

I. ENVIRONMENT

A. General parameters

- 100 by 100 meters squared area.
- 1050 Users.
- Each User is assumed to hold their device at a fixed height $h_u = 1.5$ meter.

B. User distribution

- The users are distributed under three 2-dimensional Gaussian distillation layers and one uniform distribution layer
- The center of the cluster is determined at random.
- It can take positions from 30 to 70, both in x and y coordinates.
- The users are scattered randomly with normal distribution around the cluster center, with 300, 250 and 200 MSs follows a 10, 7 and 6 meters standard deviation respectively.
- 300 MSs follow a uniform distribution

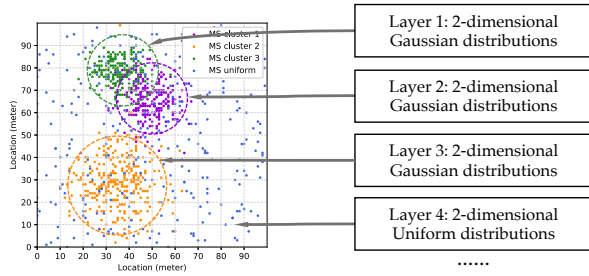


Fig. 1. Environment construction

C. Drones

- The drones are assumed to have ideal communication among themselves.
- The drones share their positions and number of users allocated with each other.
- Each drone has a directional antenna with a main lobe with an aperture angle of $\theta = 60$ degrees. The antenna is pointing downwards. An illustration of the directivity angle is presented in Fig. 2
- There is no limit on the number of users that each drone can allocate.
- Each drone is assumed to be flying at a fixed height $h_d = 30$ meters.

II. STATES

- The states are the positions of both drones.
- The environment is discretized into 121 possible positions for each drone (steps of 10 meters).
- Drones cannot assume the same position at the same time.
- The total number of states is therefore $2 \times \binom{121}{2} = 14520$.

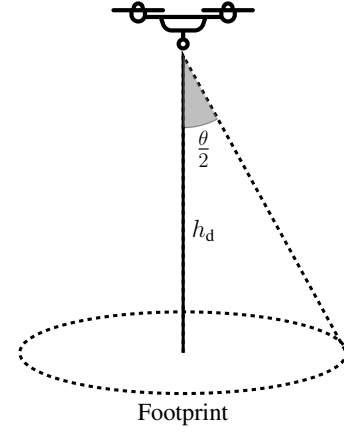


Fig. 2. Illustration of the coverage radius of a drone flying at h_d meters and with directivity angle of θ .

III. ACTIONS

- The possible actions are to move ± 1 or ± 0 (stay) step in x or y .
- If one drone would move to the same position as the other (e.g. chooses the action to move right when the other drone is one step to its right), it does not move.
- If a drone would move out of the grid (e.g. chooses to move right when at x coordinate 100), it does not move.
- If a state has not been explored, the action is also chosen at random, in order to avoid bias.
- Otherwise, actions are deterministic.

IV. POLICY

- The policy is ε -greedy.
- Meaning that each drone chooses a random action with probability ε and $\max(Q_{s_{t+1},*})$ with probability $1 - \varepsilon$.

V. REWARD

- The reward is the total number of users allocated by both drones.

A. User Association

- A user associates with a drone if its received SINR is above a threshold of 40 dB.
- The SINR for the link between user u and drone n is calculated according to

$$SINR_{n,u} = \frac{RSRP_{n,u}}{N + \sum_{i \neq n} RSRP_{i,u}}. \quad (1)$$

- The RSRP for the link between user u and drone n is calculated according to the free space path loss, as

$$RSRP_{n,u} = \frac{P_t}{16(\pi f_c d)^2}, \quad (2)$$

where c is the speed of light in meters per second, P_t is the drone transmit power in Watts and f_c is the carrier frequency in Hz.

- Any user outside the main lobe is considered to receive 0 W of power from the drone.

VI. TRAINING

- The training session is comprised of 100 episodes of 20,000 iterations each.
- The state of the drones is set randomly at the beginning of each episode.
- ε decays exponentially with the episode number, according to:

$$\varepsilon_j = e^{-j/20}, \quad (3)$$

where j is the episode number and e is Euler's constant.

A. Algorithm

The update strategy for DQN is expressed as,

$$L(\theta) = E[(\text{TargetQ} - Q(s, a; \theta))^2] \quad (4)$$

where TargetQ is present as,

$$\text{TargetQ} = r + \gamma \max_{a'} Q(s', a'; \theta) \quad (5)$$

where α is the learning rate and γ is the discount factor. A pseudo code of the implemented solution is presented as follow.

Algorithm 1: DQN implementation

```

1: Initialization
2: for every episode  $j$  do
3:    $s_1 \leftarrow$  random
4:   for Every iteration  $t$  do
5:     for Every drone  $\delta$  do
6:        $a_t \leftarrow \max_a Q(\phi(s_t), a; \theta)$  with probability  $\epsilon$ 
7:       select a random action
8:        $r_t, x_{t+1} \leftarrow$  based on the action  $a_t$ 
9:        $\phi_{t+1} = \phi(s_{t+1}) \leftarrow s_{t+1} = (s_t, a_t, x_{t+1})$ 
10:       $D \leftarrow$  Add to dataset  $D + (\phi_t, a_t, r_t, \phi_{t+1})$ 
11:       $evaQ_{t+1} \leftarrow$  evaQ learn from data in  $D$ 
12:       $s_t \leftarrow s_{t+1}$ 
13:    end for
14:     $Q \leftarrow$  fine-tune from  $evaQ$  for every  $n$  steps
15:  end for
16: end for

```

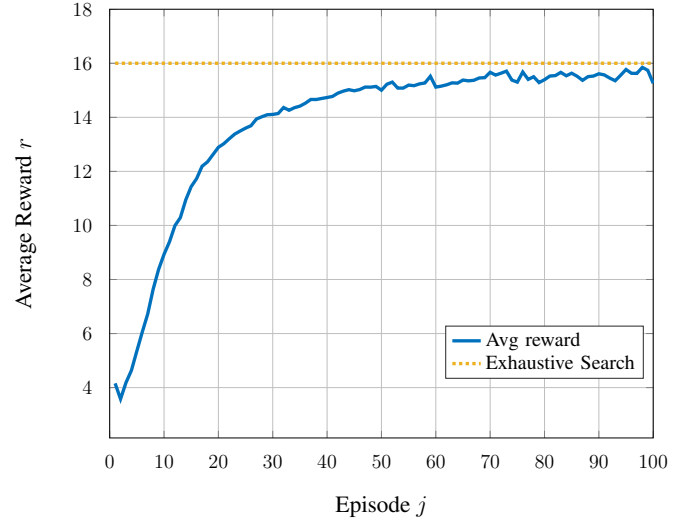


Fig. 3. Average reward per episode considering for this training session

VII. RESULTS

The average reward per episode is shown in Fig. 3.