ORIGINAL RESEARCH



A Reinforcement Learning Approach for Routing in Marine Communication Network of Fishing Vessels

Simi Surendran¹ · Alberto Montresor² · Maneesha Vinodini Ramesh³

Received: 30 May 2024 / Accepted: 6 December 2024 © The Author(s), under exclusive licence to Springer Nature Singapore Pte Ltd. 2025

Abstract

The lack of affordable communication facilities to the shore remains a fundamental problem for fishermen engaged in deepsea fishing. The Offshore Communication Network (OCN) is a wireless network of fishing vessels, whose goal is to provide Internet over the ocean. However, the dynamic nature of OCNs characterized by extreme weather, the difficulty of deploying additional infrastructure, wave-induced vessel movements, and high mobility causes significant challenges for traditional routing protocols. This paper proposes OCN-AR, a Q-learning-based adaptive routing strategy for ocean networks. The quality of the learning process relies on the reward function, which has been carefully designed to incorporate the most important features, including real-time forecasts of connectivity quality, path probability, link availability duration, and distance to the destination. The routing performance is evaluated through extensive simulations conducted under diverse conditions, including varying mobility scenarios, transmission rates, vessel rocking intensities, and node densities, and is compared against traditional protocols. The results demonstrate that OCN-AR significantly outperforms existing routing approaches, making it a reliable solution for maritime communication.

Keywords Marine communication network · Fishing vessel network · Routing protocol · Reinforcement learning

Introduction

The absence of adequate means to communicate with the shore is a fundamental impediment encountered by fishermen in deep-sea fishing. Conventional communication technologies such as cellular networks and marine radio can only provide connectivity up to a distance of about 20 km from the coast. These systems cannot be exploited since fishing vessels travel more than 100 km into the deep sea. On the other hand, satellite phone services are too expensive for fishermen. The Offshore Communication Network (OCN) proposed by Rao et al. aims at resolving this issue by

Simi Surendran simisurendran@am.amrita.edu

- ¹ Department of Computer Science and Engineering, Amrita School of Computing, Amrita Vishwa Vidyapeetham, Amritapuri, Kollam, India
- ² Department of Information Engineering and Computer Science, University of Trento, Trento, Italy
- ³ Amrita Center for Wireless Networks & Applications (AmritaWNA), Amrita School of Engineering, Amrita Vishwa Vidyapeetham, Amritapuri, Kollam, India

building a wireless network of fishing vessels that provide internet over the ocean [1]. The goal is to enable fishermen to connect to the Internet using affordable, handheld devices such as mobile phones, allowing them to communicate with the shore and other vessels.

Although marine networks share some of the features of terrestrial mobile ad-hoc and vehicular networks, they present unique characteristics and research challenges [2]. Unlike terrestrial vehicular networks, which are limited by road infrastructure, OCN benefits from greater freedom of movement in the open ocean. However, OCNs face significant communication challenges due to the inability to deploy additional infrastructure in the marine environment, the impact of extreme weather conditions on wireless signals, and the misalignment of directional links. The most critical challenge arises from the rough sea conditions, which can severely affect link quality. The topology of the network can change rapidly due to antenna orientation, the rocking movement of vessels, and propagation effects that weaken signal strength. For all the above reasons, providing uninterrupted internet connectivity is a difficult problem in OCNs.

Given these challenges, an adaptive routing mechanism capable of responding to the dynamic and unpredictable

nature of OCNs is essential for maintaining operational efficiency. Reactive routing protocols face difficulties in establishing reliable end-to-end paths due to issues in route discovery, while proactive protocols are resource-intensive, as they precompute all source-destination routes, including those that are not needed [3, 4]. Location-based routing protocols, which use predetermined parameters to select the next hop, fail to adapt quickly to the dynamic environment of OCNs. Integrating machine learning algorithms into the routing process continuously monitors the network conditions and adjusts the routing paths in real time. This adaptability is important in a marine environment where the links can change within short periods, making precomputed routes ineffective. Reinforcement learning (RL) techniques, in particular, have demonstrated potential, as agents interact with the wireless environment and make decisions based on feedback signals.

Although numerous routing protocols based on RL strategies have been proposed in the literature for terrestrial networks, they must be adapted to meet the unique requirements and challenges of OCNs [5, 6]. A well-designed reward function is essential for learning the optimal next-hop to deliver packets effectively under dynamic conditions. The contribution of this paper goes exactly in this direction. An OCN-specific reward function is proposed, taking into account the real-time connectivity quality estimate, the link availability duration, the path probability, and the geographical distance to the destination of neighbor nodes. Based on a Q-learning strategy, the proposed protocol discovers adaptive routes in a completely distributed manner. We evaluated the proposed protocol under various network conditions and compared its performance with other Q-learning-based routing schemes. In the simulation environment, the routing strategy achieved significantly higher packet delivery ratio.

The rest of the paper is organized as follows: Section Related Works reviews previous works related to routing protocols based on reinforcement learning. In Section Network Architecture, we present the architecture of OCN. Section RL Based Routing Model for OCN describes the reinforcement learning model and the formulation of the reward function. Section Results presents simulation results, followed by the concluding remarks of Sect. Conclusion and Future Scope.

Related Works

The routing problem in wireless networks has been extensively studied in the literature. In conventional *proactive* table-driven routing, algorithms predetermine all sourcedestination paths. However, in OCNs, the control traffic required to update neighborhood information is excessive, and nodes cannot predict in advance whether these paths will be needed for routing. In *reactive* routing, traditional protocols aim to discover a complete path between a source and target node, but gathering neighbor information in OCNs is challenging, often causing reactive protocols to fail in establishing a proper end-to-end path. Similarly, in *locationbased* routing, the selection of predefined parameters for the next hop cannot adapt quickly enough to the highly dynamic nature of the environment.

RL approaches have been successfully utilized to introduce intelligence in routing solutions across various wireless applications [6–9]. Mammeri has presented a review of routing protocols based on reinforcement learning [10]. Boyan et al. introduced the first Q-learning algorithm for routing in a telephone network [11]. While Q-routing outperforms non-adaptive algorithms, it may struggle to find the optimal policy under low or fluctuating network load conditions. To address this issue and improve convergence, several variants of Q-routing have been proposed [12]. Further, Routing schemes based in RL have been applied to diverse networks such as mobile ad-hoc network (MANET) [13, 14], vehicular adhoc networks (VANET) [15], wireless sensor networks (WSN) [16–18], wireless mesh networks (WMN) [19] and delay tolerant networks (DTN) [20] for performance improvement. RL-based routing has also been employed in other contexts, such as underwater sensor networks [21, 22], software-defined networks [23, 24], and information-centric networks [25].

Several model-free Q-learning based routing approaches have been proposed such as QLAODV [31], MQ routing [27], FROMS [16], SMART [32], QGrid [15], and Q-Geo [30]. Q-learning based routing for flying ad-hoc networks was discussed in [33, 34]. Lahsen-Cherif et al. proposed a Q-routing approach designed for wireless mesh networks that utilize directional antennas [35]. Meanwhile, Lee et al. introduced a multi-agent Q-learning framework specifically for UAV networks [36]. Wu et al. took a different approach by modifying the AODV protocol to develop QLAODV, a routing scheme for vehicular networks that considers vehicle movements and available channel bandwidth [31]. Nonetheless, QLAODV can encounter delays in collecting link information on multi-hop routes, which extends the time required for route adaptation. In another study, Macone et al. discussed a proactive routing technique called MQ-routing, aimed at improving node lifetime in mobile ad-hoc networks used in disaster relief scenarios [27]. This technique combines path availability and energy parameters in its optimization process. However, the proactive nature of MQ-routing makes it difficult to update the neighborhood table in rapidly changing environments. Although MQ-routing accounts for mobility and link availability in updating Q-values, it does not adequately consider link connectivity quality, which may result in high Q-values for nodes with low mobility despite poor connectivity.

Forster et al. developed the FROMs protocol, a multicast routing solution for wireless sensor networks [16]. The protocol SMART, proposed by Saleem et al., applies Q-learning in a cluster-based framework to enhance stable route selection and optimize secondary user performance in cognitive radio networks [32].

The Q-Grid protocol [15], a geographic routing scheme for vehicular networks, uses Q-learning techniques. However, it does not guarantee effective routing performance in networks with intermittent links, as it fails to utilize information on link status or availability duration. Jung et al. proposed another geographic routing protocol called Q-Geo for robotic networks based on Q-learning [30]. To ensure reliable data transfer, Q-Geo's reward function takes into account packet travel speeds, distance, link status, and location status. However, the protocol's fixed learning rate results in a uniform Q-value update rate across all network conditions. This static learning rate may not be effective in networks with time-varying and context-dependent topology changes. Table 1 summarizes the RL routing schemes applied in different types of networks along with the performance metrics.

Conventional routing approaches such as static, proactive, reactive, and geographic routing protocols can not adapt to rapid network changes and discover reliable paths. Although many adaptive RL routing algorithms have been suggested, all protocols mentioned above have been developed based on the distinctive features of each network. The learning parameters and reward function must be customized to meet the specific requirements of each network. Since OCNs operate in harsh environments and communication to the shore remains an essential factor, an adaptive model is necessary to enhance network connectivity. In previous work, we focused on packet status and signal strength when designing the reward function [37]. This paper revises the reward function by incorporating node-level connectivity quality, link availability, and distance to the destination, aiming to enhance the packet delivery ratio.

Network Architecture

The OCN is a fishing vessel network that provides Internet access well beyond 100 km from the shore. Its goal is to let fishermen access the internet through their mobile phones, using off-the-shelf apps such as WhatsApp for calls and messages. The fishing vessels forming the OCN act as edge nodes that locally analyze routing data and perform multiple network functions.

Based on the resources present in the fishing vessels, OCN nodes are classified into three types: access nodes, adaptive nodes, and supernodes.

- An *access node* is a fishing vessel that only provides a wireless access router (AR) and communicates using Wi-Fi omnidirectional antennas.
- An *adaptive node* is a vessel equipped with adaptive backhaul equipment (ABE) and one AR.
- A node equipped with two ABEs and one AR is categorized as a *supernode*.

Adaptive nodes and supernodes use 120° sectored longrange Wi-Fi links. These nodes are also termed *long range* (LR) nodes and are the ones forming the backbone network.

Table 1 Summary of RL-based routing protocols and performance metrics: p_1 :delay, p_2 : delivery ratio, p_3 : packet error rate, p_4 : Lifetime, p_5 : overhead reduction, p_6 : learning speed improvement, p_7 : path length reduction

Sl No	Routing Protocol	Application	Performance Parameters						
			p_1	p_2	p_3	p_4	p_5	p_6	p_7
1	Q-routing[11]	Static Network	х						
2	PQ-routing [12]	Static Network	х					х	
3	DRQ [26]	Static Network	х					х	
4	CRL [13]	MANET		х	х				
5	MQ routing [27]	MANET		х		х			
6	ARBR [20]	DTN	х	х					
7	FROMs [16]	WSN		х		х	х		
8	QELAR [28]	Underwater WSN				х			
9	QoE routing [29]	Multimedia Network	х		х				
10	DMARL [21]	Underwater WSN		х			х		
11	QGrid [15]	VANET		х					
12	Q-Geo [30]	Robotic Networks		х			х		
13	DCR [19]	WMN		х					
14	QL-AODV [31]	MANET		х		х			
15	SMART [32]	Cognitive Radio Networks		х			х		
16	OCN-AR	Extreme Networks		x				x	

LR Wi-Fi link's range is between 15 and 20 km. AR nodes with Wi-Fi links form a wireless mesh network.

The network architecture is shown in Fig. 1. A comprehensive description of the architecture and packet forwarding strategies is available in previous studies [1, 37–41]. The OCN architecture was validated in the Arabian Sea, starting from a coastal village in Kerala, India. For the sea trials, LR Wi-Fi equipment from Ubiquiti Networks and Cisco Linksys access routers were deployed. The onshore base station was positioned 56 ms above sea level, while the vessel's ABE was placed at 9 ms above sea level. During these sea trials, the network achieved a range of over 40 km for the first hop and 20 km for each subsequent hop.

RL Based Routing Model for OCN

We consider an OCN routing scenario comprising LR nodes, access nodes, and onshore base stations, where all nodes except onshore base stations are mobile. LR nodes have a transmission range of 15–20 km, while access nodes provide connectivity within a radius of 150–250 ms. Multi-hop communication is required to ensure messages reach their final destination. Each node is assumed to be equipped with GPS devices that provide real-time positional data. The nodes generate messages of varying priorities, including emergency, audio, and video messages. Furthermore, all nodes periodically send beacons to neighboring nodes, containing location details and other routing-related information.

Routing in OCN is modeled as an RL problem. In this framework, the entire network, comprising access nodes, LR nodes, and base stations, represents the environment for the RL agent. Each packet within the network functions as an independent agent. The agent receives rewards or penalties based on the success or failure of transmissions, using these feedback signals to learn an optimal policy for selecting the most suitable next hop. Each node in the network maintains information about its neighbors' connectivity features, such as connectivity quality, link availability duration, path probability, distance to the destination, and Q-value. These features are periodically updated through beacon messages exchanged between nodes. To determine the next hop for packet forwarding, we employ a temporal difference offpolicy Q-learning algorithm [42]. This algorithm uses an ϵ -greedy policy, where ϵ represents a small probability of introducing randomness into action selection, allowing the agent to explore new potential routing solutions. As wireless links fluctuate over time, this approach ensures that new links can be incorporated into the routing path, making it adaptable to the dynamically changing environment of OCNs.

Figure 2 depicts the routing scenario in OCN-AR. In this example, node *S* needs to establish communication with the base station and must choose an appropriate next-hop neighbor. Node *S* has three possible options for backbone connections: n_1 , n_2 , and n_3 . The actions of selecting these neighbors are denoted by a_{n_1} , a_{n_2} , and a_{n_3} , respectively. The node will choose the action with the highest Q-value with a probability



Fig. 1 Architecture of Offshore Communication Network[1]

SN Computer Science

Fig. 2 State-Action model in OCN routing: Node *S* have three possible actions: a_{n_1} , a_{n_2} , a_{n_3} . It selects the best action a_{n_1} and forwards the packet to node n_1 . Node *S* will receive a reward for action a_{n_1} and update its Q-value accordingly. Similarly, the best action will be selected from node n_1 and this process will continue until the packet reaches its destination 62



of $1 - \epsilon$, for instance, a_{n_1} . Upon executing this action, node *S* receives a reward that is calculated using both local and remote information. At node n_1 , the process repeats, with the next hop being selected based on the highest Q-value. In this scenario, the base station has the highest Q-value, leading node n_1 to choose that path. Each node maintains a Q-table that records the Q-values of its neighbors, which are updated after feedback is received. These updated Q-values are used in future packet forwarding decisions, enabling the selection of the most optimal route over time.

A Markov Decision Process with the following state, action, and reward function is used to model the routing process.

- *States*: Each packet in the network corresponds to an agent, and the agent's environment is defined as the collection of all nodes in the network. The state of an agent is represented by the node where the packet is currently located. Specifically, if a node *i* generates a new packet or receives one for forwarding, the agent's state is represented as *i*. Accordingly, the state space consists of all nodes in the network, along with their respective features.
- Actions: Let $N_i = \{n_1, n_2, ..., n_k\}$ be the neighbors of node *i*. An action in node *i* at time *t* is the selection of one of the neighbors from N_i for forwarding packet. The action space comprises the set of possible actions that all nodes can take in an OCN.
- *Reward function*: The reward function is the most important factor in determining the effectiveness of Q-routing. The environment returns a reward to the agent for indi-

cating the impact of the neighborhood selection. This reward comprises a function of local information computed in the node and a component of remote information received as feedback from other nodes, obtained through the acknowledgment scheme. The parameters to be considered at the local (node) level include the connectivity quality of neighboring nodes, the link availability duration, the probability of the path to the destination, and the distance to the destination.

The subsequent subsections provide a detailed explanation of the first three components of the proposed reward function, followed by the complete reward function and the description of the Q-learning algorithm.

Node Connectivity Quality

To analyze the characteristics of marine wireless links, we collected data on signal strength variations with distance during various sea trials [41]. This dataset is employed to examine the factors affecting connectivity within oceanic networks. Sea wave-induced movements and propagation effects significantly influence signal quality in OCN, causing frequent topological changes and posing challenges for packet routing. To address these issues, we developed a machine learning framework that utilizes both historical and real-time data to predict link connectivity probabilities.

A metric called *dynamic connectivity index (DCI)* has been defined to compute the level of node connectivity by employing this link prediction model [43]. Assume that connectivity is determined solely based on the local link quality. In this case, there is a possibility of selecting a neighbor with good link quality. However, if the chosen neighbor has poor overall connectivity, it is less likely to successfully deliver the message. *DCI* is a node-level measure of a node's connectivity to clusters and base stations. *DCI* helps to decide the most suitable next-hop and minimizes the possibility of node isolation. To define the *DCI* of a node, we compute the *DCI* of its next-hop neighbors and link probability to those neighbors as shown in Eq. 1.

$$DCI(x) = \sum_{i \in AN(x)} w_i \cdot p_i \cdot DCI(i)$$
⁽¹⁾

where w_i is the weight of next-hop *i*, p_i is the link probability from node *x* to neighbor *i*, AN(x) is the list of neighbor nodes whose link probability is greater than a threshold value 0.25 and DCI(i) is the dynamic connectivity index of neighbor *i*. A dynamic weighting scheme is used to prioritize the neighbors of a node [43].

Link Availability Duration

Consider a communication scenario between two mobile LR nodes as shown in Fig. 3. Let nodes S and A be separated by a distance d^{τ} after time τ . Let v_s and v_a be the velocity vectors of nodes S and A. The alignment of directional antennas between the LR nodes affects the communication radius. Let ϕ represent the angle of misalignment among the transmitter and receiver antennas. An increase in ϕ decreases the communication distance between nodes and varies with the environments' dynamics. z_{eff}^{τ} represents the effective communication distance of S and A after time τ as a function of ϕ .



Fig. 3 Communication scenario between two mobile LR nodes with movement vectors v_s and v_a

SN Computer Science A Springer Nature journal The displacement in x, y directions and the distance between nodes due to mobility can be computed as:

• Difference in x direction δ_x^{τ} :

$$\delta_x^{\tau} = (x_s - x_a) + (v_a \cos \theta_a - v_s \cos \theta_s)\tau$$
(2)

• Difference in y direction δ_{y}^{τ} :

$$\delta_{y}^{\tau} = (y_{s} - y_{a}) + (v_{a}\sin\theta_{s} - v_{s}\sin\theta_{s})\tau$$
(3)

• Distance between nodes after time *τ*:

$$d^{\tau} = \sqrt{\left(\delta_x^{\tau}\right)^2 + \left(\delta_y^{\tau}\right)^2}$$

= $\left[\left((x_s - x_a) + (v_a \cos \theta_a - v_s \cos \theta_s)\tau\right)^2 + \left((y_s - y_a) + (v_a \sin \theta_a - v_s \sin \theta_s)\tau\right)^2\right]^{1/2}$

For effective communication, $d^{\tau} < z_{eff}^{\tau}$. The maximum duration of link availability τ can be obtained by rewriting this equation as:

$$\tau = \frac{-(\Delta_x V_{x\delta} + \Delta_y V_{y\delta}) \pm \sqrt{(V_{x\delta}^2 + V_{y\delta}^2)(z_{eff}^{\tau})^2 - (\Delta_y V_{x\delta} - \Delta_x V_{y\delta})^2}}{V_{x\delta}^2 + V_{y\delta}^2}$$
(4)

where

$$\begin{aligned} \Delta_x &= (x_a - x_s), \\ \Delta_y &= (y_a - y_s), \\ V_{x\delta} &= v_a \cos(\theta_a + \phi_a) - v_s \cos(\theta_s + \phi_s), \\ V_{y\delta} &= v_a \sin(\theta_a + \phi_a) - v_s \sin(\theta_s + \phi_s) \end{aligned}$$

Path Probability

Information regarding the possibility of an end-to-end path between the LR nodes or the base station helps to select next-hop neighbors. There may be multiple paths with different hops and various link properties between two nodes. The connectivity probabilities of these links can be predicted using the machine learning framework in OCN. Additionally, the details of existing paths are available through the feedback mechanism of the routing algorithm. We consider all such paths and link probabilities to estimate the connectivity probability between nodes.

Let $r_1, r_2, ..., r_q$ be the routing paths exist between two nodes a and b. Consider a path r_i that consists of m links, whose availability is expressed by probabilities $l_1, l_2, ..., l_m$. Assume that such probabilities are independent of each other. Then,

the probability of the existence of the path r_i can be calculated by the product of its individual link probabilities as in Eq. 5.

$$Prb(r_i) = l_1 \cdot l_2 \cdot \ldots \cdot l_m \tag{5}$$

If any one of the links in this path breaks, then $Prb(r_i)$ becomes zero. Thus, the probability that the path w_i does not exist is represented by $1 - Prb(r_i)$.

Assuming that the chances of multiple paths between two nodes are independent, and given q potential paths connecting nodes a and b, the probability of having at least one path is

$$\mathcal{P}(a,b) = 1 - \prod_{i=1}^{q} \left(1 - Prb(r_i) \right)$$
(6)

If a neighboring node has a high probability of leading to the destination, it will attain a higher Q-value. However, in scenarios where many links experience failure, calculating this path probability becomes difficult. In such cases, the reward function excludes this factor.

Reward Function

The reward for forwarding a packet from node i to node j with action a_i at time t is defined as:

$$\mathcal{R}_{t} = \mathcal{R}_{f} + \beta_{1} \cdot DCI(j) + \beta_{2} \cdot \tau + \beta_{3} \cdot \mathcal{P}(j, dest) + \beta_{4} \cdot \frac{\Delta d(i, j)}{d(i, j)}$$
(7)

where DCI(j) is the connectivity index of node j, τ is the link availability time, $\mathcal{P}(j, dest)$ is the probability that a path exists from node j to the destination and $\Delta d(i, j)$ is the difference in distance between i and j to the destination. β_i s are the weights given to each of the node-level features. Since DCI(j) and τ are the most important features, we used a normalized weight vector $\beta = \langle 0.5, 0.25, 0.1, 0.15 \rangle$ in simulations. \mathcal{R}_f is defined as

$$\mathcal{R}_{f} = \begin{cases} +5, & \text{ACK received} \\ -5, & \text{ACK not received} \\ -10, & \text{next-hop is local maximum} \\ R_{max}, & & \text{next-hop is destination} \end{cases}$$
(8)

The node will be given a positive or negative reward based on whether or not it successfully receives the packet. Furthermore, if the chosen neighbor lacks a viable connection for forwarding, a larger negative value is assigned. Also, the maximum reward will be given when the next hop from the chosen neighbor leads directly to the destination. We set R_{max} to 100 in our simulation experiments.

OCN Adaptive Routing Algorithm

The OCN Adaptive Routing (OCN-AR) algorithm is shown in Algorithm 1. The algorithm initializes the network's state *S* and possible actions *A*. It uses an exploration coefficient $\epsilon = 0.01$ to balance exploration and exploitation in decision-making. For each node $s \in S$, the Q-value is initialized using a function based on *DCI* and link probability \mathcal{P} . Every node creates and updates a list of Q-values for its adjacent nodes. In cases where an entry is absent in the Q-table, the source node generates a new entry by utilizing details from both the target and the neighboring node. Nodes periodically broadcast beacons containing their Q-values, the current *DCI(i)*, and location information *loc*. When a node receives a beacon or acknowledgment message, it updates the reward \mathcal{R}_t and the Q-value using the Eq. 9.

$$\mathcal{Q}^{new}(s_t, a_t) = (1 - \alpha)\mathcal{Q}^{old}(s_t, a_t) + \alpha [\mathcal{R}_t + \gamma \max_a \mathcal{Q}((s_{t+1}, a))]$$
(9)

where $Q^{old}(s_t, a_t)$ is the old Q-value for the action a_t in state s_t . \mathcal{R}_t denotes the reward for the action a_t in state s_t . The maximum Q-value in the next state s_{t+1} with the best action a is denoted by $\max_a Q(s_{t+1}, a)$. The learning rate and the discount factor are represented as α and γ , respectively, with values in the range [0, 1]. After creating the Q-table, the nodes select a neighbor with the highest Q as the next-hop with probability $1 - \epsilon$. This phase is the exploitation stage of Q-learning. For the exploration of the state space, the source node arbitrarily selects any of the neighbors as next-hop with probability ϵ . LR nodes select the best neighbor only from the list of LR nodes as the transmission radius of these nodes is very large compared to access nodes. Access nodes can choose multiple hops to reach an LR node.

Algorithm 1 OCN Adaptive Routing

```
State S = Set of all nodes in the network:
Action A = Neighbor nodes ;
Exploration coefficient \epsilon = 0.01;
for every s \in S do
    for every a \in A do
         \mathcal{N} = neighbors of node s;
         Q(s_t, a_t) = DCI(\mathcal{N}_{a_t}) * \mathcal{P};
while true do
    if beacon timer expires in each node x then
      Send beacon < Q - value, DCI(x), loc >
    if receive a beacon/ack message then
        Update reward function \mathcal{R}_t and Q-value Q(s_t, a_t) as in Eq. 7, 8 and 9
    if Node x has a packet to send then
         Generate a random number p \in [0, 1];
         if \epsilon \leq p then
             next-hop = random(N)
         else
             if node type(x) = access then
                 next-hop = \operatorname{argmax}_{a \in \mathcal{N}} \{ a : Q(s, a) \}
             if node type(x) = LR then
                  next-hop = \operatorname{argmax}_{a \in \mathcal{N}(\mathcal{LR})} \{ a : Q(s, a)) \}
```

The hyperparameters used in this study are the learning rate, the discount factor and the exploration coefficient. The learning rate determines the weight assigned to new information from the environment and existing learned information. The degree of topology changes differs based on sea conditions and stages of fishing. In the roughest sea conditions, the rocking motion is substantial, resulting in frequent topology changes. While an increased learning rate enables quick adaptation to these changes, a rate that is too high can cause instability in the learning process. To address this, the learning rate is initialized within the range of 0.5 to 0.7, and an adaptive approach is proposed to manage connectivity fluctuations caused by the rocking movements of vessels. Since the degree of mobility varies between different fishing stages, a context-sensitive selection of the learning rate is applied to optimize performance for each scenario.

The discount factor is another hyperparameter to balance immediate rewards with long-term rewards. A higher discount factor makes the agent prioritize long-term routing strategies. However, in specific sea conditions, such as emergency communications or sudden environmental disruptions, shorter-term paths are often required. The discount factor is selected on the basis of the situational requirements to address these challenges effectively. The exploration coefficient in the ϵ -greedy policy is used to balance exploring

SN Computer Science A Springer Nature journal new actions and exploiting actions known to provide high rewards based on current Q-value estimates. In this study, a high exploration rate is used at the beginning to allow the agent to discover various experiences in a dynamic environment. Over time, this rate gradually reduces, allowing the agent to take advantage of the learned policies.

Results

The experimental setup to evaluate the OCN-AR protocol involved both simulations and real-world experiments to collect signal strength data under different vessel rocking conditions. A machine learning framework was then developed to predict the strength of the signal under varying sea conditions. This prediction model is used in the computation of the probability of link and *DCI*, which were the important components of the reward function in the OCN-AR protocol. In the simulation phase, NS-2 was used to model the network, which included 50 mobile nodes, comprising 12 longrange nodes and 38 access nodes, along with one onshore base station. Packet generation rates varied from 2 to 64 packets per second, with a default packet size of 512 bytes, and traffic was generated using a fixed-rate UDP source. The protocol's performance was assessed based on the packet



Fig. 4 a Data collected from sea-trial experiments on signal strength variation with distance in different vessel rocking stages. b Prediction of signal strength at distance 25 km and vessel rocking stage = 3



Fig.5 Comparison of packet delivery ratio for varying mobility degree

delivery ratio and compared with three state-of-the-art routing protocols: AODV, GPSR, and Q-Geo.

The signal strength data collected during marine experiments from various vessel rocking stages are shown in Fig. 4a. Here, each state represents a different vessel rocking stage. The sea condition is rough in higher-numbered states, and we can observe more signal deviation due to the increased impact of waves on node mobility. We developed a machine learning framework to predict the signal strength under different sea conditions using this data. Figure 4b shows a sample prediction of signal strength over a distance of 25 km. Employing this prediction model, link probability and *DCI*, which are components of the reward function, are computed.

Although the average speed of fishing vessels is 3–5 m/s, the rocking motion leads to more packet drops. This moving effect of the vessels is simulated by setting the node pause time. The shorter the pause time, the higher the degree of mobility. Figure 5 shows the packet delivery ratio for various pause times. Link failures often occur when mobility levels increase, causing AODV and GPSR performance to decline. OCN-AR shows significant improvement over Q-Geo and other protocols, taking into account link availability length, direction of movement from the destination, and neighborhood connectivity level.

We varied the transmission rate at the source nodes from 2 to 64 packets per second. The results, shown in Fig. 6a, reveal a significant drop in the packet delivery ratio when the transmission rate exceeds 16 packets per second across all protocols. Among the protocols, AODV experiences a faster degradation in performance compared to the other two adaptive protocols, primarily due to higher channel occupancy and increased packet loss from link failures. In contrast, OCN-AR performs better than Q-Geo as it considers the connectivity quality of the nodes.

One of the distinctive features of OCN is the rocking movement of the nodes, which often causes link breakages. We simulated this scenario using a random waypoint



Fig. 6 a Variation of packet delivery ratio with data transmission rate. b Variation of packet delivery ratio with rocking movement vessels due to sea waves. The rocking degree indicates the wave state of the sea



Fig. 7 a Variation of packet delivery ratio with the density of LR nodes. b Variation of packet delivery ratio with increase in learning rate

mobility model with varying velocities. We dropped a fixed percentage of packets from the source in each rocking stage to simulate the vessel movement factor. In this case, the packet delivery ratio of OCN-AR was compared to AODV and Q-Geo as shown in Fig. 6b. The performance of AODV decreases considerably with increasing rocking degrees because of the large number of link failures. OCN-AR performs effectively in high rocking conditions compared to other protocols. This improvement is because it utilizes signal strength data and estimates of link availability. The simulation scenario includes both access and LR nodes.

We tested the impact of the density of LR nodes on the packet delivery ratio and the results demonstrated in Fig. 7a. The number of LR nodes varied from 2 to 16 and we noticed an improvement in packet delivery with increased density. The learning rate is one of the critical hyper-parameters in OCN-AR that controls the rate of adaptation to the dynamic topology. Different fishing stages experience varying

SN Computer Science

mobility degrees, and hence a context-dependent learning rate is used. To tune the parameter, we varied it from 0.1 to 0.9 and observed a better performance at value 0.7 in the fish searching context, as shown in Fig. 7b. The optimal learning rate computed in other fishing stages is slightly different from this context.

Since RL-based routing continuously learns and adapts to the environment, OCN-AR performs better than traditional routing approaches. RL agents effectively explore and exploit routing strategies, to maintain reliable communication despite changes in network topology. When compared to other RL-based approaches, OCN-AR improves the performance in high-mobility scenarios due to its reward function, which incorporates metrics such as the connectivity index to assess node reliability, link availability time to prioritize stable links, path existence probability to evaluate the likelihood of successful packet delivery, and distance difference to optimize route efficiency. These elements enable OCN-AR to adapt to abrupt connectivity fluctuations caused by vessel mobility, ensuring robust and efficient routing decisions in real-time. Frequent link failures in high-mobility conditions degrade the performance of traditional protocols like AODV and GPSR. AODV suffers from faster performance degradation due to its dependence on higher channel occupancy and increased packet loss resulting from frequent disconnections. As the rocking degree increases, the performance of AODV declines significantly due to the growing number of link failures. In contrast, OCN-AR demonstrates notable improvements over AODV and Q-Geo by utilizing metrics link stability, movement direction toward the destination, and neighbourhood connectivity levels. These features enable OCN-AR to make better routing decisions. By prioritizing stable links and well-connected nodes, OCN-AR effectively minimizes packet loss.

Conclusion and Future Scope

In this paper, we addressed the routing challenges in an offshore communication network of fishing vessels and introduced OCN-AR, a Q-learning-based routing protocol designed for effective message dissemination. OCN-AR utilized a reward function that incorporates important features- real-time connectivity quality forecasts, path probability, link availability duration, and distance to the destination of OCN. The effectiveness of this routing strategy was validated through simulations using data from sea trials, which demonstrated that the reinforcement learning strategy enhances routing performance under the challenging conditions of the ocean. The packet delivery ratio of OCN-AR was evaluated across various mobility scenarios, transmission rates, vessel rocking factors, and node densities, and compared with existing protocols. The results indicate that OCN-AR is well-suited for maritime networks of fishing vessels. However, the protocol faces challenges, particularly with the high computational overhead required for updating Q-values.

The scalability of RL-based protocols is a challenge in large-scale, real-world networks, particularly when it comes to resource-constrained deployments. Real-world offshore communication networks face limitations in onboard processing power, energy availability, and network bandwidth. One of the key challenges of the current approach is the computational complexity associated with Q-value updates, especially in large-scale networks. To address this, future work could focus on developing lightweight algorithms and employing value function approximation techniques, such as deep Q-networks, to reduce the computational overhead. In the context of OCNs, RL agents in multiple vessels can not only learn from their individual experiences but also share their learned experiences with neighbouring vessels' RL agents to improve routing decisions. In the future, we plan to incorporate collaborative learning techniques among different node types within the OCN to improve network performance.

Funding We wish to confirm that there has been no significant financial support for this work that could have influenced its outcome.

Declarations

Conflict of Interest On behalf of all authors, the corresponding author states that there is no Conflict of interest.

References

- Rao SN, Ramesh MV, Rangan V. Mobile infrastructure for coastal region offshore communications and networks. In: Proc. of the IEEE Global Humanitarian Technology Conf. (GHTC);2016. pp. 99–104. IEEE.
- Yau K-LA, Syed AR, Hashim W, Qadir J, Wu C, Hassan N. Maritime networking: bringing internet to the sea. IEEE Access. 2019;7:48236–55.
- Shrivastava PK, Vishwamitra L. Comparative analysis of proactive and reactive routing protocols in vanet environment. Meas Sens. 2021;16: 100051.
- Kim B-S, Ullah S, Kim KH, Roh B, Ham J-H, Kim K-I. An enhanced geographical routing protocol based on multi-criteria decision making method in mobile ad-hoc networks. Ad Hoc Netw. 2020;103: 102157.
- Nazib RA, Moh S. Reinforcement learning-based routing protocols for vehicular ad hoc networks: a comparative survey. IEEE Access. 2021;9:27552–87.
- Okine AA, Adam N, Naeem F, Kaddoum G. Multi-agent deep reinforcement learning for packet routing in tactical mobile sensor networks. IEEE Trans Netw Serv Manag. 2024;21:2155–2169.
- Prabhu D, Alageswaran R, Miruna Joe Amali S. Multiple agent based reinforcement learning for energy efficient routing in wsn. Wirel Netw. 2023;29(4):1787–97.

- 8. Nandyala CS, Kim H-W, Cho H-S. Qtar: a q-learning-based topology-aware routing protocol for underwater wireless sensor networks. Comput Netw. 2023;222:109562.
- 9 Yang X, Yan J, Wang D, Xu Y, Hua G. Woad3qn-rp: an intelligent routing protocol in wireless sensor networks-a swarm intelligence and deep reinforcement learning based approach. Expert Syst Appl. 2024;246: 123089.
- 10. Mammeri Z. Reinforcement learning based routing in networks: review and classification of approaches. IEEE Access. 2019:7:55916-50.
- 11. Boyan JA, Littman ML. Packet routing in dynamically changing networks: a reinforcement learning approach. In: Advances in neural information processing systems; 1994. pp. 671-678. Citeseer.
- 12. Choi SP, Yeung D.-Y. Predictive q-routing: a memory-based reinforcement learning approach to adaptive traffic control. In: Advances in neural information processing systems; 1996. pp. 945-951.
- 13. Dowling J, Curran E, Cunningham R, Cahill V. Using feedback in collaborative reinforcement learning to adaptively optimize manet routing. IEEE Trans Syst Man Cybern-Part A Syst Humans. 2005;35(3):360-72.
- 14. Liu J, Wang Q, He C, Jaffrès-Runser K, Xu Y, Li Z, Xu Y. Qmr: O-learning based multi-objective optimization routing protocol for flying ad hoc networks. Comput Commun. 2020;150:304-16.
- 15. Li R, Li F, Li X, Wang Y. Qgrid: Q-learning based routing protocol for vehicular ad hoc networks. In: 2014 IEEE 33rd Int. Performance Computing and Communications Conf. (IPCCC); 2014. pp. 1-8.IEEE.
- 16. Förster A, Murphy AL. Froms: a failure tolerant and mobility enabled multicast routing paradigm with reinforcement learning for wsns. Ad Hoc Netw. 2011;9(5):940-65.
- 17. Lu Y, He R, Chen X, Lin B, Yu C. Energy-efficient depth-based opportunistic routing with q-learning for underwater wireless sensor networks. Sensors. 2020;20(4):1025.
- 18. Zhu R, Jiang Q, Huang X, Li D, Yang Q. A reinforcement-learning-based opportunistic routing protocol for energy-efficient and void-avoided uasns. IEEE Sens J. 2022;22(13):13589-601.
- 19. Chai Y, Zeng X-J. A multi-objective dyna-q based routing in wireless mesh network. Appl Soft Comput. 2021;108: 107486.
- 20. Elwhishi A, Ho P-H, Naik K, Shihada B. Arbr: Adaptive reinforcement-based routing for dtn. In: 2010 IEEE 6th Int. Conf. on Wireless and Mobile Computing, Networking and Communications; 2010. pp. 376-385.IEEE.
- 21. Li X, Hu X, Zhang R, Yang L. Routing protocol design for underwater optical wireless sensor networks: A multiagent reinforcement learning approach. IEEE Internet Things J. 2020;7(10):9805-18.
- 22. Zhang Y, Zhang Z, Chen L, Wang X. Reinforcement learningbased opportunistic routing protocol for underwater acoustic sensor networks. IEEE Trans Veh Technol. 2021;70(3):2756-70.
- 23. Lin S-C, Akyildiz IF, Wang P, Luo M. Qos-aware adaptive routing in multi-layer hierarchical software defined networks: a reinforcement learning approach. In: 2016 IEEE Int. Conf. on Services Computing (SCC); 2016. pp. 25-33. IEEE.
- 24. Chen Y-R, Rezapour A, Tzeng W-G, Tsai S-C. Rl-routing: an sdn routing algorithm based on deep reinforcement learning. IEEE Trans Netw Sci Eng. 2020;7(4):3185-99.
- 25. Ye D, Zhang M, Yang Y. A multi-agent framework for packet routing in wireless sensor networks. Sensors. 2015;15(5):10026-47.
- 26. Kumar S, Miikkulainen R. Dual reinforcement q-routing: An online adaptive routing algorithm. In: Proc. of the Artificial Neural Networks in Engineering Conf.; 1997. pp. 231-238. Citeseer.
- 27. Macone D, Oddi G, Pietrabissa A. Mq-routing: mobility-, gpsand energy-aware routing protocol in manets for disaster relief scenarios. Ad Hoc Netw. 2013;11(3):861-78.

- 28. Hu T, Fei Y. Qelar: a machine-learning-based adaptive routing protocol for energy-efficient and lifetime-extended underwater sensor networks. IEEE Trans Mob Comput. 2010;9(6):796-809.
- 29. Coutinho N, Matos R, Marques C, Reis A, Sargento S, Chakareski J. Kassler A. Dynamic dual-reinforcement-learning routing strategies for quality of experience-aware wireless mesh networking. Comput Netw. 2015;88:269-85.
- 30. Jung W-S, Yim J, Ko Y-B. Qgeo: Q-learning-based geographic ad hoc routing protocol for unmanned robotic networks. IEEE Commun Lett. 2017;21(10):2258-61.
- 31. Wu C, Kumekawa K, Kato T. Distributed reinforcement learning approach for vehicular ad hoc networks. IEICE Trans Commun. 2010;93(6):1431-42.
- 32 Saleem Y, Yau K-LA, Mohamad H, Ramli N, Rehmani MH. Smart: a spectrum-aware cluster-based routing scheme for distributed cognitive radio networks. Comput Netw. 2015;91:196-224.
- 33. Costa LAL, Kunst R, Freitas EP. Q-fanet: improved q-learning based routing protocol for fanets. Comput Netw. 2021;198: 108379.
- 34. Arafat MY, Moh S. A q-learning-based topology-aware routing protocol for flying ad hoc networks. IEEE Internet of Things J. 2021;9(3):1985-2000.
- 35. Lahsen-Cherif I, Zitoune L, Vèque V. Energy efficient routing for wireless mesh networks with directional antennas: When q-learning meets ant systems. Ad Hoc Netw. 2021;121: 102589.
- 36. Lee S, Yu H, Lee H. Multi-agent q-learning based multi-uav wireless networks for maximizing energy efficiency: deployment and power control strategy design. IEEE Internet of Things J. 2021:9:6434-6442.
- 37. Simi S, Ramesh MV. A reinforcement learning approach for improving internet connectivity in maritime network. J Adv Res Dyn Control Syst. 2018;10:598-604.
- 38 Rao SN, Raj D, Aiswarya S, Unni S. Realizing cost-effective marine internet for fishermen. In: 2016 14th Int. Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt); 2016. pp. 1-5. IEEE.
- Unni S, Raj D, Sasidhar K, Rao S. Performance measurement and 39. analysis of long range wi-fi network for over-the-sea communication. In: 2015 13th Int. Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt); 2015. pp. 36-41. IEEE.
- 40. Dhivvya J, Rao SN, Simi S. Towards maximizing throughput and coverage of a novel heterogeneous maritime communication network. In: Proc. of the 18th ACM Int. Symposium on Mobile Ad Hoc Networking and Computing; 2017. p. 39. ACM.
- 41. Surendran S, Ramesh MV, Montresor A, Montag MJ. Link characterization and edge-centric predictive modeling in an ocean network. IEEE Access. 2023;11:5031-46.
- Watkins CJ, Dayan P. Q-learning Mach Learn. 1992;8(3):279-92. 42.
- 43. Surendran S, Ramesh MV, Montag MJ, Montresor A. Modelling communication capability and node reorientation in offshore communication network. Comput Electr Eng. 2020;87: 106781.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

SN Computer Science

A SPRINGER NATURE journal