

Energy Efficiency Optimization for Subterranean LoRaWAN Using A Reinforcement Learning Approach: A Direct-to-Satellite Scenario

Kaiqiang Lin, *Student Member, IEEE*, Muhammad Asad Ullah, *Student Member, IEEE*,
Hirley Alves, *Member, IEEE*, Konstantin Mikhaylov, *Senior Member, IEEE*, and Tong Hao, *Member, IEEE*

Abstract—The integration of subterranean LoRaWAN and non-terrestrial networks (NTN) delivers substantial economic and societal benefits in remote agriculture and disaster rescue operations. The LoRa modulation leverages quasi-orthogonal spreading factors (SFs) to optimize data rates, airtime, coverage and energy consumption. However, it is still challenging to effectively assign SFs to end devices for minimizing co-SF interference in massive subterranean LoRaWAN NTN. To address this, we investigate a reinforcement learning (RL)-based SFs allocation scheme to optimize the system’s energy efficiency (EE). To efficiently capture the device-to-environment interactions in dense networks, we proposed an SFs allocation technique using the multi-agent dueling double deep Q-network (MAD3QN) and the multi-agent advantage actor-critic (MAA2C) algorithms based on an analytical reward mechanism. Our proposed RL-based SFs allocation approach evinces better performance compared to four benchmarks in the extreme underground direct-to-satellite scenario. Remarkably, MAD3QN shows promising potentials in surpassing MAA2C in terms of convergence rate and EE.

Index Terms—Subterranean LoRaWAN, non-terrestrial networks, reinforcement learning, SFs allocation, energy efficiency.

I. INTRODUCTION

THE integration of LoRaWAN-based wireless underground sensor networks and non-terrestrial networks (NTN) enables subterranean massive machine-type communications (mMTC) applications to operate in hard-to-reach or disaster rescue areas [1]. LoRa, a chirp spread spectrum modulation variation in LoRaWAN, introduces quasi-orthogonality between packets with different spreading factors (SFs) [2]. This characteristic grants LoRa its resistance to interference while offering a range of trade-offs between time-on-air (ToA), radio coverage, and energy consumption through varying SF levels. However, in subterranean mMTC scenarios, the Aloha-like media access protocol used in LoRa constrains the network capacity and collision robustness. For instance, the simulation results reported in [1], [3] illustrate a relatively low probability of successful packet delivery when a large number of end devices (EDs) are assigned with the same SF in the underground direct-to-satellite (U-DtS) connectivity. This is attributed to the frequent co-SF interference that occurs when packets featuring the same SF are simultaneously transmitted on the same channel.

To leverage LoRa quasi-orthogonality, several studies have discussed the SFs allocation techniques and evaluated the scal-

ability for terrestrial networks. Specifically, two notable one-time spatial SFs allocation schemes, namely equal-interval-based (EIB) [4] and equal-area-based (EAB) [5], were proposed to mitigate the co-SF interference and improve the packet delivery ratio. To adjust to dynamic underground environments, adaptive parameters assignment schemes were proposed in LoRaWAN, such as the adaptive data rate mechanism specified by the LoRa Alliance [6] and the path-loss-based (PLB) scheme proposed in [7]. Nevertheless, both solutions overlook the co-SF interference. Thus, their performance diminishes in practical subterranean mMTC applications. Recently, reinforcement learning (RL) has shown to be a promising paradigm for solving the SFs allocation problem in LoRaWAN. For instance, in [8], [9], authors used a single-agent RL (SARL) approach to derive the optimal SFs allocation by considering the co-SF interference for improving the network reliability and throughput. To further enhance the exploration efficiency of SARL in mMTC applications, a multi-agent RL (MARL) approach has been applied in [10] to determine the optimal SFs allocation for improving the energy efficiency of underground EDs. However, the above RL approaches only adopt the basic deep Q-network (DQN). Considering the issues related to the overestimation and imprecision of Q value in the basic DQN, the multi-agent dueling double DQN (MAD3QN) is proposed to augment the agents’ optimization capabilities [11]. Meanwhile, in [12], another mainstream MARL algorithm based on value-based and policy-based optimization, namely the multi-agent advantage actor-critic (MAA2C), is developed to provide a more efficient exploration strategy compared to DQN.

To the best of our knowledge, there have been no studies exploring the effectiveness of MAD3QN or MAA2C in optimizing SFs allocation in LoRaWAN, let alone our considered massive subterranean NTN scenarios. Notably, energy efficiency is a significant metric from both economic and sustainable perspectives for the design of such a system [13], [14]. Motivated by this, this letter utilizes the MAD3QN and the MAA2C algorithms with our developed analytical reward mechanism for SFs allocation, aiming to maximize the system’s energy efficiency characterized by the average amount of energy consumed for uplink packet delivery. The simulation results demonstrate the superiority of our approach over four well-established benchmarks in extreme U-DtS scenarios.

II. SYSTEM MODEL

For the sake of clarity, in the rest of the paper, we focus on a U-DtS example based on a very-low-Earth-orbit (VLEO) satellite. However, the presented methods and obtained results can be generalized for any underground-to-NTN connectivity scenario, including those employing unmanned aerial vehicles or high-attitude platforms. Consider the subterranean

K. Lin and T. Hao are with the College of Surveying and Geo-Informatics, Tongji University, Shanghai, China. E-mail: lkq1220@tongji.edu.cn; tonghao@tongji.edu.cn. (*Corresponding Author: Tong Hao*)

M. Asad Ullah, H. Alves and K. Mikhaylov are with the Centre for Wireless Communications, University of Oulu, Finland. M. Asad Ullah is also with VTT Technical Research Centre of Finland Ltd, Oulu, Finland. Email: Muhammad.AsadUllah@oulu.fi; Hirley.Alves@oulu.fi; konstantin.mikhaylov@oulu.fi.

This work was supported in part by National Natural Science Foundation of China (No. 42211530077 and 42074179), the Academy of Finland, 6G Flagship program (No. 346208), and the China Scholarship Council.

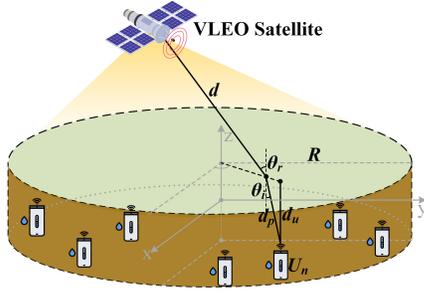


Fig. 1. The subterranean LoRaWAN NTN system taking a VLEO satellite as an example.

LoRaWAN NTN system depicted in Fig. 1 where a LoRaWAN gateway (GW) deployed on the non-terrestrial (NT) platform generates a single spot beam for covering a massive set $\mathcal{U} = \{U_n | n = 1, 2, \dots, N\}$ of underground EDs buried at the same depth d_u . Specifically, N underground EDs are distributed according to the Poisson point process (PPP) with intensity $\lambda = \frac{N}{\pi R^2 T_p}$ within a circular region of radius R . Herein, we assume that all EDs transmit an uplink packet with a period of T_p , the fixed physical layer (PHY) payload PL , the same bandwidth B , and the maximum transmit power P_t . A LoRa packet is successfully demodulated when the received signal-to-noise ratio (SNR) and signal-to-interference ratio (SIR) requirements are satisfied.

1) *Path Loss Model*: The total path loss from U_n to the NT GW comprises the air attenuation L_{air} , the fraction loss on the soil-air interface L_r , and the attenuation in underground soil L_{soil} . Consider that a ground-to-satellite link propagates in free space and the attenuation caused by the ionosphere, atmospheric gases, fog, clouds, and rain droplets can be neglected in the sub-GHz frequency band, the free space path loss model is adopted for the air path loss calculation [1]. Hence, the total path loss model is given by [1], [15].

$$g(d) = L_{air} L_r L_{soil} = \left(\frac{4\pi f_c}{c} \right)^2 (d)^\eta \left(\frac{2\beta d_p}{\exp(-\alpha d_p)} \right)^2, \quad (1)$$

where f_c is the carrier frequency, c denotes the speed of light in free space, d is the distance between U_n and the NT GW, η is the path-loss exponent, $d_p = d_u / (\cos \arcsin(1/\sqrt{\epsilon'}))$ is the length of the underground path; α and β are the attenuation constant and phase shifting constant, respectively,

$$\alpha = 2\pi f_c \sqrt{\frac{\mu_r \mu_0 \epsilon' \epsilon_0}{2} [\sqrt{1 + (\epsilon''/\epsilon')^2} - 1]}, \quad (2)$$

$$\beta = 2\pi f_c \sqrt{\frac{\mu_r \mu_0 \epsilon' \epsilon_0}{2} [\sqrt{1 + (\epsilon''/\epsilon')^2} + 1]}, \quad (3)$$

where μ_r is the soil's relative permeability, μ_0 is the free-space permeability, ϵ_0 is the free space permittivity. At the same time, ϵ' and ϵ'' are the real and imaginary parts of the soil's complex dielectric constant (CDC), i.e., $\epsilon = \epsilon' + j\epsilon''$. CDC can be calculated by the mineralogy-based soil dielectric model developed in [16]. Notably, the refraction loss on the soil-air interface L_r can be neglected in our study, implying $L_r = 1$. This is because most energy is refracted when electromagnetic waves propagate from a high-density medium (soil) to a lower-density one (air).

2) *Success Probability for SNR Guarantee*: In the absence of interference, the probability of successfully decoding a packet as a function of distance d is

$$P_{SNR}(d) = \mathbb{P} \left[\frac{P_t G_t G_r |h|^2}{g(d) \sigma_w^2} > q \right] = \exp \left(-\frac{g(d) q \sigma_w^2}{P_t G_t G_r} \right), \quad (4)$$

where G_t and G_r are the antenna gains of the underground ED and the NT GW, respectively, $|h|^2$ accounts for fading in the EDs-to-GW channel, whose coefficients are characterized by Rayleigh fading and the power follows an exponential distribution with a unit mean, σ_w^2 is the variance of the additive white Gaussian noise, and LoRa SF-specific SNR demodulation threshold $q = \{-6, -9, -12, -15, -17.5, -20\}$ dB for SF7~12 denoted by $\{SF_k | k = 1, \dots, 6\}$, respectively [2].

3) *Success Probability for SIR Guarantee*: The recent studies have demonstrated the presence of the capture effect for LoRa signals, which implies that a receiver demodulates the stronger packet under the interference of the weaker ones if the SIR is above a certain threshold [4]. Given the SIR threshold δ and the interference set (i.e., simultaneously transmitted packets featuring the same-SF) Φ , the SIR success probability according to distance d is

$$P_{SIR}(d) = \mathbb{P} \left[\frac{|h|^2 g(d)}{\sum_{i \in \Phi} |h_i|^2 g(d_i)} > \delta \right] \\ \stackrel{(a)}{=} \exp \left(\frac{-4N_k \cdot T_o A_k}{d_{max}^2 T_p N_c} \int_0^{d_{max}} \frac{\delta d_i^\eta d_i^{-\eta}}{1 + \delta d_i^\eta d_i^{-\eta}} d_i \, dd_i \right) \\ \stackrel{(b)}{=} \exp \left[\frac{-2N_k \cdot T_o A_k}{T_p N_c} {}_2F_1 \left(1, \frac{2}{\eta}; 1 + \frac{2}{\eta}; -\frac{d_{max}^\eta}{\delta d^\eta} \right) \right], \quad (5)$$

where (a) follows after using the probability generating functional of the product over PPP [4], and (b) is obtained by adopting the definition of the Gauss Hypergeometric function ${}_2F_1(\cdot)$ [17]. Furthermore, i represents the interfering signal, N_k is the number of EDs assigned by SF_k , d_{max} denotes the maximum distance between EDs and the NT GW, $T_o A_k$ is the time-on-air of SF_k , and N_c is the number of uplink channels.

4) *Packet Delivery Ratio*: The overall probability of successful packet delivery is the product of P_{SNR} and P_{SIR}

$$P_S(d) = P_{SNR}(d) P_{SIR}(d). \quad (6)$$

Fig. 2 highlights that the success probability P_S of the analytical model described in (6) agrees well with that obtained from the Monte-Carlo simulations, where $N = 1000$ (i.e., 1k) EDs transmit a 23-byte PHY payload packet in a single uplink channel with the period of $T_p = 600$ s.

5) *Energy Per Packet (EPP)*: The system's energy efficiency is characterized by EPP, which denotes the average amount of energy consumed by an ED to successfully deliver a packet to the NT GW [6]. Note that for the sake of tractability, we do not consider the downlink communication (i.e., LoRaWAN receive windows) and the energy consumption associated with it. Thus, EPP is given by

$$EPP = \frac{V_{supply} I_{tx} T_o A_k}{P_S}, \quad (7)$$

where $V_{supply} = 3.3$ V is the supply voltage of the ED, while I_{tx} is the transmit current consumption determined by P_t .

III. RL-BASED SFs ALLOCATION APPROACH

This work aims to determine the optimal SFs allocation strategy that minimizes the system's EPP in massive U-DtS scenarios. Hence, the optimization objective is formulated as

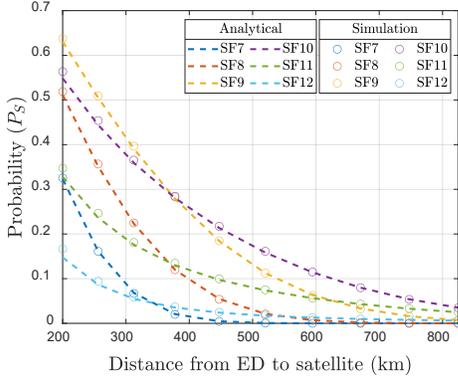


Fig. 2. Analytical and simulated overall success probability P_S versus distance from ED to satellite d . The related parameters can be found in Table I.

$$\min_{\{a_1, \dots, a_N\}} \frac{1}{N} \sum_{n=1}^N EPP_n \text{ s.t. } a_n \in \{SF_k\}, \quad (8)$$

where a_n is the selected SF configuration for the n -th ED. To achieve this, we utilize two MARL approaches, i.e., MAD3QN and MAA2C. Concretely, we consider each ED as an independent RL agent responsible for selecting and utilizing an SF configuration for packet transmission. The NT GW then applies the MARL approach to derive the SFs allocation strategy based on the received information from all agents. Finally, the NT GW broadcasts the SFs allocation results to all EDs. The MARL components are as follows:

1) *Agent*: Each agent includes $(s_n^t; a_n^t; r_n^t; s_n^{t+1})$, which implies the n -th ED in state s_n^t chooses an action a_n^t according to certain policy at step t . Then, it will receive its own r_n^t and fall into state s_n^{t+1} at the next training step.

2) *Action Space \mathcal{A}* : The action space of each agent is $A = \{SF_k\}$. Thus, the selected actions of all agents at step t can be defined as $\mathcal{A}^t = \{a_1^t, \dots, a_N^t\}$, $a \in A$.

3) *State Space \mathcal{S}* : The state observed in each agent consists of the selected SF configuration, SNR probability, SIR probability, packet delivery ratio, and EPP, which can be denoted as $s_n^t = \{a_n^t, (P_{SNR})_n^t, (P_{SIR})_n^t, (P_S)_n^t, EPP_n^t\}$. Hence, the state set of all agents at step t is $\mathcal{S}^t = \{s_1^t, \dots, s_N^t\}$.

4) *Reward \mathcal{R}* : The objective of the RL approach is to minimize the EPP; thus, the reward of each agent at step t is defined as $r_n^t = \frac{1}{EPP_n^t}$. Consequently, the reward set of all agents is $\mathcal{R}^t = \{r_1^t, \dots, r_N^t\}$. Notably, our reward mechanism, which accounts for local and global rewards, contributes to the expeditious convergence of the MARL algorithm.

A. MAD3QN Approach

The workflow of the MAD3QN approach is illustrated in **Algorithm 1**. The value-based MAD3QN is devoted to obtaining an optimal policy, which maps a state to a distribution over actions [11]. Each agent's policy is characterized by the DQN $Q(s_n^t, a; \omega_n)$ to estimate the expectation of action-value distribution, where ω_n denotes the policy network weights of the n -th agent. Hence, MAD3QN aims to search for optimal weights of each agent by minimizing the loss function, i.e.,

$$y_n^t = r_n^t + \gamma Q(s_n^{t+1}, \operatorname{argmax}_{a \in A} Q(s_n^{t+1}, a; \omega_n); \hat{\omega}_n), \quad (9)$$

$$\hat{L}_n = (y_n^t - Q(s_n^t, a_n^t; \omega_n))^2,$$

where $\gamma \in [0, 1)$ denotes a discount factor that balances the trade-off between immediate and future rewards, and $\hat{\omega}_n$

Algorithm 1 MAD3QN Approach

- 1: Initialize initial state \mathcal{S}^0 , policy network $Q(s_n^0, a; \omega_n)$ with random weights ω_n , target Q-network $Q(s_n^0, a; \hat{\omega}_n)$ with $\hat{\omega}_n = \omega_n$, replay memory \mathcal{M} , $\epsilon_t = 0$, $\epsilon_T = 0.9999$, $\gamma = 0.7$, and $m = 100$
- 2: **for** $t = 1$ to T_{max} **do**
- 3: **for** $n = 1$ to N **do**
- 4: $a_n^t = \begin{cases} \text{Randomly select } a_n^t \in A, \text{rand}() > \epsilon_t \\ \operatorname{argmax}_{a \in A} Q(s_n^t, a; \omega_n), \text{ otherwise} \end{cases}$
- 5: **end for**
- 6: Execute action \mathcal{A}^t in the environment and get $(\mathcal{S}^t, \mathcal{A}^t, \mathcal{S}^{t+1}, \mathcal{R}^t)$ by (4), (5), (6), and (7)
- 7: **for** $n = 1$ to N **do**
- 8: Store transition $(a_n^t, s_n^t, r_n^t, s_n^{t+1})$ in \mathcal{M}
- 9: **if** \mathcal{M} is full **then**
- 10: Sample random mini-batch of transitions from \mathcal{M}
- 11: Update ω_n by performing a gradient descent step on (9)
- 12: **end if**
- 13: Set state $s_n^t = s_n^{t+1}$
- 14: Update ϵ_t with $\epsilon_t = \min(\epsilon_T, \epsilon_t + 0.0002)$
- 15: Every m steps clone ω_n to $\hat{\omega}_n$
- 16: **end for**
- 17: **end for**

Output: Learned $Q(s_n^t, a; \omega_n)$

denotes the target network weights of the n -th agent. The target network is generated by cloning the current network and updating the weights after a fixed number of iterations. The network weight ω_n is updated through a gradient descent method, i.e., $\omega_n = \omega_n - \tau \nabla_{\omega_n} \hat{L}_n(\omega_n)$, where τ is the learning rate. Furthermore, compared with the single estimator in the basic DQN, the output of each agent in MAD3QN is divided into two estimators, i.e., value function and advantage function, to accelerate the convergence. Accordingly, the Q-value of each agent in MAD3QN can be presented as

$$Q(s_n^t, a_n^t; \bar{\omega}_n, \kappa_n, \nu_n) = \mathbb{V}(s_n^t; \bar{\omega}_n, \kappa_n) + \left[\mathbb{A}(s_n^t, a_n^t; \bar{\omega}_n, \nu_n) - \frac{1}{|A|} \sum_{a \in A} \mathbb{A}(s_n^t, a; \bar{\omega}_n, \nu_n) \right], \quad (10)$$

where $\bar{\omega}_n, \kappa_n, \nu_n$ are the weights of the shared convolutional encoder, value function $\mathbb{V}(\cdot)$ and advantage function $\mathbb{A}(\cdot)$, respectively, for the n -th agent.

B. MAA2C Approach

The workflow of the MAA2C approach is described in **Algorithm 2**. Unlike MAD3QN, MAA2C focuses on training the critic function $C(s_n^t; \hat{\psi}_n)$ that measures average expected return from current state s_n^t to obtain the optimal actor policy $T(a_n^t | s_n^t; \psi_n)$ of the n -th agent [12]. The ψ_n and $\hat{\psi}_n$ are the actor and critic network weights, respectively. MAA2C aims to obtain the optimal policy of each agent by minimizing the loss of actor and critic functions. The loss function of the critic network for the n -th agent is

$$z_n = r_n^t + \gamma C(s_n^{t+1}; \hat{\psi}_n) - C(s_n^t; \hat{\psi}_n), \quad (11)$$

$$\hat{L}_C = (z_n)^2.$$

Meanwhile, the loss function of the actor network is

$$\hat{L}_T = -z_n \log(T(a_n^t | s_n^t; \psi_n)). \quad (12)$$

Herein, the network weights of actor and critic are updated by a gradient descent method, i.e., $\psi_n = \psi_n - \tau \nabla_{\psi_n} \hat{L}_T(\psi_n)$ and $\hat{\psi}_n = \hat{\psi}_n - \tau \nabla_{\hat{\psi}_n} \hat{L}_C(\hat{\psi}_n)$, respectively.

Algorithm 2 MAA2C Approach

```

1: Initialize initial state  $S^0$ , actor networks  $T(a_n^0|s_n^0; \psi_n)$  with
   random weights  $\psi_n$ , critic network  $C(s_n^0; \hat{\psi}_n)$  with random
   weights  $\hat{\psi}_n$ , and  $\gamma = 0.7$ 
2: for  $t = 1$  to  $T_{max}$  do
3:   for  $n = 1$  to  $N$  do
4:     Select  $a_n^t \sim T(a_n^t|s_n^t; \psi_n)$ 
5:   end for
6:   Execute action  $\mathcal{A}^t$  in the environment and get
    $(S^t, \mathcal{A}^t, S^{t+1}, \mathcal{R}^t)$  by (4), (5), (6), and (7)
7:   for  $n = 1$  to  $N$  do
8:     Update  $\hat{\psi}_n$  by performing a gradient descent step on (11)
9:     Update  $\psi_n$  by performing a gradient descent step on (12)
10:    Set state  $s_n^t = s_n^{t+1}$ 
11:   end for
12: end for
Output: Learned  $T(a_n^t|s_n^t; \psi_n), C(s_n^t; \hat{\psi}_n)$ 

```

TABLE I
SIMULATION PARAMETERS

Parameters	Values
Operation Environments [1]	
Total nodes (N)	1k~10k
Burial depth (d_u)	0.4 m
VWC (m_v)	11.19%
Clay (m_c)	16.86%
PHY payload (PL)	23 Bytes
Report period (T_p)	600 s
Traffic pattern	Periodic
VLEO Satellite Configuration [18]	
Elevation angles (E)	$10^\circ \leq E \leq 90^\circ$
Orbital height (H)	200 km
Link distance (d)	200 km~825 km
Coverage radius (R)	822 km
Radio Configuration [1]	
Carrier frequency (f_c)	915 MHz
Antenna gains (G_t and G_r)	(2.15 dBi, 35 dBi)
SIR threshold (δ)	6 dB
Path loss exponent (η)	2
Uplink channel number (N_c)	1
Transmit power (P_t)	20 dBm
Transmit current (I_{tx})	133 mA
BW (B)	125 KHz
Agent Setting	
Approach	MAD3QN MAA2C
	Actor Critic
Learning rate (τ)	0.001 0.001 0.001
Input layer	linear, $ s_n $ linear, $ s_n $ linear, $ s_n $
Hidden layer	ReLU, 16 ReLU, 16 ReLU, 16
Output layer	linear, $ A $ softmax, $ A $ linear, 1

IV. SIMULATION RESULTS AND ANALYSIS

A. Simulation Settings

To evidence the performance of our proposed approach, the extreme subterranean LoRaWAN NTN scenario, i.e., U-DtS, is considered in the simulation. Herein, we envision a LoRaWAN GW deployed on a VLEO satellite, which provides one spot beam covering the real-life center-pivot irrigation farm [1], [18]. The specific simulation parameters are listed in Table I. To ensure the execution of the algorithm, we assume that our system is capable of compensating for Doppler effects and accurately predicting the positions of satellites [19]. For the MARL parameters, we set that the reply memory size is $|\mathcal{M}| = 16$, the mini-batch size is 4, and the maximum training episode is $T_{max} = 6000$. We optimize the weights of all DQNs with the Adam optimizer. Besides, the hyper-parameter

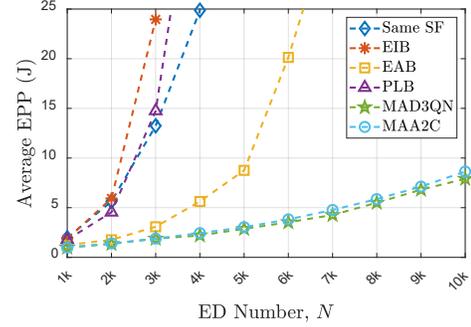


Fig. 3. The comparison of the average EPP versus the number of EDs for proposed RL-based SFs allocation scheme and the same SF [3], EIB [4], EAB [5], and PLB [7] based techniques.

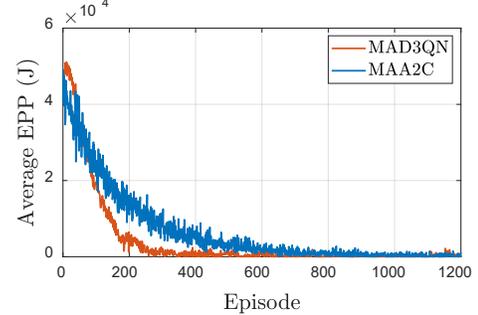


Fig. 4. The convergence performance of MAD3QN and MAA2C under $N = 5k$ EDs.

settings, the activation functions, and the neuron numbers of each deep neural network layer are summarized in Table I.

For performance comparison, four benchmarks are implemented: (1) Same SF [3]: all EDs are assigned by the same SF; (2) EIB [4]: the coverage area is partitioned into six equal-width (i.e., $\frac{R}{|A|}$) annuli for SFs allocation; (3) EAB [5]: the coverage area is segmented into six equal-area (i.e., $\frac{\pi R^2}{|A|}$) annuli for SFs allocation; (4) PLB [7]: annuli are determined based on the path loss model described in (1) and the SF-specific SNR threshold q .

B. Simulation Results

We first present the average EPP (i.e., $\frac{1}{N} \sum_1^N EPP_n$) under different approaches as a function of the number of EDs N in Fig. 3. One can observe that the average EPP increases with the number of EDs due to the high probability of co-channel co-SF interference in a denser network. Our proposed MAD3QN and MAA2C approaches yield a significant improvement in the average EPP compared to the other four benchmarks, and such a performance gain becomes more pronounced with an increase in N . Meanwhile, the average EPP of MAD3QN is slightly lower than that of MAA2C, implying better energy efficiency. Specifically, the average EPP of MAD3QN and MAA2C is 2.46 J and 2.61 J, respectively, at $N = 5k$.

Fig. 4 depicts the average EPP of MAD3QN and MAA2C under $N = 5k$ EDs concerning the training episodes, from which we can observe that MAD3QN performs better than MAA2C in terms of convergence rate. For instance, the EPP of MAD3QN converges to a stable value after around 300th training episodes, while it takes MAA2C nearly 600th training episodes to converge.

Fig. 5 demonstrates the SFs distribution for $N = 5k$ EDs under different approaches. In Fig. 5(a), the same SF

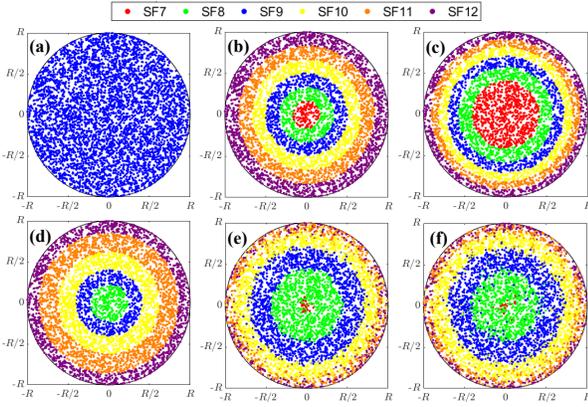


Fig. 5. SF distribution of $N = 5k$ EDs under (a) Same SF, (b) EIB, (c) EAB, (d) PLB, (e) MAD3QN, and (f) MAA2C approaches. The color of point represents the selected SF configuration, e.g. red for SF7.

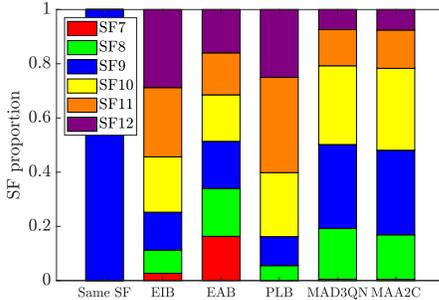


Fig. 6. The statistical proportions of each SF for Figs. 5(a)~(f), respectively.

scheme allocates all EDs with SF9, which balances ToA and propagation capability. However, many EDs configured with the same SF can result in a lower P_{SIR} . Meanwhile, P_{SNR} of the EDs at the edge degrades due to the limited link budget. Figs. 5(b) and (d) highlight that EIB and PLB allocate the larger SF (i.e., SF11 and SF12) to the peripheral EDs for a higher P_{SNR} . However, the significant number of EDs operating at higher SF levels leads to a considerable increase in EPP. This is primarily due to the elevated transmit power consumption and the higher collision probability caused by the extended ToA. Despite EAB equally assigning SFs among all EDs to mitigate the same-SF interference, as depicted in Fig 5(c), the EDs configured with SF7 cannot establish reliable U-DtS connectivity, which results in the worse P_{SNR} and the increased average EPP. Figs. 5(e) and (f) reveal that our proposed MARL approaches assign the EDs near the inner ring by SF8~10 while allocating SF11 or SF12 to the EDs at the edge. This allocation strategy aims to accomplish the robust link with the lower transmit power consumption and to reduce the co-SF interference probability, thereby improving the average EPP. Consequently, our proposed approaches exhibit notably superior performance compared to all the four benchmarks. The share of EDs using each SF under the discussed approaches are illustrated in Fig. 6.

V. CONCLUSION

This letter investigates the effectiveness of MARL for optimizing SFs allocation in massive U-DtS scenarios. After developing an analytical model to characterize packet delivery ratio and using it as our reward mechanism, we utilize the MAD3QN and MAA2C approaches to optimize SFs allocation for improving the system's energy efficiency. Through a com-

parison with the four benchmarks in a realistic farm case, our numerical results reveal that the proposed approaches exhibit the lowest average EPP, where MAD3QN slightly outperforms MAA2C in terms of the average EPP and convergence rate. Note that our proposed MARL approach is universal, and can be generalized for other subterranean LoRaWAN NTN applications. The future work will focus on developing a strategy for broadcasting the derived SF configuration to each ED by considering in more detail the mobility of NTN and the reliable downlink communication [20], [21].

REFERENCES

- [1] K. Lin *et al.*, "Subterranean mMTC in Remote Areas: Underground-to-Satellite Connectivity Approach," *IEEE Commun. Mag.*, vol. 61, no. 5, pp. 136–142, 2023.
- [2] O. Georgiou and U. Raza, "Low Power Wide Area Network Analysis: Can LoRa Scale?" *IEEE Wireless Commun. Lett.*, vol. 6, no. 2, pp. 162–165, 2017.
- [3] M. A. Ullah *et al.*, "Enabling mMTC in Remote Areas: LoRaWAN and LEO Satellite Integration for Offshore Wind Farm Monitoring," *IEEE Trans. Ind. Informal.*, vol. 18, no. 6, pp. 3744–3753, 2022.
- [4] A. Mahmood *et al.*, "Scalability Analysis of a LoRa Network Under Imperfect Orthogonality," *IEEE Trans. Ind. Informal.*, vol. 15, no. 3, pp. 1425–1436, 2019.
- [5] J.-T. Lim and Y. Han, "Spreading Factor Allocation for Massive Connectivity in LoRa Systems," *IEEE Commun. Lett.*, vol. 22, no. 4, pp. 800–803, 2018.
- [6] J. Park *et al.*, "EARN: Enhanced ADR With Coding Rate Adaptation in LoRaWAN," *IEEE IoT J.*, vol. 7, no. 12, pp. 11 873–11 883, 2020.
- [7] K. Lin and T. Hao, "Adaptive Selection of Transmission Configuration for LoRa-based Wireless Underground Sensor Networks," in *Pro. 2021 IEEE Wireless Commun. Netw. Conf.*, 2021, pp. 1–6.
- [8] R. M. Sandoval *et al.*, "Optimizing and Updating LoRa Communication Parameters: A Machine Learning Approach," *IEEE Trans. Netw. Serv. Manag.*, vol. 16, no. 3, pp. 884–895, 2019.
- [9] I. Ilahi *et al.*, "LoRaDRL: Deep Reinforcement Learning Based Adaptive PHY Layer Transmission Parameters Selection for LoRaWAN," in *Proc. IEEE 45th Conf. Local Comput. Netw. (LCN)*, 2020, pp. 457–460.
- [10] G. Zhao *et al.*, "Optimizing energy efficiency of LoRaWAN-based wireless underground sensor networks: A multi-agent reinforcement learning approach," *Internet of Things*, vol. 22, p. 100776, 2023.
- [11] H. Xiang *et al.*, "Multi-Agent Deep Reinforcement Learning-Based Power Control and Resource Allocation for D2D Communications," *IEEE Wireless Commun. Lett.*, vol. 11, no. 8, pp. 1659–1663, 2022.
- [12] Y. Shao *et al.*, "Graph Attention Network-Based Multi-Agent Reinforcement Learning for Slicing Resource Management in Dense Cellular Network," *IEEE Trans. Veh. Technol.*, vol. 70, no. 10, pp. 10 792–10 803, 2021.
- [13] Z. Lin *et al.*, "Secrecy-Energy Efficient Hybrid Beamforming for Satellite-Terrestrial Integrated Networks," *IEEE Trans. Commun.*, vol. 69, no. 9, pp. 6345–6360, 2021.
- [14] Z. Lin *et al.*, "SLNR-Based Secure Energy Efficient Beamforming in Multibeam Satellite Systems," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 59, no. 2, pp. 2085–2088, 2023.
- [15] K. Lin and T. Hao, "Experimental Link Quality Analysis for LoRa-Based Wireless Underground Sensor Networks," *IEEE IoT J.*, vol. 8, no. 8, pp. 6565–6577, 2021.
- [16] V. L. Mironov *et al.*, "Physically and Mineralogically Based Spectroscopic Dielectric Model for Moist Soils," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 7, pp. 2059–2070, 2009.
- [17] J. M. de Souza Sant'Ana *et al.*, "LoRa Performance Analysis with Superposed Signal Decoding," *IEEE Wireless Commun. Lett.*, vol. 9, no. 11, pp. 1865–1868, 2020.
- [18] J. Walsh *et al.*, "Drag reduction through shape optimisation for satellites in Very Low Earth Orbit," *Acta Astronautica*, vol. 179, pp. 105–121, 2021.
- [19] T. Janssen *et al.*, "A Survey on IoT Positioning Leveraging LPWAN, GNSS, and LEO-PNT," *IEEE IoT J.*, vol. 10, no. 13, pp. 11 135–11 159, 2023.
- [20] Z. Xu *et al.*, "Enhancement of Satellite-to-Phone Link Budget by Using Distributed Beamforming," 2023. [Online]. Available: arXiv:2308.04818
- [21] Z. Xu *et al.*, "Enhancement of Direct LEO Satellite-to-Smartphone Communications by Distributed Beamforming," 2023. [Online]. Available: arXiv:2308.05055