

Appendix:

I. ENVIRONMENT

A. General parameters

- 100 by 100 meters squared area.
- 1050 Users with three clusters
- The number of ABSes is 2, 4, 8.

B. User distribution

- The users are distributed around a local cluster.
- 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 mean $\mu = 0$ and standard deviation $\sigma = 20$.
- An example of the user distribution can be seen in Fig. 1. Total three clusters are simulated.

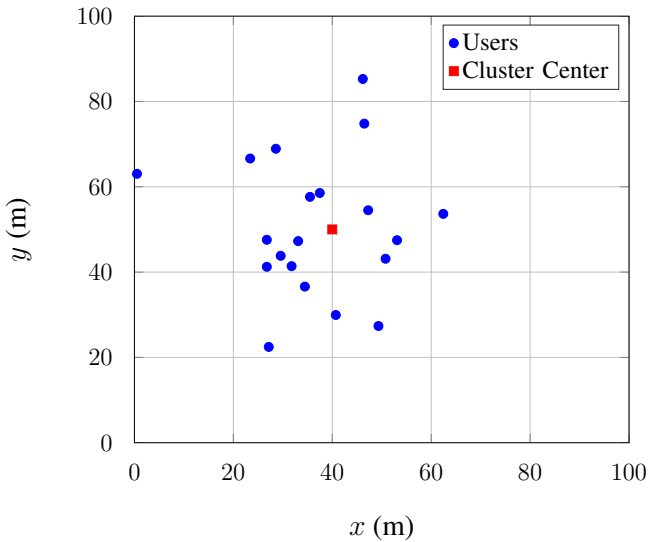


Fig. 1. Example of user distribution for one cluster in one training session.

C. Airbone Base Stations (ABSes)

- The ABSes are assumed to have ideal communication among themselves.
- The ABSes share their positions and number of users allocated with each other.
- Each ABS 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 ABSes can allocate.
- Each ABS is assumed to be flying at a fixed height $h_d = 30$ meters.

II. STATES (for two ABSes Cases)

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

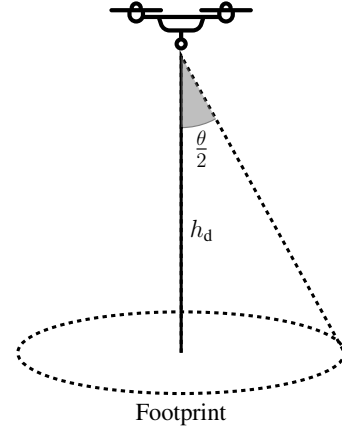


Fig. 2. Illustration of the coverage radius of a ABSes flying at h_d meters and with directivity angle of θ .

III. ACTIONS

- The possible actions are to move ± 1 step in x or y .
- If one ABS would move to the same position as the other (e.g. chooses the action to move right when the other ABS is one step to its right), it does not move.
- If a ABS 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 ABS chooses a random action with probability ϵ and maxQ with probability $1 - \epsilon$.

V. REWARD

- The reward is the total number of users allocated by all ABSes.

A. User Association

- A user associates with a ABSes if its received SINR is above a threshold of 40 dB.
- The SINR for the link between user u and ABS 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 ABS n is calculated according to the free space path loss, as

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

where c is the speed of light in meters per second, P_t is the ABSes 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 ABS.

VI. TRAINING

- The typical training session is comprised of 10 episodes of 2000 iterations each. At DQN we use 1000 episodes •
- The state of the ABSes 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 (SARSA)

The update strategy for SARSA is expressed as

$$Q(s_t, a_t, \delta) = Q(s_t, a_t, \delta) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}, \delta) - Q(s_t, a_t, \delta)), \quad (4)$$

where α is the learning rate and γ is the discount factor. A pseudo code of the implemented solution is presented in Algorithm 1.

Algorithm 1: SARSA implementation

```

1 Initializations
2 for Every episode  $j$  do
3    $s_1 \leftarrow \text{random}(1, 14520)$ 
4   for Every iteration  $t$  do
5     for Every ABSes  $\delta$  do
6        $a_t \leftarrow \text{chooseAction}(Q_{s_t, *}, \varepsilon_j, \delta)$ 
7        $s_{t+1} \leftarrow \text{takeAction}(s_t, a_t, \delta)$ 
8        $a_{t+1} \leftarrow \text{chooseAction}(Q_{s_{t+1}, *}, \varepsilon_j, \delta)$ 
9        $r_t \leftarrow \text{computeReward}(s_{t+1})$ 
10       $Q(s_t, a_t, d) \leftarrow \text{updateQ}()$ 
11       $s_t \leftarrow s_{t+1}$ 
12    end
13  end
14 end
```

VII. RESULTS

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

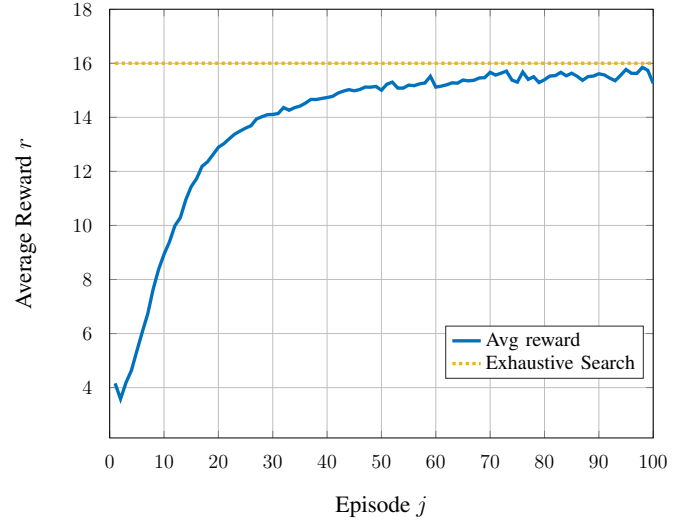


Fig. 3. Average reward per episode considering for this training session