

# Improving the Software-Defined Wireless Sensor Networks Routing Performance Using Reinforcement Learning

Muhammad Usman Younus<sup>✉</sup>, *Member, IEEE*, Muhammad Khurram Khan<sup>✉</sup>, *Senior Member, IEEE*, and Abdul Rauf Bhatti

**Abstract**—Software-defined networking (SDN) is an emerging architecture used in many applications because of its flexible architecture. It is expected to become an essential enabler for the Internet of Things (IoTs). It decouples the control plane from the data plane, and the controller manages the whole underlying network. SDN has been used in wireless sensor networks (WSNs) for routing. The SDN controller uses some algorithms to calculate the routing path; however, none of these algorithms have enough ability to obtain the optimized routing path. Therefore, reinforcement learning (RL) is a helpful technique to select the best routing path. In this article, we optimize the routing path of SDWSN through RL. A reward function is proposed that includes all required metrics regarding energy efficiency and network Quality-of-Service (QoS). The agent gets the reward and takes the next action based on the reward received, while the SDWSN controller improves the routing path based on the previous experience. However, the whole network is also controlled remotely through the Web. The performance of the RL-based SDWSN is compared with SDN-based techniques, including traditional SDN and energy-aware SDN (EASDN), QR-SDN, TIDE and non SDN-based techniques, such as *Q*-learning and RL-based routing (RLBR). The proposed RL-based SDWSN outperforms in terms of lifetime from 8% to 33% and packet delivery ratio (PDR) from 2% to 24%. It is envisioned that this work will help the engineers for achieving the desired WSN performance through efficient routing.

**Index Terms**—Energy optimization, Internet of Things (IoT), reinforcement learning (RL), RL-based WSN, routing, software-defined wireless sensor network (SDWSN), wireless sensor networks (WSNs).

Manuscript received April 11, 2021; revised June 23, 2021 and June 28, 2021; accepted July 12, 2021. Date of publication August 3, 2021; date of current version February 21, 2022. This work was supported by King Saud University, Riyadh, Saudi Arabia, through Researchers Supporting Projects under Grant RSP-2021/12. (*Corresponding author: Muhammad Usman Younus.*)

Muhammad Usman Younus is with Ecole Doctorale Mathématiques, Informatique, Télécommunications de Toulouse, University Paul Sabatier, 31330 Toulouse, France, and also with the Department of Electrical and Computer Engineering, COMSATS University Islamabad, Sahiwal Campus, Sahiwal 57000, Pakistan (e-mail: usman1644@gmail.com).

Muhammad Khurram Khan is with the College of Computer and Information Sciences, King Saud University, Riyadh 11653, Saudi Arabia (e-mail: mkhurram@ksu.edu.sa).

Abdul Rauf Bhatti is with the Department of Electrical Engineering and Technology, Govt. College University Faisalabad, Faisalabad 38000, Pakistan (e-mail: bhatti\_abdulrauf@gcuf.edu.pk).

Digital Object Identifier 10.1109/JIOT.2021.3102130

## I. INTRODUCTION

THERE are tiny sensor nodes in wireless sensor networks (WSNs) that may be stationary or mobile nodes deployed in a dynamic environment. Each sensor node consists of a small power source, transmission, and processing units [1]. The WSN is an application-oriented information-centric network used in many applications, including defense, environmental monitoring (i.e., light, temperature, humidity, and vibration), military, and health [2], [3].

Much of the recent research work in the WSN area focuses on low-cost and low-power networking solutions to execute the cooperative and collaborative tasks under rigorous computation and energy constraints. Therefore, the wireless sensor nodes operate for a long time without replacing their batteries in many applications [4]. Thus, the sensor nodes' energy consumption is considered extremely crucial in the wireless network design. Routing is the core networking activity in WSNs that routes the sense data from the source (the sensor node) to the destination (the sink). It significantly impacts the network performance, such as energy consumption, delay, latency, and packet delivery ratio (PDR). In order to make WSNs more efficient, routing strategies, such as software-defined networking (SDN) [5], [6] and reinforcement learning (RL) [7], can be an excellent choice to get an optimized routing path in WSNs.

SDN is an emerging architecture that manages the network efficiently. It has three planes that include a data plane, a control plane, and an application plane. It decouples the control plane from the programmable data plane. The idea behind it is to manage and utilize network resources efficiently. In SDN, the network is controlled through a centralized controller (the control plane) that can globally view the whole underlying network and packet forwarding devices on the data plane. The control plane is responsible for routing, traffic management, and fault recovery; however, the data plane manages the packet delivery. SDN uses the well-defined interfaces between each plane, such as the southbound interface between the data plane and the control plane, while the northbound interface between the control plane and application plane. It also uses the OpenFlow communication protocol used by the SDN controller to communicate with data plane devices, i.e., sensors, switches, and routers. The use of SDN in the Internet of Things (IoT) is increasing day-by-day because of its efficient structure

TABLE I  
COMPARISON OF OUR CURRENT WORK AND THE PREVIOUS WORK

Parameters	Previous research work [10]	Current research work
QoS parameters used in Reward function	X	✓
QoS based algorithms for both Controller and Node side	X	✓
PDR comparison	X	✓
Results comparison with previous RL-SDWSN techniques	X	✓
IoT based control network	X	✓

that can efficiently manage and control the billion of network devices [8]. Software-defined IoT uses in a different networks, such as edge networking, access networking, core networking, and data center networking. IoT is also comprised of several wireless devices managed through SDN; however, the IoT network is still facing some issues and challenges related to security and scalability [9]. SDN is also used in WSN, known as a software-defined WSN (SDWSN), that makes it robust and well organized. However, SDWSN still has some limitations, such as finding the best routing path. It can be solved through a learning technique called RL. RL is an efficient method for a real-time path optimized path in reducing energy consumption and delay, and increasing PDR.

In our previous contribution [10], we used RL to optimize the energy consumption of SDWSN. The reward was calculated through estimate node life time (ENLT) and path estimate life time (PELT), but the quality of service (QoS) parameters have not considered in [10] that may reduce the network performance. To calculate the energy consumption of SDWSN, a different model was used. However, the comparison of the proposed work has compared with basic techniques only, instead of the latest work. The main focus of [10] was to optimize the energy consumption of SDWSN. The comparison of our current and the previous research work [10] is given in Table I.

RL is a type of machine learning (ML), in which a learner is known as an agent who learns the optimal policy ( $\pi$ ) based on rewards and experiences. It satisfies the Markov property, called the Markov decision process (MDP) [7], [11], which gives the concept of RL. RL agent in the environment can be explained as  $(S, A, P, R)$ , where  $S$  is a set of states  $\{s_1, s_2, s_3, \dots, s_m\}$ ,  $A$  is the set of actions  $\{a_1, a_2, a_3, \dots, a_m\}$  that agent goes from one state to another state, while  $P$  is the state transition probability. The agent  $A$  will get positive or negative feedback after each round, called reward  $R(s'|s, a)$ , and select the next action according to the reward received, as shown in Fig. 1. The reward can be maximized by optimizing the policy  $\pi : A \leftarrow S$ .

The combination of both SDWSN and RL can efficiently manage the network and improves the network performance. The RL-based SDWSN architecture is shown in Fig. 2. Some artificial intelligence (AI) techniques (Deep RL) are used

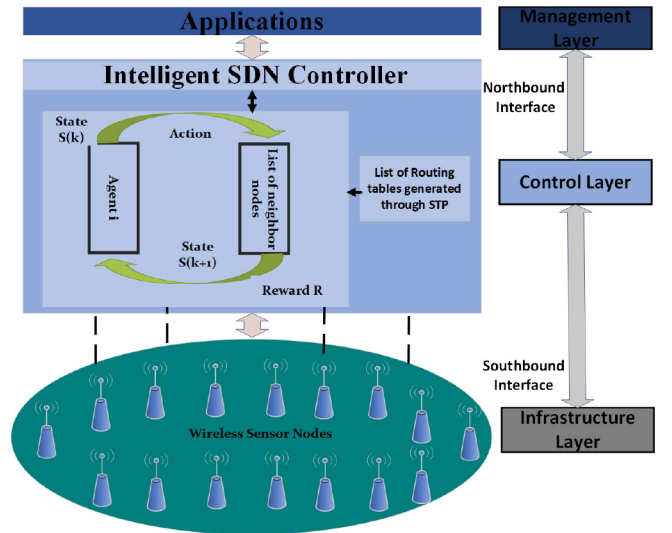


Fig. 1. RL-based SDWSN architecture [10].

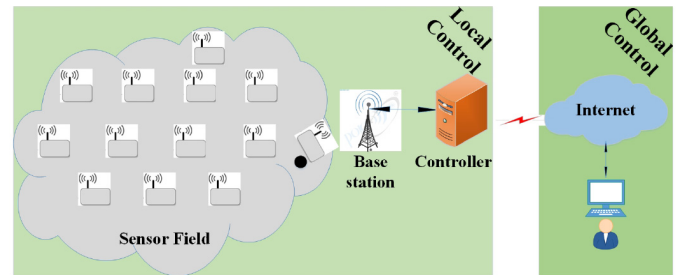


Fig. 2. WSN framework for IoT application.

in SDN for the network self-learning that can control the network efficiently [12]. However, the WSN network can be controlled remotely by the IoTs [13], [14]. It is a new paradigm that enables any device to connect with another device through the Internet. The IoT architecture consists of low-processing devices called nodes, the local controller that collects the data from nodes, and the cloud layer that monitors and controls the nodes remotely known as the global control mechanism. The WSN framework for IoT application is shown in Fig. 3. For the local controlling and optimization, we use the SDN architecture and RL for optimizing the routing path. However, we are also controlling the sensor node globally through the Web (Internet) as shown in Section VI. The contribution of this article is summarized as follows.

- 1) A QoS-based reward function, including different QoS parameters, is proposed to optimize the routing that selects the best path from the routing list. As a result, it reduces the network's energy consumption and improves the network lifetime (LT) and PDR.
- 2) Two algorithms are proposed to improve the network performance, such as energy efficiency, PDR, and, etc. These algorithms are applied to an intelligent SDN controller, which can efficiently control the data plane devices.

- 3) The loop-free communication is established through the spanning tree protocol (STP).
- 4) A real-time experimental platform is developed using Raspberry Pi, where SDWSN routing based on the RL experiment is carried out.
- 5) A Web-based dashboard is also developed to control and analyze the sensor nodes data remotely.
- 6) The proposed RL-based SDWSN technique is compared to RL-based and non RL-based techniques/algorithms on a real-testbed.

The remainder of this article is organized as follows: Section II provides the literature review related to RL-based SDWSN and non RL-based routing (RLBR) techniques. In Section III, RL is explained. However, the reward function is also explained in detail in Section III. Section IV provides the detail of the energy consumption model used for the experimental work. In Section V, the proposed methodology and algorithms for the SDN controller and sensor node are explained. Section VI provides the details of the experimental testbed and also compares the proposed RL-based SDWSN algorithm with different routing protocols. Finally, this article is concluded in Section VII.

## II. RELATED WORK

Routing can optimize the performance of WSN, such as energy efficiency and QoS parameters. In optimal routing, we adopt the SDN architecture and propose some algorithms to improve SDN-based network performance through routing. While, sometimes, the SDN-based network has poor real-time performance. To resolve this issue, RL can play an essential role in optimizing network performance. This section divides the routing approaches into two different groups: 1) SDN-based routing approaches and 2) RL-based SDN routing approaches.

### A. Non RL-Based Routing Techniques

Junli *et al.* [15] proposed an energy-aware routing algorithm for the SDN-based WSN network. The author adopts the extreme assumption for implementation that includes the SDN controller knows the initial status data (i.e., position, energy, and neighboring nodes) of all sensor nodes. The controller also has direct access (e.g., the controller can send a routing table directly to each node) to all the nodes, which is practically impossible. A distance-based routing path is established by using a Dijkstra algorithm and changes it if any node runs out of energy. An energy-aware routing protocol is proposed for the SDN network known as energy-aware SDN (EASDN) [16]. This article proposes three algorithms, including neighbor discovery algorithm, status data collection algorithm, and controller operation phase algorithm. In EASDN, the controller uses two parameters of distance and residual energy to calculate the routing path. The SDN controller sends a new routing table whenever he observes any node runs out of energy or is below the user's threshold. Raspberry pi is used to perform a real-time experiment as a sensor node. However, Li *et al.* [17] proposed a traffic control scheme that monitors the traffic behavior and prevents

traffic congestion through deep packet inspection (DPI). In another work [18] proposed a novel green datapath framework was proposed to reduce energy consumption by looking for energy-efficient routing paths for ternary content addressable memory (TCAM)-based SDN. TCAM provides search operations for packet switching networks and is used by networking devices to speed-up the networking tasks, for example, packet classification. This study of green datapath focuses on TCAM usage in hybrid SDN networks by introducing dynamic voltage and frequency scaling (DVFS) power management technique. Hence, TCAM memory deploys the enriched routing functionality and also provides speed-up interaction between switches and controllers. However, in [19], the energy-efficient architecture is proposed for SDWSN. It reduces data packet generation by content awareness and adaptive data broadcast of cached data. While, an SDN-based routing protocol is developed for the multihop wireless network [20]. The routing path develops through residual energy and a minimum number of hop counts for getting the shortest path. The model was developed using OPNET. The simulation results are compared to old routing protocols, such as optimized links state routing protocol (OLSR), and *ad hoc* on-demand distance vector (AODV). A software-defined energy-aware routing (SD-EAR) is proposed in [21] to reduce the network's energy consumption. The network is divided into different zones or clusters, and each zone is controlled through the SDN controller. The SDN controller knows each zone topology, selects the energy-efficient path based on the global view, and knows each node's residual energy. While in [22], an energy-efficient routing algorithm is developed for SDWSNs. The proposed algorithm selects the control node to assign different tasks dynamically. The NP-hard definition is used to select control nodes. For tackling the NP-hard problem, an improved particle swarm optimization (PSO) algorithm is proposed. However, PSO often becomes more complicated in solving these types of issues.

SDN-based optical networks support huge IP traffic with great energy efficiency. Therefore, an SDN-based cross-network is proposed in [23] for joint energy efficiency designs from wired to wireless networks, potentially facilitating a large-scale deployment. This approach uses SDN to focus on designing the improved communication protocols, considering that the sensors are not constrained by the battery-based power, as energy efficiency has always been one of the important goals of cellular wireless communications. A great work made by Wu [24] focuses on the design of green cloud radio access network (C-RAN) through jointly assuming the power consumption of remote radio head (RRH) and transport network. Ejaz *et al.* [25] proposed an energy-efficient mechanism for transferring the power wirelessly to SDWSN. The proposed process specifies the minimum number of energy transmitters. The optimization problem is formulated to find the minimum number of energy transmitters that reduce the network's energy consumption. However, SDN can also improve WSN management. Ndiaye *et al.* [26] addressed the many traditional and SDN-based protocol. This article gives the inside detail of SDN-based management techniques for WSN. The benefits of SDN-based network management are discussed, some

challenges of SDN faces are also explained. For multiprotocol label switching (MPLS) [27], the energy-aware routing and resource management model is developed using SDN. First, the SDN controller sets up the multipath and then uses these pre-multipath paths (PMPs) for routing. However, depending on traffic conditions, the SDN controller saves energy by turning the PMP on or off. Load balancing and route resizing are often used to utilize network resources and effectively minimize network energy consumption. However, global power management is achieved through SDN in [28]. In the proposed technique, the traffic is rerouted and adjust the network workload in different links. A network topology is constructed according to routers connection. The integer linear programming model (0-1) minimizes the integrated chassis and line-cards power. Two algorithms (alternative greedy algorithm and global greedy algorithm) are proposed for efficient link utilization and packet delay reduction.

WSN network reliability and traffic load management are also a critical problem discussed in [29] and [30]. A system called improved software-defined WSN (Improved SD-WSN) is proposed to solve the WSN reliability problem. The proposed structure tackles heterogeneous network management and network coverage problem, which can increase network reliability, where the traffic load minimization (TLM) problem in SDWSN is resolved through the flow splitting optimization (FSO) algorithm. In [30], two methods are used to reduce the traffic load issue; the first one is used to select the optimal relay sensor nodes, while the second one is responsible for the transmission of optimal splitting flow.

A recent study in [31] confirms the great potential of information and communications technologies (ICTs) to achieve social and environmental advantages. Such systems make a significant contribution for maintaining user connectivity requirements, ensuring the desired user experience through the deployment of new frameworks. Despite the abundance of energy-saving technologies, ICTs systems are critical regarding the current and future energy consumption of telecommunication networks. Most of the evidence speak against flattening or reducing ICT power. Great efforts have been made to improve one or more specific features of energy-efficient-based systems. Energy efficiency has gained great importance in the green spectrum. Green topics are multidisciplinary in nature, encompassing energy-related issues and the broader context of the environmental impact of the IoTs. Green IoTs can be integrated into standalone lighting systems [32], from green energy harvesting to smart energy management, which means significant energy savings. In machine-to-machine (M2M) systems, some machine nodes are battery-powered and need to operate for long periods of time without battery replacement. Therefore, it is very important for such nodes to be energy efficient. To that end, researchers have come up with several techniques to make these nodes energy efficient. To reduce the energy consumption, a wake-up/paging strategy is proposed in [33] by activating only those nodes ready to transmit and receive. Using a more efficient XML interface, core links, and protocol buffers, a security technique is used to control sleeping devices from attackers.

### B. RL-Based Routing Techniques for Software-Defined Wireless Sensor Networks

RLBR protocols are used to improve network performance. An RLBR protocol is defined in [34], called RLBR. The authors suggested the reward function in RLBR, which uses neighbors' distance, residual energy, and hop count to sink to measure the reward. It optimizes WSN routing, which helps to increase the network LT. For the first time, Watkins proposed an estimation action function called  $Q$ -function in [35] to reduce the time delay in the network. After each iteration, the  $Q$ -function was modified and converged to an optimum point after some iterations. The combination of RL and SDN, however, provides a stronger solution. A prototype is proposed for improving the energy efficiency and adaptability of SDNWSN by using RL [36]. For monitoring the environmental application, the prototype considers the energy, computational capabilities, and radio resources to optimize the WSN performance. RL is not adequately used to optimize SDWSN performance. However, in the SDN-based network, the SDN switches face the service placement problem that increases end-users' service costs. An algorithm named  $Q$ -placement [37] is suggested to reduce the end-user expense through RL. It guarantees the efficiency of the network and increases its convergence rate. It does not, however, concentrate on the issue of energy optimization for the network. For multimedia-based SDNs, an RL-based LearnQoS system is proposed in [38]. Services based on the video are important because it has become an integral part of the end-user life. By employing policy-based network management (PBNM), the proposed LearnQoS framework increases the video QoS. In our previous paper [10], we proposed a RL-based technique for SDWSN that can optimize the energy consumption. Furthermore, different reward functions are also proposed that enhance the network performance in terms of energy. However, the reward function is considering only energy-efficient parameters. Thus, this article is just focusing on the energy consumption of the WSN networks. QoS is not considered that may reduce the network performance. In [39], QoS-aware adaptive routing is proposed to enhance the SDN-based network throughput and reduce the delay and packet loss ratio. The proposed routing protocol consists of multilayers hierarchical architecture. There are three kinds of controllers to control the network: 1) super controller; 2) domain controller; and 3) slave controller. The network's performance is enhanced through RL. The designed reward function optimizes the QoS and achieves some metrics, such as QoS-provisioning packet forwarding, time-efficiency, and fast convergence rate. Time-relevance deep RL (TIDE) is another technique presented in [40] to improve the network QoS. An AI-based layer is introduced to control the network associated with the SDN controller. In TIDE, only QoS parameters are considered to make routing decisions. However, in [41], the QR-SDN routing approach provides multiple routing paths while preserving flow integrity. An other algorithm [42], that is based on  $Q$ -learning, sets up the  $Q$ -tables. The goal of the proposed algorithm is to find the best link for data forwarding. But the main concern is QoS. The reward value is fixed against the QoS performance of each link, so that if the QoS

performance of the link is between 0% and 30% then the reward value will be 50, if the QoS performance of the link is between 31% and 60% then the reward value will be 100; otherwise, it will be 150. In the future, deep learning will be combined with  $Q$ -learning, which is called deep  $Q$ -learning. In the modern era of information technology, most services mainly rely on cloud infrastructure. The network flow includes the mice and elephant flows in the cloud environment. For the management of these flows, SDN-based approaches are used to allocate the network resources efficiently. However, due to TCAM's limited capacity in OpenFlow enabled switches, the management of flow is a big issue. Still, it is challenging to find the useful forwarding rule in the flow table, whose rules need to be processed by the SDN controller. It is required to investigate the useful flow items in the flow table and exclude the remaining flow items. By doing this, control plane overhead can be reduced. For this reason, it is helpful to use RL in [43]. To minimize the overhead between the SDN controller and the switch, the author proposes an algorithm based on RL and an algorithm based on deep RL. Cybersecurity in SDN also becomes a critical issue. Different forms of causative attacks are studied using RL to learn how to react. RL-based algorithm is proposed in [44] to prevent the SDN network. However, the proposed algorithm's basic purpose is to make the SDN platform's autonomous defense system, as IoT is a new paradigm to interconnect the large-scale wireless devices and has a wide range of application prospects. In [45], AI-based solutions have been addressed related to spectrum access and random access. In addition, deep learning algorithms can make IoT smarter and friendly, as the development of AI is crucial for IoT, which can be strongly supported in different ways. It also facilitates AI's development direction by proposing a DQN for conducting efficient online training for RL. A learning-based algorithm protects WSN from any external attack in [44]. Each sensor node plays a critical role in WSN, and it should be protected from external attacks. Mostafaei and Obaidat [44] proposed a self-protected learning algorithm (SPLA) that can resist external attacks. In the proposed algorithm, sensing graph (SG) is used to model the problem, determining the minimum number of the nodes that can be protected. Each node is equipped with a learning automation. SPLA tries to find the minimum number of nodes that can be activated to protect the network. However, in [46], a decentralized swift vigilance (DeSVig) framework is proposed to identify adversarial attacks in an industrial AI system. Due to decentralization, DeSVig improves the effectiveness of recognizing abnormal inputs.

### III. $Q$ -LEARNING

$Q$ -learning is a model-free learning approach first proposed by Watkins [35]. It is used to estimate  $Q'(s_t, a_t)$ , called  $Q$ -function, and its learning technique is called  $Q$ -learning. In  $Q$ -learning, a learner is known as an agent that selects an action according to the current state by interacting with the environment and getting either positive or native feedback based on the action called reward and calculates the  $Q$ -value. In the next state  $s_{t+1}$ , the agent selects an action based on the previous

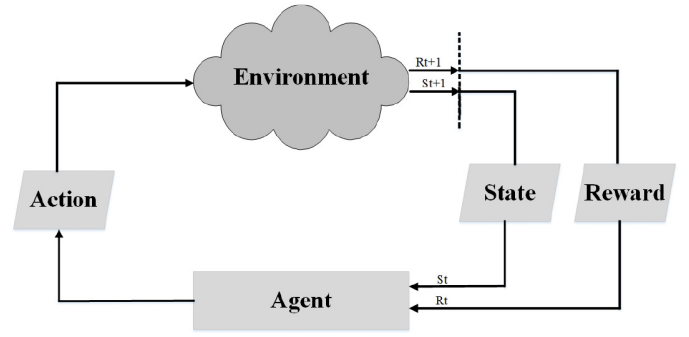


Fig. 3. RL model.

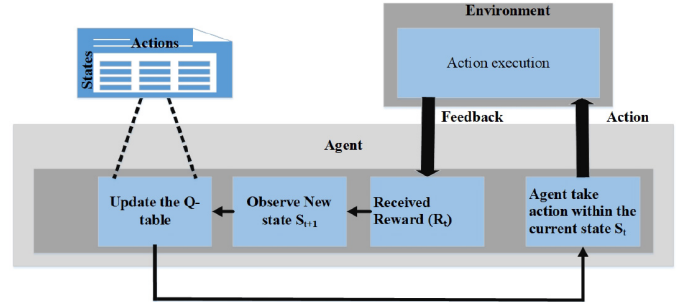


Fig. 4.  $Q$ -learning block diagram.

reward. After some rounds, the agent is trained and converges to the optimal point. In  $Q$ -learning, the  $Q$ -value is updated after each iteration, as shown in Fig. 4. The update  $Q$ -value is defined in the following:

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

$Q(s_t, a_t)$  is the  $Q$ -value that the agent takes an action at in state  $s_t$  and gets an immediate reward  $R_t$ , where  $\max_{a'} Q'(s_{t+1}, a_{t+1})$  is the maximum value that can get in the next state  $s_{t+1}$ .  $\alpha$  is the learning rate that determines the what extend the newly acquired information update to the old information and the range is  $(0 < \alpha \leq 1)$ , where  $\gamma$  is the discount factor that determines the importance of future reward and its range is  $(0 < \gamma \leq 1)$ .

#### A. $Q$ -Routing Protocol

It comes from the  $Q$ -table, which uses the  $Q$ -learning method to route the data packets [35]. In  $Q$ -routing, the first  $Q$ -matrix of node  $i$  is initialized. The initialization of the  $Q$ -matrix may be random. Then node  $i$  sends a packet  $P$  to neighbor node  $j$ . Node  $i$  selects the forwarder node  $j$  with the lowest  $Q$ -value because of low distance up to destination  $d$ . Low distance has a low delivery delay (end-to-end delay), which is routing cost.  $Q_i(d, j)$  is the delivery delay calculated from (3) and  $t_j$  is the estimated remaining time in the trip that node  $i$  immediately gets back  $j$ 's and it can be found by

$$t_j = \min_{k \in \text{Ng}(j)} Q_j(d, k). \quad (2)$$

Let  $Q_i(d, j)$  be the estimated time that node  $i$  takes to deliver the packet  $P$  up to node  $j$  using a distance, and it can be



**Algorithm 1: Q-Routing Algorithm**


---

```

begin Initialized the Q-matrix (Q(x,y))
while (Until the terminal state is reached) do
    if Is data packet ready to send then
        Calculate the Q-Value.
        Select next hop  $j$  with lowest Q-value.
        Send the packet to selected forwarder node  $j$ .
        Calculate the delivery delay time using Eqns (2)
        and (3).
        Update the node  $i$  delivery delay time.
    end
end

```

---

calculated by

$$Q_i(d, j) = (1 - \alpha) * Q_i(d, j) + \alpha * (q_i + Tx + t_j). \quad (3)$$

$q_i$  is the time to spend packet  $P$  in the queue of node  $i$ ,  $\alpha$  denotes the learning period, and  $Tx$  represents the time used for transmission between node  $i$  and node  $j$ .

The  $Q$ -routing algorithm is given in Algorithm 1.

### B. Q-Learning-Based Routing in SDWSN

Consider a network containing  $n$  sensor nodes, and each sensor senses the environment and selects one neighboring node  $j$  to send the data packet up to the destination  $d$ . For selecting the best optimal path,  $Q$ -learning uses the previous experience to select routing paths. In the next epoch, the sensor node  $i$  selects the best path to send the destination's data packet. We use  $Q$ -learning for the route selection in SDWSN. In the SDN-based network, the controller controls the whole network. It selects the best path according to the implemented algorithm. We use the STP for getting all possible routing paths at the controller. SDN controller selects a routing table from the list of paths with  $Q$ -learning and sends the nearest neighboring node.  $Q$ -learning includes the state, action, reward, and  $Q$ -value, which are defined as follows.

1) *States*: Let  $S$  be the set of states  $S = \{s_i, s_j, s_k, \dots, s_n\}$ , which means that after each round, the controller selects another routing table from the routing table list.

2) *Actions*: An action is an act where any path from the set of routing paths  $\{p_1, p_2, p_3, \dots, p_n\}$  is selected by the agent and sends it to the neighboring node.

3) *Reward Function*: The reward  $R_t$  is the feedback after agent action  $a_i$ , which can be positive or negative. It has much importance in  $Q$ -learning because the next action will be based on the  $Q$ -value, which varies according to the received reward. The optimal strategy will be obtained after each iteration. The reward function uses different metrics to calculate the reward. In [35], Watkins only used packet delay in the reward function that is not enough for the energy-efficient network. This article uses different reward function metrics, including the distance to the sink, the number of hops from neighboring nodes to destination/sink, node residual energy, and packet success ratio. Each metric's weight is considered in the proposed reward and took the sum of all nodes reward defined in the following:

$$R_t = \sum_{i=0}^N (w_1 * d(s_i) + w_2 * H(s_i) + w_3 * E(s_i) + w_4 * p(s_i)) \quad (4)$$

$$d(s_i) = \frac{\text{Dist\_to\_Negb}(s_i)}{\text{Dist\_to\_Sink}(s_i)} \quad (5)$$

where  $\text{Dist\_to\_Negb}$  is the distance from the neighboring node to the destination/controller and  $\text{Dist\_to\_Sink}(s_i)$  is the maximum distance to the sink

$$H(s_i) = \frac{H_j(s_i)}{H_{\max}(s_i)}. \quad (6)$$

$H(s_i)$  is the ratio of the number of hops from neighboring nodes to the destination and the maximum possible hop counts

$$E(s_i) = \frac{E_R(s_i)}{E_{\text{total}}(s_i)}. \quad (7)$$

$E(s_i)$  is the ratio of remaining energy  $E_R(s_i)$  to the total energy  $E_{\text{total}}(s_i)$

$$p(s_i) = \frac{p_{\text{ack}}(s_i)}{p_{\text{send}}(s_i)}. \quad (8)$$

$p(s_i)$  is the ratio of packet acknowledgment  $p_{\text{ack}}(s_i)$  received by the neighboring node to the total packet  $p_{\text{send}}(s_i)$  sent.

All these factors are used in the reward function to calculate the reward, which is given in (5)–(8).

The intelligent SDN controller calculates the reward after each round. After getting a reward, the  $Q$ -value is updated. The new  $Q$ -value of the state-action ( $s_t, a_t$ ) pair is determined by

$$Q(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha[R_t + \gamma \max Q(s_{t+1}, a_{t+1})]. \quad (9)$$

## IV. ENERGY CONSUMPTION MODEL

Each sensor node contains its small power supply that is required for transmission and reception. Each node consumes energy during transmission and reception. In order to calculate the energy consumption during data communication, we consider the first-order energy model [47]. In this model, two types of channels are used to calculate path loss: one is a free space model and the second is multipath fading [16], [22]. The model selection is based on the distance between the transmitter and the receiver. If the transmitter–receiver distance is less than or equal to the threshold  $d_0$ , then the free space model is selected; otherwise, multipath is selected. The following equations measure the energy consumption due to the transmission of each packet:

$$E_{tx}(b_l, d) = \begin{cases} b_l * E_{fs} * d^2 + b_l * E_{\text{elec}}, & d \leq d_0 \\ b_l * E_{mp} * d^4 + b_l * E_{\text{elec}}, & d > d_0. \end{cases} \quad (10)$$

$E_{tx}$  is the energy consumption due to transmission,  $b_l$  is the total packet length in bytes, and  $d$  is the distance between the sender and the receiver. Since  $E_{fs}$  denotes the energy consumption due to free space,  $E_{mp}$  is the energy consumption of multiple paths. While  $E_{\text{elec}}$  refers to the energy consumption due to circuit processing before the transmission of a data packet and after the reception of a data packet. As  $d_0$  is the distance threshold, it can be found by the following equation:

$$d_0 = \sqrt{(E_{fs}/E_{mp})}. \quad (11)$$

The energy consumption due to the reception of a data packet is calculated by

$$E_{rx}(b_l) = b_l * E_{\text{elec}}. \quad (12)$$

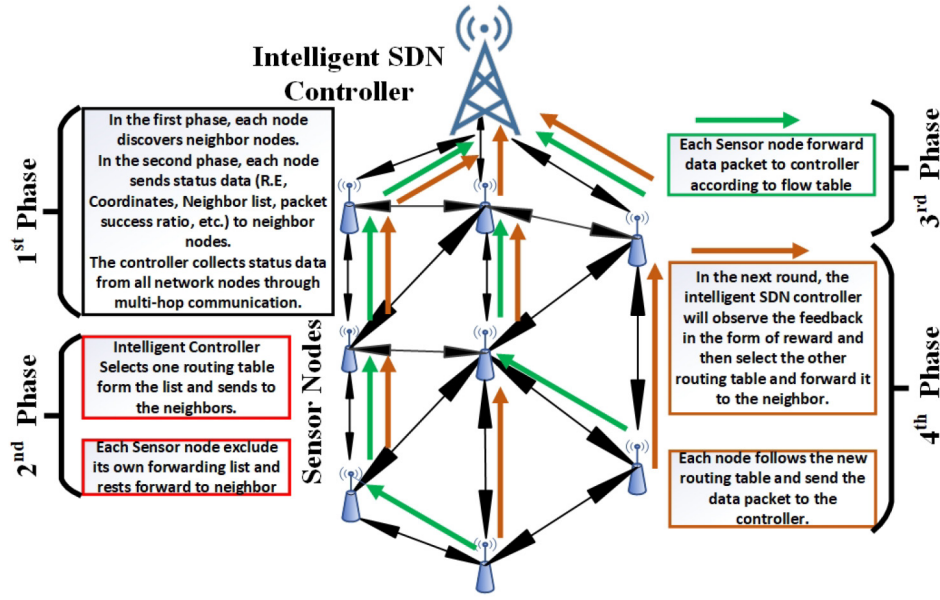


Fig. 5. Establish the route according to receive the flow table from the controller example.

$E_{rx}$  is the energy consumption due to the receiving bytes.

## V. PROPOSED METHOD

In the WSN, each sensor node collects data from the environment according to the network design and communication range and sends it to the sink through a single/multi-hop. In traditional WSN, each node broadcasts a control packet for getting the neighboring node information. The network consumes a lot of energy because of broadcasting, and also generic algorithms cannot optimize the path. SDN is an emerging architecture in which the data plane is separated from the control plane. It manages the network through a centralized controller that can view the network globally. However, the SDN controller algorithms are inefficient in real-time routing path optimization. RL is an efficient technique for learning that can optimize the real-time routing route. The best choice for optimizing the WSN routing path is the combination of both SDN and RL. This article uses RL on the SDN controller to choose the routing list's best routing path and change the routing path if it finds the path to be bad.

We propose a reward function in our energy optimization approach including all the required parameters related to network performance. A reward function consists of the distance to sink, the number of hops to sink, the residual energy, and the packet success rate. To optimize the WSNs routing path, this algorithm is divided into two parts: one for the intelligent SDN controller and the other for the sensor nodes. On the controller side, neighboring nodes are found after initializing the controller node that collects the entire network's status data, as shown in the first phase of Fig. 5. The controller generates all possible routing paths through STP [48], [49]. It selects one routing table from the list with  $Q$ -learning and sends it to the neighboring nodes, as shown in the second phase of Fig. 5. After each epoch, the controller gets the

status data of the sensor nodes and calculates the reward. The intelligent SDN controller changes the routing path according to system feedback in terms of reward, as shown in the third phase of Fig. 5. If the reward is negative, it will decrease the network's performance and change the path; otherwise, it maintains the same path. The controller also continuously monitors each node's remaining energy. If any node has energy less than the threshold, it is excluded from the node list, recalculates the routing path list using STP, selects a routing table from the list, and sends it to the neighboring node. The algorithm for the SDWSN controller is shown in Algorithm 1.

However, on the node side, the first neighbor discovery period starts, and each node broadcasts the Hello packet to find the neighbors, as shown in Algorithm 3. This process continues up to a specific time and a maximum number of an acceptable neighbor threshold. After the neighbor discovery period, each node shares the status data with its neighboring nodes. The status data of each node reach their destination (controller) through multihop communication. Each node receives a routing table from the SDN controller and sends a data packet according to the routing table provided by the intelligent SDN controller (i.e., Algorithm 2). The sensor node calculates the energy consumption using a mathematical model at the end of each round; and sends the residual energy status to the controller through a data packet. In the case of the relay node, it first checks the node's remaining energy; if it is greater than the threshold, it accepts; otherwise, it sends a low energy message to the controller and disconnects from the network. While the computational complexity of the energy-efficient algorithm of both controller and node side is also examined, the controller side runs  $n$  times operation until the last node dies, and the complexity of the whole controller side operation is  $O(n)$ . While, on the node side, the first operation remains valid up to two threshold approaches  $t_{nbr}$  and  $NBR_{max}$ . This operation runs  $n$  times while the two inside operations take  $n$  times that receives the status data of neighboring

**Algorithm 2:** Algorithm for SDN Controller

---

**Input** Network status data including total number of edges, Vertex,  $G=(V, E)$ , STP, Reward function, Learning rates.

**Output** Set of routing paths.

Initialize the controller.

Assign the IP to controller.

Controller discover the neighboring nodes.

Collect the status data from all sensor nodes within the threshold and  $N_{max}$ .

$N_{max}$  is the maximum number of nodes that can be possible neighbors; however, the threshold is the time threshold.

Calculate the routing table using STP.

$RT \leftarrow \{x_1, x_2, x_3, \dots, x_n\}$ .

$SDNController \leftarrow RT$ .

Initially, the Q-value is considered as the worst case where all the nodes die without sending any data.

Select one routing table randomly from the routing table list.

Calculate the Reward by using equation 4.

Calculate the Q-Value using equation 9.

Update the Q-value.

**while** (Received node data upto last node die) **do**

**if** ( $E_{node}^{residual} < TH$ ) **then**

        Exclude that node from the list.

        Recalculate the routing tables using STP.

$RT \leftarrow \{x_1, x_2, x_3, \dots, x_m\}$ .

        Select one routing path from the list and send it to the neighboring node.

**end**

**else**

$RT \leftarrow \{x_1, x_2, x_3, \dots, x_n\}$ .

$SDNController \leftarrow RT$ .

        Estimate the PELT.

        Calculate the Reward by using equation 4.

        Calculate the Q-Value using equation 9.

        Select one routing table with highest Q-value.

        Update the Q-value.

**end**

**end**

---

nodes, and other operations calculate the energy consumption of node up to “Required Energy for Tx load.” The whole operation contains  $n(n + n)$  iterations. So the overall node side complexity is  $O(n^2)$ . The global optimization problem is NP-hard, and the majority of previous solutions have been heuristic algorithms that are slower than the proposed algorithms.

## VI. EXPERIMENTAL RESULTS AND DISCUSSION

This section will discuss the performance of the proposed routing technique RL-SDWSN in terms of lifetime and PDR. The performance of RL-SDWSN is compared with previous protocols, such as Q-routing, RLBR,  $T\_SDWSN$ , EASDN,

**Algorithm 3:** Algorithm for Sensor Nodes

---

**Input** Routing paths that receive from controller.

**Output** Each node  $\{S_1, S_2, \dots, S_n\}$  sends the collected information to the controller related to energy and QoS. Initialize the nodes  $S = \{S_1, S_2, \dots, S_n\}$ . Assign the IP to controller.

**if** ( $RE > E_{Threshold}$ ) **then**

$S \leftarrow parametersettingfromcontroller$

    Searchtheneighboringnode

**while** ( $time < t_{nbr}$ ) &  $len(My Neighbor list) < NBR_{max}$  **do**

**if**  $t_{lbt} < t_{max\_allowed}$  **then**

            Broadcast the Hello packet

**end**

**On** (Reception of Hello packet for neighbor)

**if** ( $Nbr\ node\ id\ does\ not\ exist\ in\ "Nbr\ list"$ ) **then**

            Add into the neighboring list.

**end**

**end**

    Calculate the energy consumption of nodes after each Tx and Rx control and data traffic.

**while** ( $t < t_{threshold}$ ) and ( $response\ number < N_{max}$ ) **do**

        Send the status data to neighboring node.

**On** Reception of Status data packet **do**

**if** ( $source\ address\ does\ not\ exist\ in\ the\ "list"$ ) **then**

                Add the sender node address into list.

                ++ Status data response number.

**end**

**end**

**while** ( $RE < Required\ Energy\ for\ Tx\ load$ ) **do**

        Tx & Rx the both control and data traffic.

        Calculate the energy consumption of nodes after each Tx and Rx control and data traffic.

**end**

    Send the notification of low energy to controller through neighboring node and disconnect.

**end**

---

QR-SDN, and TIDE implementation is given in the next section. Q learning and RLBR are based on RL, but  $T\_SDWSN$  and EASDN depend on the SDN architecture. However, QR-SDN and TIDE use RL to control the flow in an SDN-based network. Q-routing, RLBR,  $T\_SDWSN$ , and EASDN routing protocols are used to optimize the energy consumption of WSN; however, QR-SDN and TIDE are used to optimize the network QoS. In our experimental setup, we use raspberry pi for real-time experimental work and create a wireless ad hoc network. Each raspberry pi (node) contains a wireless card using the IEEE 802.11ac standard. The topology used for experimental work is shown in Fig. 6. All nodes are placed in the same location, communicate using logical distance, and access the neighboring node based on the distance threshold.

In our experimental work, the communication between the nodes is real time, but the energy consumption is calculated by simulation. The reason to use the simulation for the calculation



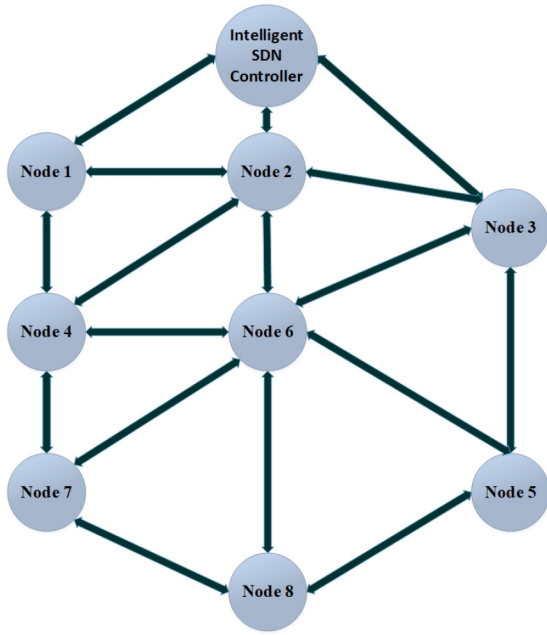


Fig. 6. Topology used for experimental work.

TABLE II  
EXPERIMENTAL PARAMETER TABLE

Parameters	Value
Initial Energy	5J
$E_{elc}$	50 nJ/bit/m <sup>2</sup>
$E_{fs}$	100 pJ/bit/m <sup>2</sup>
$E_{mp}$	0.0013 pJ/bit/m <sup>4</sup>
Data packet size	100 bytes
$\alpha$	0.3, 0.5
$\gamma$	0.5

of energy consumption is that raspberry pi could not directly calculate the energy consumption. It requires an external module to calculate energy consumption (i.e., MoPi). However, the SDN controller energy consumption is not considered here, and it is assumed that the controller has infinite energy. The simulation parameters are given in Table II. In the parameter setting, we took two different  $\alpha$  learning rates to observe the impact on network performance. In each  $\alpha$  experiment, we consider three definitions of the lifetime (the first node dies, 50% node dies, and the last node dies) and two different deployment areas (100 m  $\times$  150 m and 200 m  $\times$  300 m). Experimental settings, platforms, and performance metrics are explained in the following sections.

#### A. Experimental Platform

We did our experimental work on a real testbed using Raspberry Pi and developed the program using Python 3.0. Raspberry Pi is a low-powered, low-cost, and small-size single-board computer. It is used as a sensor node to transmit data up to the controller/sink and as a controller to control the underlying network that collects data from the sensor nodes in the SDWSN network. Raspberry pi has become more attractive because of its small size and used in many real-time applications such as WSN, cloud computing applications, robotic projects, and so on. In our experiments, we used

the Raspberry pi B+ model, which has a powerful 1.4 GHz  $\times$  Cortex-A53 CPU and ARMv8 microcontroller with 1-Gb RAM. On the software side, it supports a variety of operating systems, including Debian Linux-based operating system recommended by the Raspberry pi foundation, and also optimizes the Raspberry pi hardware.

Compared to other models, the specifications of Pi B+ are marvelous. It also contains the wireless LAN card having IEEE 802.11ac for wireless communications. In our experimental work, we use Raspberry pi WLAN to create the wireless *ad hoc* connection. Another advantage of Raspberry pi is that it has a secondary memory, storing large amounts of data on a micro SD card. It has been deployed as an intelligent sensor node, but in our experimental work, we use it as an intelligent SDN controller, and the rest of the nodes are considered as normal nodes that collect the data from the environment and send it to the SDN controller.

#### B. Performance Metrics

In this article, we consider two metrics for results evaluation, as describe below.

##### 1) Network Lifetime:

- 1) The time until the first node runs out of energy.
- 2) The time until the X% nodes runs out of energy.
- 3) The time until the last node run out of energy.

2) *Packet Delivery Ratio*: The percentage of packet successfully deliver up to a destination is called PDR. It can be represented as

$$\text{PDR} = \left[ \frac{P_{\text{TDP}}}{P_{\text{TSP}}} \right] * 100. \quad (13)$$

$P_{\text{TDP}}$  is the total delivered packet, and the  $P_{\text{TSP}}$  is the total send packets from the source to the destination.

3) *Lifetime*: We performed our experimental work in two different scenarios. In the first scenario, the deployment area is 100 m  $\times$  150 m, and in the second scenario, it is 200 m  $\times$  300 m. We consider three different lifetime definitions: 1) the first node dies; 2) 50% node dies; and 3) the last node dies. The experimental work for a lifetime is also performed with two different parameters (e.g.,  $\alpha = 0.3$  and 0.5). In the first part of the experimental work, we compare the lifetime of five different prior methods with three different lifetime definitions by considering two different areas (100 m  $\times$  150 m and 200 m  $\times$  300 m). The proposed technique RL-based SDWSN is compared with four previous techniques. In the second scenario of lifetime description, we compare all previous techniques with our proposed technique with different learning rates ( $\alpha$ ).

First, we use two dimensions to calculate the network's lifetime, as shown in Fig. 7. In Fig. 7(a), all three-lifetime definitions are taken into account, and the area of 100 m  $\times$  150 m is taken for the experiment. The proposed RL-based SDWSN technique provides a better lifetime in all definitions. In the first node dies definition, the proposed technique RL-based SDWSN has a higher lifetime than  $Q$ -learning, RLBR, traditional SDWSN, EASDN, QR-SDN, and TIDE. RL-based SDWSN is better than  $Q$ -learning because it considers all parameters (distance from the source to destination, number of hops up to the destination, remaining energy, and QoS)

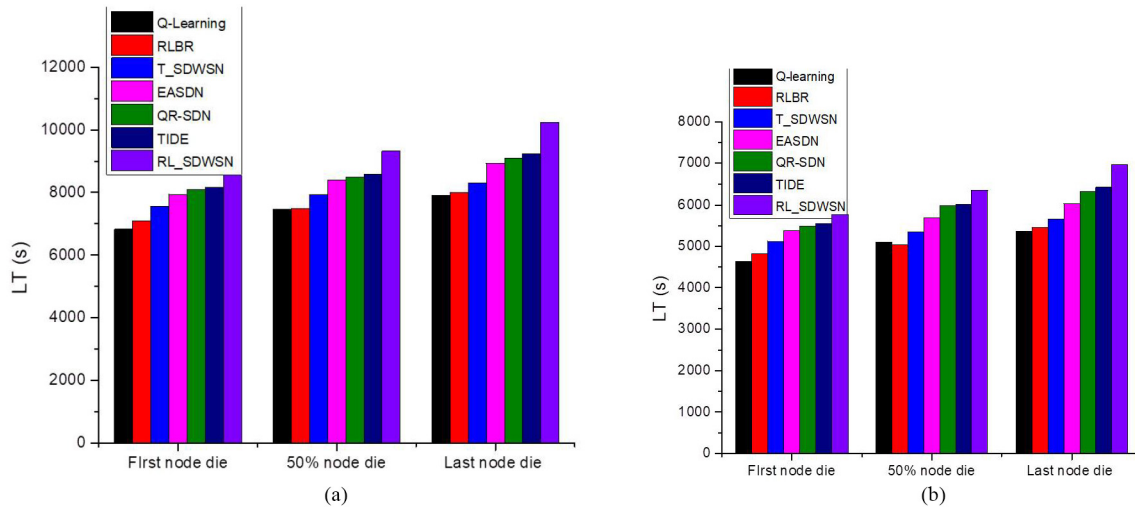


Fig. 7. Network LT. (a) Area 100 m × 150 m. (b) Area 200 m × 300 m.

related to energy optimization. But in *Q*-learning, it only considers the delay parameter and wastes a lot of network energy during path finding. *Q*-learning also uses the control packet to obtain new information about neighboring nodes. Usually, the control packet frequency is greater than the data packet frequency, causing the network LT to be shortened. However, in RLBR, it also has a shorter lifetime than RL-based SDWSN. In RLBR, a large number of control packets update the neighboring table. It is also not considered all the required parameters to find the optimized path in energy optimization. Both traditional SDN and EASDN give less lifetime than RL-based SDWSN because of nonefficient SDN controller algorithms. In traditional SDN, only the distance is considered for the calculating the routing path and also considered some assumption (i.e., the controller can access to all nodes within one hop, and also knows the initial status data of all nodes, such as node energy, position, and neighboring node list) in a real-time network, that is not possible. SDN controller cannot find the optimized path based on distance only; for these reasons, the traditional SDN has a shorter lifetime than RL-based SDWSN. In EASDN, it takes into account two parameters (i.e., distance and residual energy) to calculate the routing path. Network quality is also not taken into account when calculating routing paths. However, both QR-SDN and TIDE networks are based on the SDN architecture that is controlled through RL. These techniques do not focus on network energy consumption issues, and no energy measurement parameter is considered in the reward function. In contrast, RL-based SDWSN considers all required parameters in its reward function to find the best optimal path and change the path whenever the SDN controller observes another best path. Therefore, SDWSN based on RL has a higher lifetime than *Q*-learning, RLBR, traditional SDN, EASDN, QR-SDN, and TIDE, which are approximately 25%, 20%, 13%, 8%, 6%, and 5.5%, respectively. In the second-lifetime definition (50% node die), the proposed technique RL-based SDWSN is also outperforming as compared to SDN and RL-based techniques, as described earlier. RL-based SDWSN gives a higher lifetime than *Q*-learning, RLBR, traditional SDWSN, and EASDN,

which is approximately 25%, 24%, 17%, and 11%, respectively, as shown in Fig. 7(a). However, in the third-lifetime definition, the proposed RL-based SDWSN technique also has a higher lifetime as compared to SDN and RL-based techniques. RL-based SDWSN gives a higher lifetime than *Q*-learning, RLBR, traditional SDWSN, EASDN, QR-SDN, and TIDE, which is approximately 33%, 28%, 23%, 14%, 12%, and 11%, respectively.

In the second scenario, we enhanced the deployment area 100 m × 150 m to 200 m × 300 m and observed the lifetime performance of the network. Lifetime decreases in this scenario compared to the small area scenario because the increased distance between nodes results in more energy consumption. The energy consumption of nodes depends on the length of the data packets and the distance between the nodes, as discussed in the energy consumption model section. However, the proposed technique improves the lifetime as compared to SDN and RL-based techniques, as shown in Fig. 7(b). With the exception of RLBR in the second life cycle definition (half node die), all techniques behave similarly to the first scenario. RL-based SDWSN has a higher lifetime than *Q*-learning, RLBR, traditional SDWSN, EASDN, QR-SDN, and TIDE approximately 8% to 33% in all definitions.

4) *Comparison of Lifetime With Difference Learning Rate:* In the RL-based network, the learning rate has a significant effect on network performance. This section compares the two learning rates and looks at their impact on the network's lifetime. For this experimental work, we also consider three definitions of lifetime and two different areas. In the first scenario, the area of 100 m × 150 m is selected. As can be seen from Fig. 8, the learning rate is influencing the lifetime of the network. It increases the network LT of the RL-based network, while the SDN-based network remains unchanged due to the lack of learning. All RL-based techniques (i.e., *Q*-learning, RLBR, and RL-SDWSN) improve the lifetime as the learning rate increases. Fig. 8(a)–(c) shows that  $\alpha = 0.5$  gives a higher lifetime because it quickly learns the network and can quickly optimize the routing path. In all definitions of a lifetime, when  $\alpha$  changes from 0.3 to 0.5, the life cycle increases by 3%–5%.

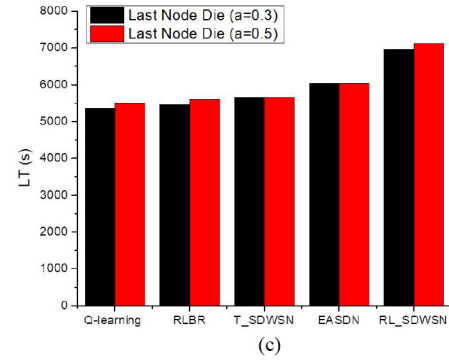
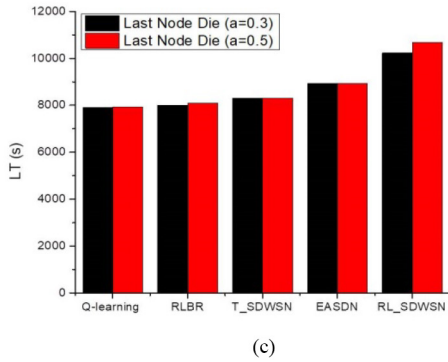
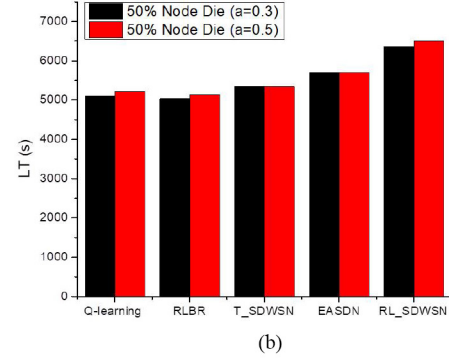
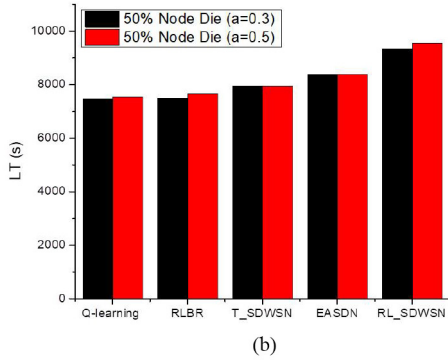
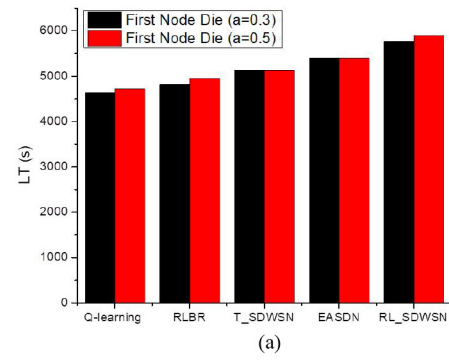
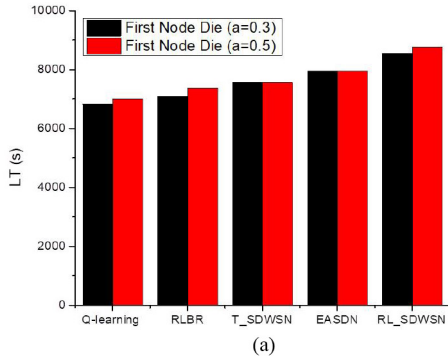


Fig. 8. Network LT (Area 100 m  $\times$  150 m) with different learning rates. (a) First node die. (b) 50% nodes die. (c) Last node die.

Fig. 9. Network LT (Area 200 m  $\times$  300 m) with different learning rates. (a) First node die. (b) 50% nodes die. (c) Last node die.

In the second scenario, we change the deployment area from 100 m  $\times$  150 m to 200 m  $\times$  300 m to see its impact on lifetime. As can be seen from Fig. 9(a)–(c), as the area increases, the lifetime of the network also decreases because of the distance between nodes increases, which increases the energy consumption of the network. However, we can also see from these results that the learning rate also affects the network LT. It increases the lifetime of RL-based techniques by 3% when  $\alpha$  increases from 0.3 to 0.5. We also observed that the proposed technique RL-based SDWSN has a higher lifetime than SDN and RL-based networks.

5) *Packet Delivery Ratio*: This section discusses the packet success rate, how many packets reach their destination, and

how much are lost in transmission. To test PDR, we also use two different areas. In Fig. 10(a), the area of 100 m  $\times$  150 m is used. It can be seen that RL-SDWSN has a higher PDR than *Q*-learning, RLBR, traditional SDN, EASDN, QR-SND, and TIDE. In *Q*-learning, each node tries to learn the best path by itself. At the beginning of the communication, many packets are lost because the network needs time to establish a path from the source to the destination. During this period, each node broadcasts control packets to get neighbor information to estimate the best possible path. However, in this case, the network becomes congested, and a large number of packets are lost due to congestion.

Similarly, in RLBR, the network is decentralized, and each node tries to learn the optimal path, but in order to learn, each

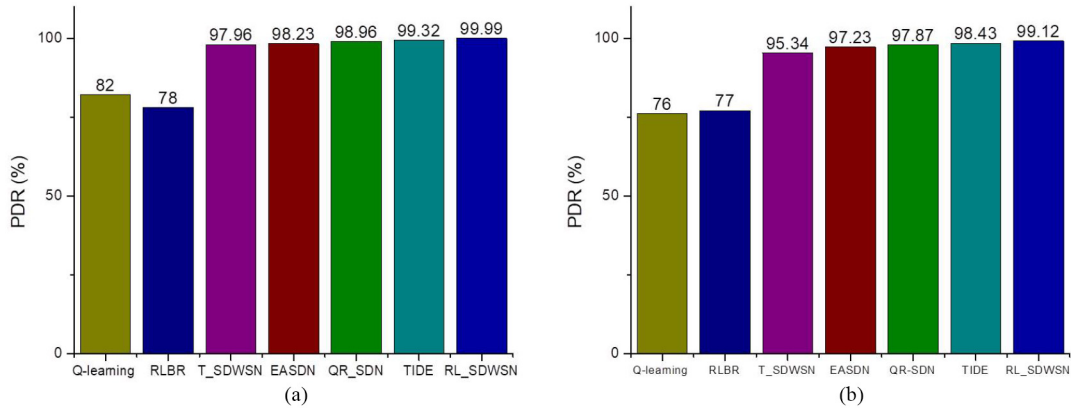


Fig. 10. PDR. (a) PDR of 100 m × 150 m area. (b) PDR of 200 m × 300 m area.

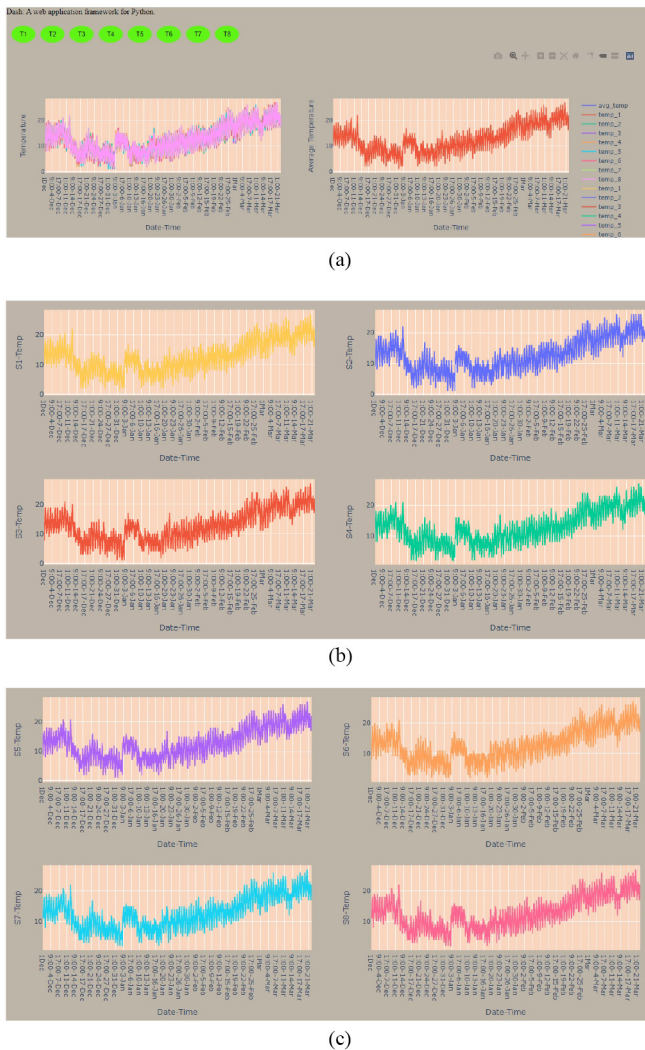


Fig. 11. Web-based sensor node monitoring dashboard 1. (a) Single dashboard for all sensor nodes. (b) Sensor nodes 1–4 dashboard. (c) Sensor nodes dashboard.

node broadcasts control packets to obtain the state information of neighboring nodes. In the RL-based network, it is necessary to obtain the latest status data of neighboring nodes. The network became congested due to the large number of control packets being exchanged. As a result, a large number of

packets are lost. However, in the SDN-based network, it has better PDR as compared to the RL-based network, but the RL-based SDWSN performance is better than that of traditional SDWSN and EASDN. In the SDN-based network, the controller can view the underlying network globally and control the network. However, both QR-SND and TIDE are also SDN-based networks. In QR-SDN, the design of the reward function is not too much affecting the network data flow control that results in losing more packets during communication and giving less PDR as compared to the proposed RL-based SDWSN. In contrast, TIDE performance is approximately approaching RL-based SDN performance because it is taken into account in the reward function that improves the network performance in terms of PDR.

The SDN controller manages network traffic to avoid congestion; that is why the SDN-based network gives better PDR than the RL-based network. However, the SDN controller algorithm sometimes does not work well. The controller did not consider all the necessary factors to avoid congestion. In RL-based SDWSN, the agent selects a routing path after considering all factors (i.e., the packet loss rate between two nodes, number of hops up to the destination, and distance up to sink). These factors affect network performance in terms of the packet success rate. In RL-based SDWSN, the agent has the ability to select the best routing path to achieve low congestion and high PRR. In the second scenario, the node deployment area is 200 m × 300 m. In Fig. 10(b), it can be seen that the RL-based SDWSN also has better PDR compared with all other techniques, and its behavior is also similar. However, the performance of RLBR is a little low compared with *Q*-learning. In terms of the delivery ratio, all methods are giving low PDR in this scenario (200 m × 300 m) than the first scenario (100 m × 150 m) because the distance between two nodes increases, and the loss rate also increases with distance. From Fig. 10(b), we can see that the RL-based SDWSN has little impact compared with other technologies. However, RL-based techniques are more effective in terms of PDR.

### C. Web-Based Sensor Node Monitoring System

IoT-based networks control the network devices (i.e., sensor nodes) remotely through the Web. We also developed a



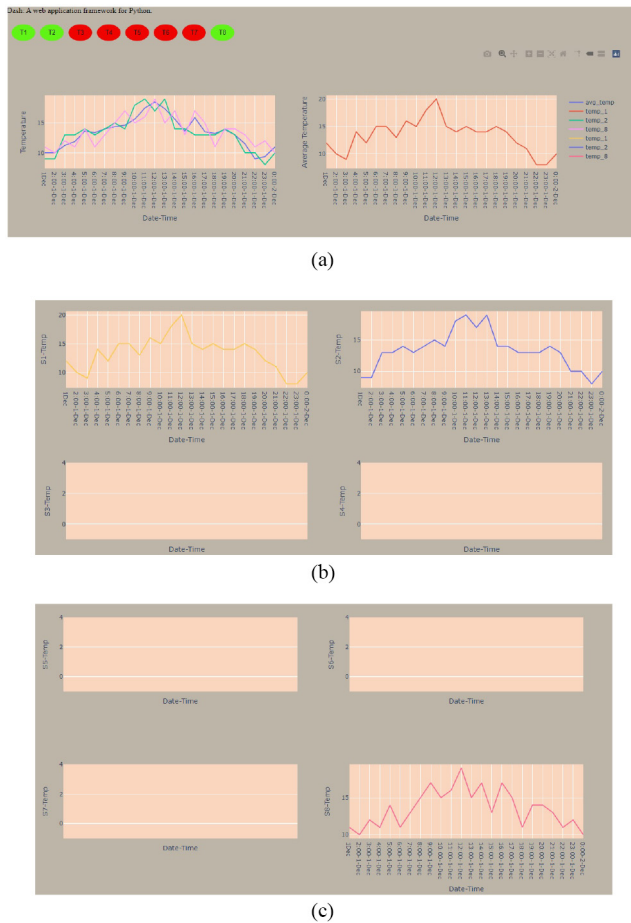


Fig. 12. Web-based sensor node monitoring dashboard 2. (a) Single dashboard for all sensor nodes. (b) Sensor nodes 1–4 dashboard. (c) Sensor nodes 5–8 dashboard.

Web-based dashboard for controlling the sensor nodes. In our network, eight sensor nodes have been deployed that measure the temperature. The dashboard is shown in Figs. 11 and 12. In Fig. 11, four months' temperature is monitored through the Web. In Fig. 11(a), two different windows are shown. The first window shows the temperature of all deployed sensor nodes, while the second window shows the average temperature of all sensor nodes. However, Fig. 11(b) shows the individual temperature of sensor nodes 1–4 and Fig. 11(c) shows the temperature from sensor nodes 5 to 8.

In Fig. 12, some sensor nodes are off (i.e., 3–6) that is shown in red color; however, sensor node 1, 2, and 8 are on and the data (temperature) of these sensor nodes are shown in Fig. 12(a). However, Fig. 12(b) and (c) shows the individual data of all sensor nodes. In Fig. 12(b), we can see that only sensor nodes 1 and 2 are active and data is also showing on the dashboard and, similarly, sensor nodes 5–7 are inactive, and the windows of these sensor nodes are empty. As we can see from both Figs. 11 and 12, network devices can be easily managed, observed, and controlled remotely through the Web.

## VII. CONCLUSION

WSNs is increasingly used in our daily life. This article elaborates the need for an efficient IoT-based WSN framework because of its great importance in various application environments. It is an intimidating task to achieve

the desired WSN performance through efficient routing. Therefore, we use RL and SDN combinations for efficient WSN routing, leading to improved routing decisions. SDN controller learns the routing path through RL and takes the next action according to the previously received reward. We used RL to select the best routing path from the routing list obtained by STP. The performance of RL-based SDWSN is compared with the existing SDN-based techniques, including the traditional SDN and energy-efficient technique used for routing. The proposed RL-based SDWSN technique provides a better lifetime of 8%–33% on average in all scenarios. We also compared the results of two different learning rates. Our proposed RL-based SDWSN method also increased the average network LT from 8% to 12% in all scenarios. Furthermore, we compared the performance of PDR that gives a higher delivery ratio as compared to SDN and RL-based techniques. Some other metrics, such as latency, delay, and throughput, also greatly impact network performance. In the future, we plan to optimize the parameters that are mentioned above by proposing some new algorithms and implementing the large-scale network. Also, we intend to compare the performance of the real-time network with the simulation work. However, we have also plan to detect the abnormal traffic of the network by using RL.

## ACKNOWLEDGMENT

The authors would like to thank Dr. S. H. Bukhari for reviewing and giving valuable comments to improve the manuscript's quality.

## REFERENCES

- [1] G. Wang, J.-K. Zhang, M. G. Amin, and K. M. Wong, "Nested cooperative encoding protocol for wireless networks with high energy efficiency," *IEEE Trans. Wireless Commun.*, vol. 7, no. 2, pp. 521–531, Feb. 2008.
- [2] A. De La Piedra, F. Benitez-Capistros, F. Dominguez, and A. Touhafi, "Wireless sensor networks for environmental research: A survey on limitations and challenges," in *Proc. IEEE Eurocon*, 2013, pp. 267–274.
- [3] R. Jafari, A. Encarnacao, A. Zahoor, F. Dabiri, H. Noshadi, and M. Sarrafzadeh, "Wireless sensor networks for health monitoring," in *Proc. IEEE 2nd Annu. Int. Conf. Mobile Ubiquitous Syst. Netw. Services*, 2005, pp. 479–481.
- [4] X. Lu, P. Wang, D. Niyato, D. I. Kim, and Z. Han, "Wireless charging technologies: Fundamentals, standards, and network applications," *IEEE Commun. Surveys Tuts.*, vol. 18, no. 2, pp. 1413–1452, 2nd Quart., 2015.
- [5] D. Kreutz, F. M. Ramos, P. E. Verissimo, C. E. Rothenberg, S. Azodolmolky, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proc. IEEE*, vol. 103, no. 1, pp. 14–76, Jan. 2015.
- [6] M. U. Younus, S. ul Islam, I. Ali, S. Khan, and M. K. Khan, "A survey on software defined networking enabled smart buildings: Architecture, challenges and use cases," *J. Netw. Comput. Appl.*, vol. 137, pp. 62–77, Jul. 2019.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 2018.
- [8] S. Bera, S. Misra, and A. V. Vasilakos, "Software-defined networking for Internet of Things: A survey," *IEEE Internet Things J.*, vol. 4, no. 6, pp. 1994–2008, Dec. 2017.
- [9] K. Sood, S. Yu, and Y. Xiang, "Software-defined wireless networking opportunities and challenges for Internet-of-Things: A review," *IEEE Internet Things J.*, vol. 3, no. 4, pp. 453–463, Aug. 2016.
- [10] M. U. Younus, M. K. Khan, M. R. Anjum, S. Afridi, Z. A. Arain, and A. A. Jamali, "Optimizing the lifetime of software defined wireless sensor network via reinforcement learning," *IEEE Access*, vol. 9, pp. 259–272, 2020.
- [11] T. Hu and Y. Fei, "QELAR: A machine-learning-based adaptive routing protocol for energy-efficient and lifetime-extended underwater sensor networks," *IEEE Trans. Mobile Comput.*, vol. 9, no. 6, pp. 796–809, Jun. 2010.



- [12] H. Yao, T. Mai, X. Xu, P. Zhang, M. Li, and Y. Liu, "NetworkAI: An intelligent network architecture for self-learning control strategies in software defined networks," *IEEE Internet Things J.*, vol. 5, no. 6, pp. 4319–4327, Dec. 2018.
- [13] Q. Zhu, R. Wang, Q. Chen, Y. Liu, and W. Qin, "IoT gateway: Bridging wireless sensor networks into Internet of Things," in *Proc. IEEE/IFIP Int. Conf. Embedded Ubiquitous Comput.*, 2010, pp. 347–352.
- [14] C. Alcaraz, P. Najera, J. Lopez, and R. Roman, "Wireless sensor networks and the Internet of Things: Do we need a complete integration?" in *Proc. 1st Int. Workshop Security Internet Things (SecIoT)*, 2010, pp. 1–9.
- [15] F. Junli, W. Yawen, and S. Haibin, "An improved energy-efficient routing algorithm in software define wireless sensor network," in *Proc. IEEE Int. Conf. Signal Process. Commun. Comput. (ICSPCC)*, 2017, pp. 1–5.
- [16] M. U. Younus, S. W. Kim, and S. U. Islam, "Proposition and real-time implementation of an energy-aware routing protocol for a software defined wireless sensor network," *Sensors*, vol. 19, no. 12, p. 2739, 2019.
- [17] G. Li, M. Dong, K. Ota, J. Wu, J. Li, and T. Ye, "Deep packet inspection based application-aware traffic control for software defined networks," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2016, pp. 1–6.
- [18] H. Huang, S. Guo, J. Wu, and J. Li, "Green datapath for TCAM-based software-defined networks," *IEEE Commun. Mag.*, vol. 54, no. 11, pp. 194–201, Nov. 2016.
- [19] S. Buzura, B. Iancu, V. Dadarlat, A. Peculea, and E. Cebuc, "Optimizations for energy efficiency in software-defined wireless sensor networks," *Sensors*, vol. 20, no. 17, p. 4779, 2020.
- [20] J. Wang, Y. Miao, P. Zhou, M. S. Hossain, and S. M. M. Rahman, "A software defined network routing in wireless multihop network," *J. Netw. Comput. Appl.*, vol. 85, pp. 76–83, May 2017.
- [21] A. Banerjee and D. Hussain, "SD-EAR: Energy aware routing in software defined wireless sensor networks," *Appl. Sci.*, vol. 8, no. 7, p. 1013, 2018.
- [22] W. Xiang, N. Wang, and Y. Zhou, "An energy-efficient routing algorithm for software-defined wireless sensor networks," *IEEE Sensors J.*, vol. 16, no. 20, pp. 7393–7400, Oct. 2016.
- [23] S. Fu, J. Wu, H. Wen, Y. Cai, and B. Wu, "Software defined wireline-wireless cross-networks: Framework, challenges, and prospects," *IEEE Commun. Mag.*, vol. 56, no. 8, pp. 145–151, Aug. 2018.
- [24] J. Wu, "Green wireless communications: From concept to reality [industry perspectives]," *IEEE Wireless Commun.*, vol. 19, no. 4, pp. 4–5, Aug. 2012.
- [25] W. Ejaz, M. Naeem, M. Basharat, A. Anpalagan, and S. Kandeepan, "Efficient wireless power transfer in software-defined wireless sensor networks," *IEEE Sensors J.*, vol. 16, no. 20, pp. 7409–7420, Oct. 2016.
- [26] M. Ndiaye, G. P. Hancke, and A. M. Abu-Mahfouz, "Software defined networking for improved wireless sensor network management: A survey," *Sensors*, vol. 17, no. 5, p. 1031, 2017.
- [27] M. R. Celenioglu, S. B. Goger, and H. A. Mantar, "An sdn-based energy-aware routing model for intra-domain networks," in *Proc. IEEE 22nd Int. Conf. Softw. Telecommun. Comput. Netw. (SoftCOM)*, 2014, pp. 61–66.
- [28] R. Wang, Z. Jiang, S. Gao, W. Yang, Y. Xia, and M. Zhu, "Energy-aware routing algorithms in software-defined networks," in *Proc. IEEE Int. Symp. World Wireless Mobile Multimedia Netw.*, 2014, pp. 1–6.
- [29] Y. Duan, W. Li, X. Fu, Y. Luo, and L. Yang, "A methodology for reliability of WSN based on software defined network in adaptive industrial environment," *IEEE/CAA J. Automatica Sinica*, vol. 5, no. 1, pp. 74–82, Jan. 2018.
- [30] G. Li, S. Guo, Y. Yang, and Y. Yang, "Traffic load minimization in software defined wireless sensor networks," *IEEE Internet Things J.*, vol. 5, no. 3, pp. 1370–1378, Jun. 2018.
- [31] J. Lorincz, A. Capone, and J. Wu, "Greener, energy-efficient and sustainable networks: State-of-the-art and new trends," *Sensors*, vol. 19, no. 22, p. 4864, 2019.
- [32] C. Estevez and J. Wu, "Recent advances in green Internet of Things," in *Proc. 7th IEEE Latin-Amer. Conf. Commun. (LATINCOM)*, 2015, pp. 1–5.
- [33] N. Gligorić, S. Krčo, D. Drajić, S. Jokić, and B. Jakovljević, "M2M device management in LTE networks," in *Proc. IEEE 19th Telecommun. Forum (TELFOR)*, 2011, pp. 414–417.
- [34] W. Guo, C. Yan, and T. Lu, "Optimizing the lifetime of wireless sensor networks via reinforcement-learning-based routing," *Int. J. Distrib. Sensor Netw.*, vol. 15, no. 2, pp. 1–9, 2019.
- [35] C. J. C. H. Watkins, "Learning from delayed rewards," Cambridge Univ., Cambridge, U.K., Rep., 1989.
- [36] R. Huang, X. Chu, J. Zhang, and Y. H. Hu, "Energy-efficient monitoring in software defined wireless sensor networks using reinforcement learning: A prototype," *Int. J. Distrib. Sensor Netw.*, vol. 11, no. 10, 2015, Art. no. 360428.
- [37] Z. Zhang, L. Ma, K. K. Leung, L. Tassiulas, and J. Tucker, "Q-placement: Reinforcement-learning-based service placement in software-defined networks," in *Proc. IEEE 38th Int. Conf. Distrib. Comput. Syst. (ICDCS)*, 2018, pp. 1527–1532.
- [38] A. Al-Jawad, P. Shah, O. Gemikonakli, and R. Trestian, "LearnQoS: A learning approach for optimizing QoS over multimedia-based SDNs," in *Proc. IEEE Int. Symp. Broadband Multimedia Syst. Broadcast. (BMSB)*, 2018, pp. 1–6.
- [39] S.-C. Lin, I. F. Akyildiz, P. Wang, and M. Luo, "QoS-aware adaptive routing in multi-layer hierarchical software defined networks: A reinforcement learning approach," in *Proc. IEEE Int. Conf. Services Comput. (SCC)*, 2016, pp. 25–33.
- [40] P. Sun, Y. Hu, J. Lan, L. Tian, and M. Chen, "Tide: Time-relevant deep reinforcement learning for routing optimization," *Future Gener. Comput. Syst.*, vol. 99, pp. 401–409, Oct. 2019.
- [41] J. Rischke, P. Sossalla, H. Salah, F. H. Fitzek, and M. Reisslein, "QR-SDN: Towards reinforcement learning states, actions, and rewards for direct flow routing in software-defined networks," *IEEE Access*, vol. 8, pp. 174773–174791, 2020.
- [42] C. Fang, C. Cheng, Z. Tang, and C. Li, "Research on routing algorithm based on reinforcement learning in SDN," *J. Phys. Conf.*, vol. 1284, no. 1, 2019, Art. no. 012053.
- [43] T.-Y. Mu, A. Al-Fuqaha, K. Shuaib, F. M. Sallabi, and J. Qadir, "SDN flow entry management using reinforcement learning," *ACM Trans. Auton. Adapt. Syst.*, vol. 13, no. 2, pp. 1–23, 2018.
- [44] H. Mostafaei and M. S. Obaidat, "Learning automaton-based self-protection algorithm for wireless sensor networks," *IET Netw.*, vol. 7, no. 5, pp. 353–361, 2018.
- [45] H. Song, J. Bai, Y. Yi, J. Wu, and L. Liu, "Artificial intelligence enabled Internet of Things: Network architecture and spectrum access," *IEEE Comput. Intell. Mag.*, vol. 15, no. 1, pp. 44–51, Feb. 2020.
- [46] G. Li, K. Ota, M. Dong, J. Wu, and J. Li, "DeSVig: Decentralized swift vigilance against adversarial attacks in industrial artificial intelligence systems," *IEEE Trans. Ind. Informat.*, vol. 16, no. 5, pp. 3267–3277, May 2020.
- [47] W. R. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "Energy-efficient communication protocol for wireless microsensor networks," in *Proc. IEEE 33rd Annu. Hawaii Int. Conf. Syst. Sci.*, 2000, p. 10.
- [48] P. Winter, "An algorithm for the enumeration of spanning trees," *BIT Numer. Math.*, vol. 26, no. 1, pp. 44–62, 1986.
- [49] M. Chakraborty, S. Chowdhury, J. Chakraborty, R. Mehera, and R. K. Pal, "Algorithms for generating all possible spanning trees of a simple undirected connected graph: An extensive review," *Complex Intell. Syst.*, vol. 5, no. 3, pp. 265–281, 2019.