reading.ac.uk/116499/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1109/VTC2022-Spring54318.2022.9860535>

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

**CentAUR**

Central Archive at the University of Reading

Reading's research outputs online

# A Reinforcement Learning-based Assignment Scheme for EVs to Charging Stations

Mohammad Aljaidi\*, Nauman Aslam\*, Xiaomin Chen\*, Omprakash Kaiwartya<sup>†</sup>,  
Yousef Ali Al-Gumaei\*, Muhammad Khalid<sup>‡</sup>

\*Department of Computer and Information Sciences, Northumbria University, Newcastle upon Tyne, UK

<sup>†</sup>School of Science and Technology, Nottingham Trent University, UK

<sup>‡</sup>Department of Computer Science and Technology, University of Hull, UK

Email: \*{mohammad.jaidi,nauman.aslam,xiaomin.chen,yousef.al-gumaei}@northumbria.ac.uk,

<sup>†</sup>omprakash.kaiwartya@ntu.ac.uk, <sup>‡</sup>m.khalid@hull.ac.uk

**Abstract**—Due to recent developments in electric mobility, public charging infrastructure will be essential for modern transportation systems. As the number of electric vehicles (EVs) increases, the public charging infrastructure needs to adopt efficient charging practices. A key challenge is the assignment of EVs to charging stations in an energy efficient manner. In this paper, a Reinforcement Learning (RL)-based EV Assignment Scheme (RL-EVAS) is proposed to solve the problem of assigning EV to the optimal charging station in urban environments, aiming at minimizing the total cost of charging EVs and reducing the overload on Electrical Grids (EGs). Travelling cost that is resulted from the movement of EV to CS, and the charging cost at CS are considered. Moreover, the EV's Battery State of Charge (SoC) is taken into account in the proposed scheme. The proposed RL-EVAS approach will approximate the solution by finding an optimal policy function in the sense of maximizing the expected value of the total reward over all successive steps using Q-learning algorithm, based on the Temporal Difference (TD) learning and Bellman expectation equation. Finally, the numerous simulation results illustrate that the proposed scheme can significantly reduce the total energy cost of EVs compared to various case studies and greedy algorithm, and also demonstrate its behavioural adaptation to any environmental conditions.

**Index Terms**—Electric vehicle assignment, charging station, Q-learning, temporal difference, Bellman expectation equation, energy consumption, energy cost, electrical grids.

## I. INTRODUCTION

**B**EING one of the fastest growing sources of energy demand and carbon dioxide (CO<sub>2</sub>) emissions, the transportation sector is under great pressure to be decarbonized through deploying EVs. At the same time, various parameters, such as the technological innovation in battery efficiency and electric drivetrain, led to dramatically increase the EVs penetration in recent years in urban environments [1]. Charging EVs at home is an alternative solution for the EV users, however, it takes too much time (6 to 8 hours) for each charging process. Therefore, establishing a conveniently available fast-charging stations is essential to increase the satisfaction of EV users, and alleviating peak charging loads. Using fast-charging stations can charge EV batteries at least 12 times faster [2].

The adoption rate of EVs mainly depends on the presence of wide range of charging stations in metropolitan areas to help EV users to charge their vehicles during their journeys [3]. Furthermore, the selection of best fast-charging stations for EVs is an important factor that affects

not only the spread of EVs but also minimize the total energy consumption of EVs to reach charging stations and increase the sustainability of transportation, knowing that the energy demand in the transportation sector will increase by 54% until 2035 [4].

In recent years, more attention has been paid to suggesting learning-based techniques for finding the optimal charging stations for EVs in metropolitan environments. Several algorithms and technique have been utilized to solve this problem. Machine learning is a sub-field of artificial intelligence (AI) that provides systems with the ability to automatically learn from experience without being explicitly programmed and finally solve the most complex problems [5]–[7]. RL technique is one of the machine learning categories that has been utilized recently in several studies to solve the problem of assigning EV to optimal CS. In [8], an optimal charging station selection algorithm based model-free RL technique has been proposed, to minimize the total travel time of EVs charging demands from origin to CS using the selected optimal route, taking into account unknown future demands and dynamically changing traffic conditions. The charging station selection was formulated as a Markov decision process (MDP) model with unknown transition probability.

A congestion control in CSs allocation with Q-learning has been introduced in [9]. In this paper, the authors considered the travel time and the queuing time in CSs in their model, to build joint-resource congestion game which utilises the interaction between the EVs and resources. Q-learning algorithm was applied to the model in order to solved the problem. In [10], a deep RL (DRL) based EV charging navigation has been proposed to adaptively learn the optimal technique for EV charging navigation without any prior knowledge of uncertainties. Feature states have been extracted from stochastic data. The authors investigated the DRL-based charging navigation from a single perspective of EV owner, considering randomness in smart grid and intelligent transportation data. To the best of our knowledge, the existing literature did not take into account the variation in the charging rate between CSs, in order to assign EVs to the optimal CSs. Our proposed scheme is unique in that, as we take into account the differences in charging rate among all available CSs, the charging cost at CS, the travelling cost of EV to reach CS, and also the capacity of the CS. We argue that all of these metrics have

a significant impact on the decision of assignment EVs to CSs.

The main contributions of the present paper are summarized as follows:

- A RL-scheme for assignment of EVs to the optimal CSs in metropolitan environments is proposed in this paper. The proposed scheme considers the energy consumption cost that is resulted from the movement of EV towards CS (travelling cost), and the total expected cost to fully charge EV at CS (total energy cost). EV's battery SoC in the process of finding the optimal CS is also taken into account.
- Q-learning algorithm has been utilized to solve this problem based on maximizing the cumulative reward of the EV during learning process by reducing the total cost of charging EVs.
- Reduce the load on the overwhelmed EGs, by assuming different rewards for the available CSs in the study area. The reward at each CS is determined based on the electricity price that offered by electrical grids (EGs) to CSs, these prices vary according to the load at EGs. Higher load EGs will offer a higher price to the CSs associated with them as a form of punishment. which in turn forces EV users to look for other options to charge their vehicles.

The rest of the paper is organized as follows. EV charging station assignment problem formulation and optimization model are presented in Section II. Then, a RL based approach is proposed in Section III. Section IV shows the numerical results of our proposed approach. Finally, Section V draws a conclusion.

## II. EV CHARGING STATION ASSIGNMENT PROBLEM

This paper proposes a new RL scheme to assign EVs to the CSs in metropolitan environment based on minimizing the total expected cost of charging EVs. The notations used in this paper are listed in Nomenclature list. In addition, parameters and variables are explained in the where they are first used.

### NOMENCLATURE

#### Sets and Index

$\mathcal{A}$	Set of actions.
$\mathcal{M}$	Set of CSs.
$\mathcal{S}$	Set of states.
$j$	index of CSs.

#### Parameters and Variables

$\alpha, \gamma$	The learning rate and discount factor, respectively.
$\mu_{ev}, \mu_j$	The longitude of $EV$ and $CS_j$ , respectively.
$\Phi$	A threshold value that restricts the maximum amount of energy consumption that EV can consume to reach CS
$\pi^*(s)$	The optimal policy.
$\psi_{ev}, \psi_j$	The latitude of $EV$ and $CS_j$ , respectively.

$\varphi_j$	The charging rate at $CS_j$ .
$\vartheta_{ev,j}$	The energy consumption cost that is resulted from the movement of EV towards CS.
$\xi_{ev,j}$	The total expected energy to fully charge EV at CS.
$\zeta_{ev}$	The EV battery capacity.
$a, a'$	The action that the agent takes in the current state, and target state, respectively.
$C_{ev,j}$	The overall cost of charging EV.
$d_{ev,j}$	The total distance that EV travels to reach CS.
$E_{ev,j}$	The overall energy.
$M$	The number of CSs in the study area.
$Q^*(s, a)$	The optimal state-action value function.
$Q(s, a)$	A state-action value function, i.e., Q-value function.
$r$	The immediate reward.
$s, s'$	The current state and target state, respectively.
$x_{ev,j}$	A binary decision variable shows that the $EV$ selects $CS_j$ for charging.

#### A. Problem Formulation

The problem has been formulated as shown in the following sections:

1)  $EV$ : The EV is represented by a single agent that moves in the environment trying to find the optimal CS considering the possible actions and reward at each state. The  $EV$  has one attribute; ( $p^{ev}$ ), where  $p^{ev}$  is the position (coordinates) of  $EV$ .

2)  $CS$ s: Define the CS set as  $\mathcal{M} = \{1, \dots, u, \dots, M\}$ . The cardinality of  $\mathcal{M}$  is  $M$ , i.e., there are  $M$  CSs in the investigated area.  $CS_j$  in  $\mathcal{M}$  has two attributes;  $b^u, r^u$ , where  $b^u$ , and  $r^u$  are the position and reward of the  $CS$ . The reward at  $CS_j$ , depends on the charging rate of  $CS_j$ .

3)  $Cost$ -based  $EV$  assignment: The proposed strategy uses a RL technique to assign EV to the best CS based on minimizing the total cost of charging EVs, i.e.,  $C_{ev,j}$ . To calculate the total expected cost of charging EV, two factors should be investigated; the energy consumption cost that is resulted from the travelling of EV towards CS, i.e.,  $\vartheta_{ev,j}$ , and the total expected energy to fully charge EV at CS, i.e.,  $\xi_{ev,j}$ , considering the charging rate at each CS, which is usually determined based on the electricity price offered by EGs to the CS owners. Fig. 1 depicts the interaction between the three entities of the system, where Entity 1 represents EV, Entity 2 is the CS, and Entity 3 is the EG. The electricity price that offered by EGs to CSs is different due to the load and location of each EG. CSs connected to the same EG have the same electricity price, and therefore the same charging rate. The overall energy, i.e.,  $E_{ev,j}$ , can be calculated as follows:

$$E_{ev,j} = \vartheta_{ev,j} + \xi_{ev,j} \quad (1)$$

To calculate  $\vartheta_{ev,j}$ , we need to calculate the amount of energy that the EV consumes per km to reach CS, i.e.,  $\delta_{ev}$  [11], [12], and also calculate the total distance that the EV travels towards CS, i.e.,  $d_{ev,j}$ . In this work, we assume that  $\delta_{ev}$  is 0.16 kWh/km, with an average speed 40 km/h [13], [14] The following illustrates how  $d_{ev,j}$ ,  $\vartheta_{ev,j}$ , and  $\xi_{ev,j}$  are calculated:

$$d_{ev,j} = \sqrt{(\psi_{ev} - \psi_j)^2 + (\mu_{ev} - \mu_j)^2} \quad (2)$$

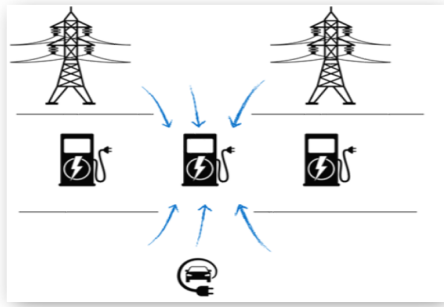


Fig. 1. The charging station interaction system

where  $\psi_{ev}$ ,  $\psi_j$  are the latitude of  $EV$  and  $CS_j$ , and  $\mu_{ev}$ ,  $\mu_j$  are longitude, respectively.

$$\vartheta_{ev,j} = \delta_{ev} * d_{ev,j} \quad (3)$$

where  $\delta_{ev}$  is the amount of energy that the  $EV$  consumes per km to reach  $CS$

$$\xi_{ev,j} = (\zeta_{ev} - SoC) + \vartheta_{ev,j} \quad (4)$$

where  $\zeta_{ev}$  denotes the capacity of the battery, and  $SoC$  is the battery state of charge.

The total cost of charging  $EV$  can be calculated as follows:

$$C_{ev,j} = E_{ev,j} \times \varphi_j \quad (5)$$

where  $\varphi_j$  represents the charging rate at  $CS_j$ .

### B. Optimization Problem

The corresponding optimization problem of our proposed approach can be written as:

$$\min_X \sum_{j=1}^M C_{ev,j} x_{ev,j} \quad (6)$$

$$s.t. \quad \sum_{j=1}^M x_{ev,j} = 1 \quad (7)$$

$$x_{ev,j} \in \{0, 1\}, \quad \forall j \quad (8)$$

$$\vartheta_{ev,j} < \Phi, \quad \forall j \quad (9)$$

where the variable  $x_{ev,j}$  is used as a decision variable with binary values  $\{0,1\}$  as shown in (8), to indicate whether the charging station  $j$  is selected by the  $EV$  or not,  $x_{ev,j}$  is equal to 1 if the  $j$  is selected, otherwise it is equal to 0. Constraint (7) restricts that only one charging station is selected as destination. Constraint (9) shows that the total amount of energy consumption that  $EV$  needs to reach  $CS$  should not reach to a certain threshold in order to maintain the  $EV$  battery's  $SoC$ .

## III. REINFORCEMENT LEARNING APPROACH

Typically, RL techniques can be processed under two categories: off-policy and on-policy [15]. In particular, an off-policy learning method (e.g., Q-learning) earns an optimal target policy independent of the behavior policy used during exploration process as long as the different states are explored enough times. Whereas on-policy learning method finds the optimal policy taking into account the actual actions taken over the exploration process, which means that the target policy is the same as the behavior policy used in exploration process. Q-Learning technique will be used in this work to address the problem of assignment  $EV$  to  $CS$ .

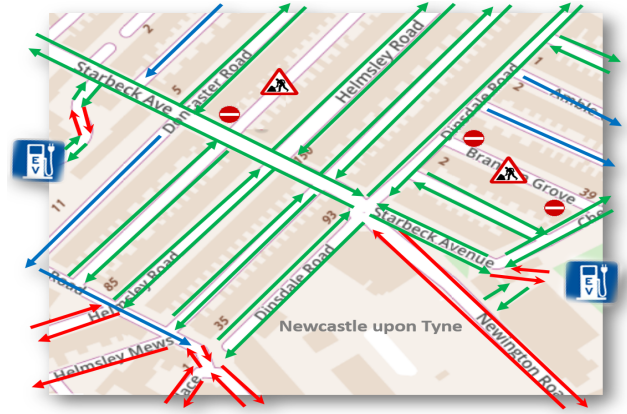


Fig. 2. Sample from the urban area of Newcastle upon Tyne city, UK

### A. Q-Learning-based EV Charging Station Assignment

In this section, a RL technique is employed to solve our optimization problem (6-9), based on Q-Learning Algorithm technique. Q-learning, an incremental technique for dynamic programming, is appropriate for solving such kind of problems. Q-learning is an agent-based technique in which the AI agent interacts with its environment and adapts its actions based on rewards or penalties received in response to its actions [16]. Mainly, there are three basic elements in the algorithm: environment, state, and action. We will introduce the algorithm after setting the elements.

1) *Environment, State, and Action Set*: The environment is an essential element in Q-learning, in which the AI agent selects its actions according to corresponding rewards. In our scenario, the environment should involve roads between  $EV$ s and those available  $CS$ s.

Fig. 2 shows part of the Newcastle upon Tyne, UK. The directions (actions) the  $EV$  is allowed to take in the study area, are represented by three colors of arrows, where the green and blue arrows show that the  $EV$  can move in only four directions; South, North, East, and West, the difference between them is that the green arrows denote that the  $EV$  can move in the streets in both directions, but the blue arrows shows that these streets can be used only in one direction. While, the red arrows indicate that the  $EV$  can move in the other four directions; South-East, South-West, North-East, and North-West and also in both directions. Road works signs indicate that these streets are closed and can't be used by vehicles, which means that the  $EV$  user needs to find other streets to reach  $CS$ . The  $CS$ s are distributed in fixed locations in the study area as shown in Fig. 2.

The state set in our scenario can be denoted by  $\mathcal{S}$ , and defined as following:

$$\mathcal{S} = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\} \quad (10)$$

where the dimension of this grid is  $(n \times n)$ , Each cell in this grid represents a state of the  $EV$ . For each cell, it has an incremental reward once it is on the path chosen by the  $EV$ . In this way, the system records the accumulated reward of the  $EV$  based on the actions taken in the environment.

The action set, i.e.,  $\mathcal{A}$ , of the  $EV$  in the grid world denotes the way how the  $EV$  can moves to change its

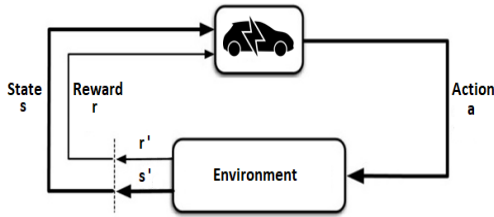


Fig. 3. Example of a single movement of EV

state. In our scene, the directions that the EV is allowed to use are included in the following set:

$$\mathcal{A} = \{South, North, East, West, South - East, South - West, North - East, and North - West\} \quad (11)$$

The energy consumption cost of EV movement towards CS, i.e.,  $\vartheta_{ev,j}$ , is mainly dependent on the distance between the locations of EV and CS. To minimize this cost, we need to reduce the covered distance that the EV needs to travel to reach the CS. To achieve this, the EV earns punishment (penalty) for every hop in the grid. The rewards associated with CSs depends on the charging rates at these CSs. Accordingly, CS with higher charging rate, the lower its reward. The reward increases as the charging rate decreases. Fig. 3 shows an Illustrative example of a single movement of EV in the environment.

2) *Q-Learning Algorithm*: The Q-learning algorithm is based on an action-value function or Q-function. The Q-learning algorithm has two input parameters; state and action. The objective function of this work is minimizing the total energy cost of EV. Basically, the optimal strategy is determined by recursively updating action-value function based on the Bellman equation, with the TD learning algorithm technique as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \quad (12)$$

where  $Q(s, a)$  represents the value of the state-action function, i.e, Q-value function  $(s, a)$ ,  $\alpha$  and  $\gamma$  denotes the learning rate and discount factor between 0 and 1, and  $r$  is the immediate reward value received as the result of taking action  $a$  in state  $s$ . In the environment, there are many different Q-value functions according to the different actions and policies that can be used in the learning process. The optimal Q-value function is the value which yields maximum Q-value compared to all other Q-values that have been acquired during the learning process. So, mathematically the optimal Q-value function, i.e, state-value function, can be expressed as:

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) \quad (13)$$

where  $Q^*(s, a)$  denotes the optimal value-action function.

The purpose of the EV charging navigation process, is to find the optimal policy  $\pi^*$  over all feasible policies that EV can select during the learning process, which minimizes the cost or maximizes the reward. Therefore, Once we have  $Q^*(s, a)$ , EV can act optimally based on the optimal strategy as shown in the following greedy strategy:

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a) \quad (14)$$

where  $\pi^*(s)$  is the optimal policy that EV can act in the environment, which achieves the optimal value-action function.

Q-learning algorithm uses  $\epsilon$ -greedy policy for the action selection step, which is also called behavior policy, to ensure a high level of balance between exploration and exploitation, i.e, exploration-exploitation trade-off, as well

as to improve the learning level of the agent during the direct interaction with the environment. A simple strategy that has been proposed to deal with this problem is the  $\epsilon$ -greedy (with  $0 \leq \epsilon < 1$ ), with greater corresponding to greater probability of exploration. the value of  $\epsilon$  has a significant impact on the performance and complexity of the Q-learning algorithm. Details can be seen in Algorithm 1 as follows:

---

#### Algorithm 1 Training Process of Q-Learning algorithm

---

**Input:**  $\alpha, \gamma, \epsilon, Q$  (terminal-state),  $s \in S, a \in \Lambda, N, M$

**Output :**  $X^{opt}$ , Optimal  $Q^*(s, a)$ , and  $\pi^*(s)$  for EV

---

**Initialization:**

- 1: Initialize  $Q(s, a)$  arbitrary
  - 2: Initialize fixed CS locations
  - 3:  $K$ =maximum number of episodes
  - 4: Initialize random  $s$  for the EV
  - 5: **for**  $eps = 1$  to  $K$  **do**
  - 6:   Select  $a \in \Lambda$  for EV using  $\epsilon$ -greedy policy
  - 7:   Execute the action  $a$
  - 8:   Receive immediate reward  $r$
  - 9:   Observe the new state  $s'$
  - 10:   Select  $a'$  in  $s'$  for the EV using Eq. (12)
  - 11:   Update  $Q(s, a)$  value in Q-table
  - 12:    $s \leftarrow s'$
  - 13:   **if** ( $s$  is terminal or  $\vartheta_{ev,j} \geq \Phi$ ) **then**
  - 14:     Start new episode
  - 15:   **else**
  - 16:     Select  $a \in \Lambda$
  - 17:   **end if**
  - 18: **end for**
  - 19: Return  $Q^*$ -values, and  $\pi^*$  for the EV
  - 20: Return  $X^{opt}$  considering Eqs. (6)-(9)
- 

## IV. EXPERIMENTS

In this section, a comparison between the proposed approach (RL-EVAS) and greedy strategy, and also different case studies are carried out within the proposed environment to demonstrate the effectiveness and feasibility of the proposed approach. In Section IV-A, the details of experimental setup are presented. The training process and simulation results are discussed in Section IV.

### A. Experimental Setup

The performance of the proposed approach is demonstrated within an area  $25 \times 25$  grid map. The agent (EV) earns -1 as penalty for each movement in the environment. The rewards that the EV earns when reaching CSs varies depending on the charging rate of each CS which mainly depends on the electricity rates as mentioned earlier. Barriers have been placed on some of the roads that EV takes in the directions leading to CSs, which in turn force the EV to search for another available roads. In this work, we assume that the rewards of CSs are two values 60 and 80, depending on the charging rate at CS. As mentioned earlier, The reward increases as the charging rate decreases. An inverse relationship between the charging rate and the given reward associated with each CS. Each episode terminates, if the EV reached to the CS or reached the threshold value of the travelling energy consumption ( $\Phi$ ). All the following experimental results have been performed by Python 3.10 on Windows 2010

TABLE I. RL-EVAS parameters

Parameter	Value
Environment	25×25 grids
$M$	4
Penalty	-1
Rewards	60, 80
$\zeta_{ev}$	62 kWh
$\varphi$	\$0.15, \$0.35
SoC	$\zeta_{ev} \times 60\%$ kWh
$\alpha$	0.1
$\epsilon$	0.1
$\gamma$	0.6
$\delta_{ev}$	0.16 kWh/km
$\Phi$	1.6 kWh

Pro 64bits, V.20H2, Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz (8 CPUs), 1.80 GHz. The simulation parameters related to the proposed scheme are presented in Table I.

### B. Results

The performance of RL-EVAS is evaluated with respect to two criteria:

- The maximum cumulative reward of the value function of the learned policy, reflecting the proposed objective function of this approach.
- The total energy cost of EV, considering the cost that resulting from the movement of EV towards CS, and the cost of charging the EV at CS.

To this end, comparisons between RL-EVAS and different case studies, and also between the RL-EVAS and greedy strategy will be conducted in the proposed scheme.

1) *Case Studies*: The following proposed case studies demonstrate the feasibility and effectiveness of the proposed approach.

- **Case A**: With reduced the number of actions. In this case, we assume that the number of actions that the EV needs to perform to interact with the environment at each state is just 4, rather than 8 actions for RL-EVAS, as shown in Eq. (11). As shown in Fig. 2, the green and blue arrows represent the possible actions that the EV can perform in Case A. While the possible actions that the EV can select in RL-EVAS are the green, blue and red arrows which give the EV to perform 8 actions. Finally, the other parameters remain the same as the RL-EVAS.
- **Case B**: With the increase in the number of obstacles standing in the way of the EV towards CSs. In this case, we assume that the number of barriers is increased by 25%, 50%, and 75% compared to the RL-EVAS with the same number of episodes. While the other parameters remain the same as the RL-EVAS.

Fig. 4 and Table II show the comparison results between RL-EVAS and case A. It can be seen that the distance, travelling energy consumption, and travel time are less in RL-EVAS compared to case A, this leads to minimize the total amount of energy needed to charge the EV. As a result, the total energy cost is minimized as shown in Fig. 4, this is due to the assumption that the possible actions that EV can select in RL-EVAS is more compared to case A. Another observation in Fig. 4, is that the cumulative reward for RL-EVAS is more compared to case

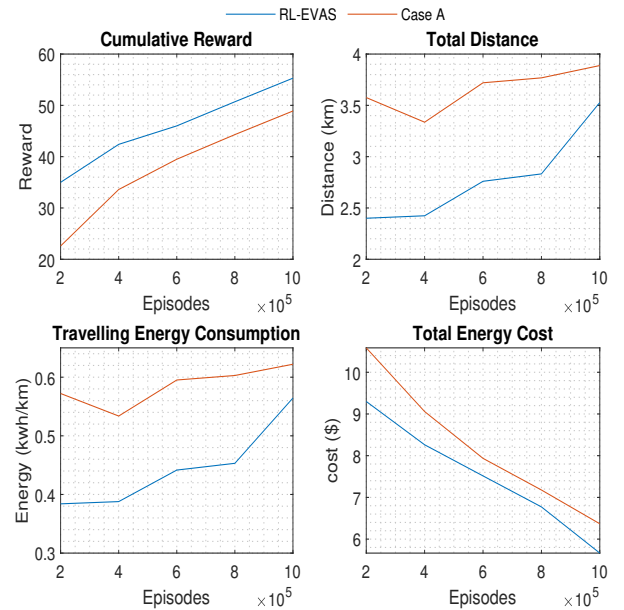


Fig. 4. Comparison between RL-EVAS and case A

A, and the reason behind this is that the total number of timestep (the hops that EV needs to reach CS) is less compared to the case A as shown in Table II. Each timestep is considered as an additional penalty for the EV. Accordingly, the cumulative penalty will finally be deducted from the reward that EV receives when it reaches CS. Based on the foregoing, it is noticeable that the cumulative reward is inversely related to the timestep.

Fig. 5 and Tables III - V, show the comparison results between RL-EVAS and case B. In this case, we assume three different scenarios. In the first scenario, we assume that the number of obstacles is increased by 25%, then in the second scenario we assume that the number is increased by 50%, while in the last scenario, we assume that the number is increased by 75%. The reason behind this assumption is that street conditions are not fixed and can change for several reasons, including, but not limited to, the works that may occur in the study area as shown in Fig. 2.

It is easy to notice that that the total timestep, total distance, travelling energy consumption, and the total energy that is required to fully charge EV are less in RL-EVAS compared to the all scenarios in case B, and it can be observed that all of these parameters are increased according to the proposed scenario in case B as shown in Fig. 5, the reason for this is the assumption of increasing the percentage of streets that the EV cannot use in the study area to access CS. To overcome this challenge, the EV will try to find other possible routes to reach CS, which in turn increases the number of hops (timestep) that the EV must take to reach CS. Consequently, it increases the distance between the location of EV and the location of the CS at any charging decision point, the energy consumed by EV en route to the CS, and also the total energy that is required to fully charge EV. It is also noticeable that the timestep, distance, travelling energy consumption and total energy decrease when the number of episodes is higher as shown in Fig. 5. This is due to the fact that the performance of the EV in the study area improves

TABLE II. Comparison between RL-EVAS and Case A in terms of Timestep, Total Energy and Travel Time

Episodes	RL-EVAS			Case A		
	Timestep	Total Energy (kwh/km)	Travel Time (m)	Timestep	Total Energy (kwh/km)	Travel Time (m)
$2 \times 10^5$	5	37.584	3.6	7.45	37.77216	5.364
$4 \times 10^5$	5.05	37.58784	3.636	6.95	37.73376	5.004
$6 \times 10^5$	5.75	37.6416	4.14	7.75	37.79519999	5.57985
$8 \times 10^5$	5.9	37.65312	4.248	7.85	37.80288	5.652
$10 \times 10^5$	7.35	37.76448	5.292	8.1	37.82208	5.831985
<b>Average</b>	<b>5.81</b>	<b>37.646208</b>	<b>4.9516</b>	<b>7.622</b>	<b>37.785215998</b>	<b>5.486367</b>

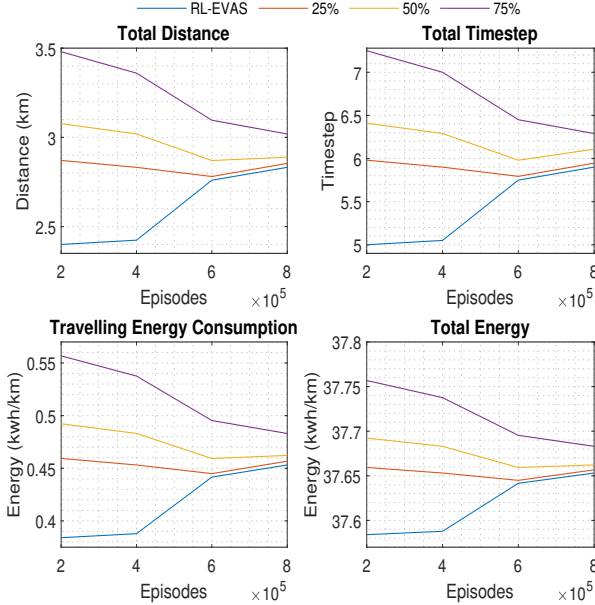


Fig. 5. Comparison between RL-EVAS and case B

TABLE III. Comparison between RL-EVAS and Case B when the number of obstacles is increased by 25%

Episodes	RL-EVAS		Case B	
	Cumulative Reward	Total Energy Cost (\$)	Cumulative Reward	Total Energy Cost (\$)
$2 \times 10^5$	35	9.3	34.27	9.8864
$4 \times 10^5$	42.4	8.2613759	41.8	8.3891
$6 \times 10^5$	46	7.515455	43.2	7.67741
$8 \times 10^5$	50.7	6.772416	48.9	6.78912
<b>Average</b>	<b>43.525</b>	<b>7.9623117</b>	<b>42.0425</b>	<b>8.1855075</b>

with the increase in the number of episodes, since the EV will have a higher chance of finding the optimal CS at any charging decision point, even though the number of available streets decreases.

As shown in Tables III - V, the cumulative reward that has been achieved in RL-EVAS is higher compared to all scenarios of case B, which means that the EV in RL-EVAS has selected CSs with higher reward rather than selecting CSs with lower rewards as the EV performed in case B. As mentioned in Section III-B, the charging rate at CS decreases as the given reward increases, and this is the reason why the total energy cost for charging EV at selected CSs in RL-EVAS is less compared to all scenarios in case B as shown in the above-mentioned tables.

Tables VI, shows the comparison between the RL-EVAS and case B, in terms of the travel time that the EV requires to reach CS. It is seen that the travel time for RL-EVAS is less compared to the all scenarios in case B, and the reason behind this is that the number of the streets that

TABLE IV. Comparison between RL-EVAS and Case B when the number of obstacles is increased by 50%

Episodes	RL-EVAS		Case B	
	Cumulative Reward	Total Energy Cost (\$)	Cumulative Reward	Total Energy Cost (\$)
$2 \times 10^5$	35	9.3	33.67	10.0879
$4 \times 10^5$	42.4	8.2613759	38.9	8.4662
$6 \times 10^5$	46	7.515455	44	7.62554
$8 \times 10^5$	50.7	6.772416	47.9	6.82147
<b>Average</b>	<b>43.525</b>	<b>7.9623117</b>	<b>41.1175</b>	<b>8.2502775</b>

TABLE V. Comparison between RL-EVAS and Case B when the number of obstacles is increased by 75%

Episodes	RL-EVAS		Case B	
	Cumulative Reward	Total Energy Cost (\$)	Cumulative Reward	Total Energy Cost (\$)
$2 \times 10^5$	35	9.3	32.87	10.32453
$4 \times 10^5$	42.4	8.2613759	37	8.7986
$6 \times 10^5$	46	7.515455	42	7.70554
$8 \times 10^5$	50.7	6.772416	46.3	6.92374
<b>Average</b>	<b>43.525</b>	<b>7.9623117</b>	<b>39.5425</b>	<b>8.4381025</b>

the EV cannot use to find CS with a lower charging rate is increased compared to RL-EVAS, which means that the EV needs to travel longer distance, which in turn increases the travel time to reach CS.

2) *A Comparison between RL-EVAS and Greedy strategy:* Reinforcement learning algorithms have been used in many previous studies to solve the problem of charging EVs. In [17], the authors have proposed a RL approach to find the optimal scheduling and pricing for EV CSs, and they demonstrated the proposed approach by comparing it with a greedy heuristic that assigns EVs to nearest CSs.

TABLE VI. Comparison between RL-EVAS and Case B in terms of the travel time (m)

Episodes	RL-EVAS	Case B		
		25%	50%	75%
$2 \times 10^5$	3.6	4.30618	4.6152	5.22
$4 \times 10^5$	3.636	4.248	4.5288	5.04
$6 \times 10^5$	4.14	4.17096	4.3056	4.644
$8 \times 10^5$	4.248	4.28134	4.33318	4.52837
<b>Average</b>	<b>3.906</b>	<b>4.25162</b>	<b>4.4457</b>	<b>4.858075</b>

TABLE VII. Comparison between RL-EVAS and the Greedy Strategy in terms of Total Energy and Total Energy Cost

Episodes	RL-EVAS		Greedy Strategy	
	Total Energy kwh/km	Total Energy Cost (\$)	Total Energy kwh/km	Total Energy Cost (\$)
$2 \times 10^5$	37.584	9.3	37.56096	10.51706
$4 \times 10^5$	37.58784	8.26137	37.55788	9.78
$6 \times 10^5$	37.6416	7.51545	37.53484	10.6891
$8 \times 10^5$	37.65312	6.77241	37.49184	8.6871
$10 \times 10^5$	37.76448	5.66467	37.47801	9.0247
<b>Average</b>	<b>37.64621</b>	<b>7.50278</b>	<b>37.52471</b>	<b>9.73959</b>



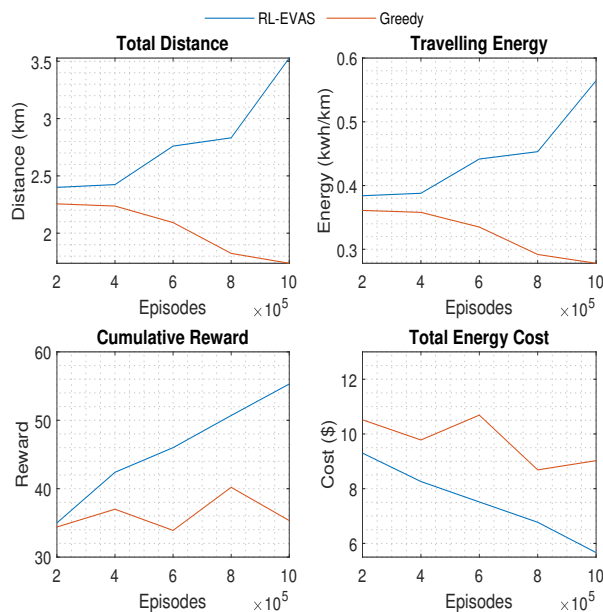


Fig. 6. Comparison between RL-EVAS and Greedy strategy

In this work, in addition to the proposed case studies that have been assumed before, we also compare RL-EVAS with greedy strategy to demonstrate our proposed approach, the comparison between the two strategies is done based on the total distance, travelling energy consumption, total energy that is required to fully charge EV, cumulative reward, and total energy cost. As mentioned before, the performance of RL-EVAS is evaluated with the respect into two criteria; the cumulative reward, and total energy cost. Fig. 6 and Table VII show the results of the comparison between RL-EVAS and the greedy approach.

Although the greedy algorithm has achieved comparable results in terms of total distance, travelling energy consumption, total energy that is required to fully charge EV as shown in Fig. 6 and Table VII. However, RL-EVAS was able to achieve the desired results, in terms of maximizing the cumulative reward and minimizing the total energy cost. The reason behind this, is that the greedy algorithm has selected the CS based only on the distance, which in turn reduced the energy consumption to reach CS, and the total energy required to fully charge EV. However, the greedy algorithm did not take into account the variance in the rewards that have been associated with each individual CS, and also the charging rate at each CS, which in turn led to reduce the cumulative reward and increase the total energy cost. On the contrary, RL-EVAS has taken into account the two parameters, thus achieved the goal of the system, which is maximizing the cumulative reward and minimizing the total energy cost.

## V. CONCLUSION

This paper has proposed a RL-based assignment scheme for EVs to CSs in urban environment. Several parameters were considered, including the distance between the locations of EV and CS at any charging decision point, the EV travelling energy consumption, charging rate at each CS, EV battery's SoC, and the total energy cost of EV. Experimental results demonstrate that the proposed scheme can approximate the solution of finding the optimal policy

in the sense of maximizing the expected value of the total reward, and minimizing the total energy cost of EV using Q-learning algorithm, outperforming all proposed case studies and greedy approach. Our future works will focus on improving the scalability of the proposed scheme for practical applications, and using the Deep Q-Network (DQN) algorithm to solve this problem.

## REFERENCES

- [1] T. Zhang, W. Chen, Z. Han, and Z. Cao, "Charging scheduling of electric vehicles with local renewable energy under uncertain electric vehicle arrival and grid power price," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 6, pp. 2600–2612, 2013.
- [2] Y. Xiong, J. Gan, B. An, C. Miao, and A. L. C. Bazzan, "Optimal electric vehicle fast charging station placement based on game theoretical framework," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 8, pp. 2493–2504, 2017.
- [3] M. Aljaidi, N. Aslam, and O. Kaiwartya, "Optimal Placement and Capacity of Electric Vehicle Charging Stations in Urban Areas: Survey and Open Challenges," in *2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT)*. IEEE, 2019, pp. 238–243.
- [4] M. R. Mozafar, M. H. Moradi, and M. H. Amini, "A simultaneous approach for optimal allocation of renewable energy sources and electric vehicle charging stations in smart grids based on improved GA-PSO algorithm," *Sustainable cities and society*, vol. 32, pp. 627–637, 2017.
- [5] S. Vandael, B. Claessens, D. Ernst, T. Holvoet, and G. Deconinck, "Reinforcement learning of heuristic ev fleet charging in a day-ahead electricity market," *IEEE Transactions on Smart Grid*, vol. 6, no. 4, pp. 1795–1805, 2015.
- [6] A. Chiş, J. Lundén, and V. Koivunen, "Reinforcement learning-based plug-in electric vehicle charging with forecasted price," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 5, pp. 3674–3684, 2016.
- [7] F. Ruelens, B. J. Claessens, S. Vandael, B. De Schutter, R. Babuška, and R. Belmans, "Residential demand response of thermostatically controlled loads using batch reinforcement learning," *IEEE Transactions on Smart Grid*, vol. 8, no. 5, pp. 2149–2159, 2016.
- [8] K.-B. Lee, M. A. Ahmed, D.-K. Kang, and Y.-C. Kim, "Deep reinforcement learning based optimal route and charging station selection," *Energies*, vol. 13, no. 23, p. 6255, 2020.
- [9] L. Zhang, K. Gong, and M. Xu, "Congestion control in charging stations allocation with q-learning," *Sustainability*, vol. 11, no. 14, p. 3900, 2019.
- [10] T. Qian, C. Shao, X. Wang, and M. Shahidehpour, "Deep reinforcement learning for ev charging navigation by coordinating smart grid and intelligent transportation system," *IEEE Transactions on Smart Grid*, vol. 11, no. 2, pp. 1714–1723, 2019.
- [11] M. Aljaidi, N. Aslam, X. Chen, O. Kaiwartya, and Y. A. Al-Gumaei, "Energy-efficient EV Charging Station Placement for E-Mobility," in *IECON 2020 The 46th Annual Conference of the IEEE Industrial Electronics Society*. IEEE, 2020, pp. 3672–3678.
- [12] M. Aljaidi, N. Aslam, X. Chen, O. Kaiwartya, and M. Khalid, "An Energy Efficient Strategy for Assignment of Electric Vehicles to Charging Stations in Urban Environments," in *2020 11th International Conference on Information and Communication Systems (ICICS)*. IEEE, 2020, pp. 161–166.
- [13] T. Ma and O. A. Mohammed, "Optimal charging of plug-in electric vehicles for a car-park infrastructure," *IEEE Transactions on Industry Applications*, vol. 50, no. 4, pp. 2323–2330, 2014.
- [14] F. V. Cerna, M. Pourakbari-Kasmaei, R. A. Romero, and M. J. Rider, "Optimal delivery scheduling and charging of evs in the navigation of a city map," *IEEE Transactions on Smart Grid*, vol. 9, no. 5, pp. 4815–4827, 2017.
- [15] M. Liu, Y. Wan, F. L. Lewis, and V. G. Lopez, "Adaptive optimal control for stochastic multiplayer differential games using on-policy and off-policy reinforcement learning," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 12, pp. 5522–5533, 2020.
- [16] M. Dabbaghjamesh, A. Moeini, and A. Kavousi-Fard, "Reinforcement learning-based load forecasting of electric vehicle charging station using q-learning technique," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 6, pp. 4229–4237, 2020.
- [17] S. Wang, S. Bi, and Y.-J. A. Zhang, "A reinforcement learning approach for ev charging station dynamic pricing and scheduling control," in *2018 IEEE Power & Energy Society General Meeting (PESGM)*. IEEE, 2018, pp. 1–5.