

# RL-Based Interference Mitigation in Uncoordinated Networks with Partially Overlapping Tones

Mrugen Deshmukh\*, Md Moin Uddin Chowdhury\*, Sung Joon Maeng\*, Alphan Şahin†, İsmail Güvenç\*

\*Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC, USA

†Electrical Engineering Department, University of South Carolina, Columbia, SC, USA

Email: {madeshmu,mchowdh,smaeng,iguvenc}@ncsu.edu, asahin@mailbox.sc.edu

**Abstract**—Partially-overlapping tones (POT) are known to help mitigate co-channel interference in uncoordinated multi-carrier networks by introducing intentional frequency offsets (FOs) to the transmitted signals. In this paper, we explore the use of POT with reinforcement learning (RL) in dense networks where multiple links access time-frequency resources simultaneously. We propose a novel framework based on Q-learning, to obtain the FO for the multi-carrier waveform used for each link. In particular, we consider filtered multi-tone (FMT) systems that utilize Gaussian, root-raised-cosine (RRC), and isotropic orthogonal transform algorithm (IOTA) based prototype filters. Our simulation results show that the proposed scheme enhances the capacity of the links by at least 30% in additive white Gaussian noise (AWGN) channel at high signal-to-noise ratio (SNR), and even more so in the presence of severe multi-path fading. For a wide range of interfering link densities, we demonstrate substantial improvements in the outage probability and multi-user efficiency facilitated by POT, with the Gaussian filter outperforming the other two filters.

**Index Terms**—Device-to-device (D2D), uncoordinated networks, Internet-of-things (IoT), partial overlapping, Q-learning

## I. INTRODUCTION

Internet-of-things (IoT) is a rapidly growing technology that is expected to connect billions of devices in the near future. IoT devices can be densely located in area on the order of  $10^3$  per  $\text{km}^2$  [1]. Hence, they can form very dense uncoordinated networks and face spectral efficiency and throughput challenges due to complicated multi-user interference scenarios. Partially-overlapping tones (POT) have recently gained attention for mitigating the co-channel interference (CCI) in uncoordinated networks. By introducing an intentional frequency offset (FO) for each link equal to a fraction of the frequency spacing between two subcarriers, POT facilitates a reduction of multi-user interference among neighboring links.

The related research on POT in the literature is largely based on using partially overlapping channels (using all the available spectrum) to improve throughput in wireless networks. For instance, partial overlapping channels are used in [2] to improve the throughput in a remote wireless device-to-device (D2D) network using unmanned aerial vehicles (UAVs), in [3] broadly considering wireless mesh networks (WMNs) scenarios, and in [4] for a WMN scenario specific to IoT devices. However, these studies do not consider the effect of the waveform

used in the physical layer (PHY). Individual subcarriers and several schemes for waveform are considered in [5], which lays a theoretical foundation for the analysis of using different waveform types in POT. In [6], the authors explore POT for the cellular networks and propose an algorithm, called *Play n Wait*, to assign FOs sequentially where one user scans for the best possible FO while all other users have to wait for several seconds. This assumption may not hold for a practical uncoordinated network. In this study, we focus FO assignment problem for POT as in [6] and aim at addressing this challenge with reinforcement learning (RL).

In wireless networks, RL is commonly used to solve resource allocation [7], [8] as well as power allocation problems [9]. For example, in [7], energy efficiency maximization problem of a hybrid-powered dense network is studied, considering an actor-critic RL technique. In [8], the authors introduce RL-based decentralized resource allocation techniques while taking strict delay constraints into account inherent in vehicle-to-vehicle networks. A distributed RL algorithm is utilized for maintaining fairness and quality of service in dense heterogeneous networks in [9]–[11]. In [12], authors propose a deep Q-learning based algorithm that solves an optimization problem considering power control, beamforming, and interference coordination for sub-6 GHz and above-6 GHz bands.

In this paper, we develop an *offline* Q-learning algorithm to assign intentional FO to each link. Since the algorithm is trained beforehand and available to all the users within the network, there is no time delay in assigning these FOs. For POT, we consider filtered multi-tone (FMT), which is a subset of filter-bank multicarrier, that allows for every subcarrier being filtered individually while allowing complex modulation symbols [13], forming a grid-like structure in the time-frequency plane that partially overlapping subcarriers can exploit. We utilize Gaussian, root-raised-cosine (RRC), and isotropic orthogonal transform algorithm (IOTA) based prototype filters with the FMT POT framework, and study the associated trade-offs for using them with POT. With numerical results, we demonstrate the benefits of using POT with FMT in terms of improvements in capacity, multi-user efficiency (ME), and outage probabilities considering common propagation channels.

The rest of the paper is organized as follows. Section II discusses the system model. Section III describes our reinforcement learning based algorithm, while Section IV describes the

This research is supported by the National Science Foundation (NSF) CNS through the award number 1814727.

formulation of the ME in POT. Section V demonstrates the simulation results related to the proposed learning scheme. We finalize our paper in Section VI with some concluding remarks.

## II. SYSTEM MODEL

Consider  $U$  links uniformly distributed in an area where each link consists of a transmission point (TP) and a reception point (RP). Among  $U$  links, consider any one to be a link of interest with index  $u$  for  $u \in \{1, \dots, U\}$ . The rest of the TPs are then considered to be *aggressors* and the RP of the link of interest is the *victim*. We model the transmitted signal from the TP of the  $u$ th link as

$$s_u(t) = \sum_{l=-\infty}^{\infty} \sum_{n=0}^{N-1} X_{ln}^u g_{ln}(t), \quad (1)$$

where  $X_{ln}^u$  is the information symbol to be transmitted from the  $u$ th TP,  $l$  is the time index,  $n$  is the subcarrier index,  $N$  is the total number of subcarriers, and  $g_{ln}(t)$  is the synthesis function [14] that maps  $X$  to the time-frequency domain in a lattice structure as  $g_{ln}(t) = p(t - l\tau_0)e^{j2\pi n\nu_0 t}$ , where  $p(t)$  is the prototype filter being used,  $\tau_0$  is the time spacing between two consecutive symbols and  $\nu_0$  is the spacing between any two subcarriers.

The received signal at the  $u$ th RP can be calculated as

$$y_u(t) = \sum_{i=1}^U \int_{\tau_{i,u}} h_{i,u}(\tau_{i,u}, t) s_i(t - \tau_{i,u}) dt + w(t), \quad (2)$$

where  $h_{i,u}(t)$  is the channel impulse response between the desired RP of  $u$ th pair and TP of  $i$ th pair (one of which is the desired signal),  $s_i(t)$  is the transmitted signal from the TP of the  $i$ th pair, and  $w(t)$  is the additive white Gaussian noise (AWGN). The information symbol  $\tilde{X}_{mk}^u$  can be obtained by calculating the projection of  $y_u(t)$  onto the analysis function  $\gamma_{mk}(t) = p(t - m\tau_0)e^{j2\pi k\nu_0 t}$  as [14]

$$\tilde{X}_{mk}^u = \langle y_u(t), \gamma_{mk}(t) \rangle, \quad (3)$$

where  $m$  and  $k$  are the time and subcarrier indices at the receiver, respectively. By grouping terms related to the interference, (3) can be written as

$$\begin{aligned} \tilde{X}_{mk}^u &= \underbrace{G_u X_{mk}^u A_{mkmk}^u}_{\text{desired symbol}} + \underbrace{G_u \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} X_{mn}^u A_{lnmk}^u}_{\text{self-interference}} + \\ &\underbrace{\sum_{i \neq u} G_i \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} X_{ln}^i A_{lnmk}^i}_{\text{co-channel interference}} + \underbrace{W_u}_{\text{AWGN}}, \end{aligned} \quad (4)$$

where  $G_u$  and  $G_i$  are the channel gains at the desired signal and the  $i$ th aggressor, respectively,  $X_{mk}^i$  is the symbol of the  $i$ th aggressor,  $A_{mkmk}^u$  and  $A_{lnmk}^i$  represent the coefficients obtained through the corresponding ambiguity functions of the

desired signal and  $i$ th aggressor, respectively, which can be calculated as [5]:

$$A_{nlmk}^i = \int_{\tau} \int_{\nu} \int_t g_{ln}(t - \tau) e^{j2\pi \Delta f_i (t - \tau)} \gamma_{mk}(t) e^{j2\pi \nu t} dt d\tau d\nu, \quad (5)$$

where  $\Delta f_i$  is the intentional FO given to the  $i$ th aggressor.

For POT, the amount of CCI and the amount of self-interference (SI) in (4) are adjusted through  $\Delta f_i$ . By sacrificing the orthogonality of the pulses on the desired link through more time-dispersive filters, an intentional FO prevents aggressors' transmit pulses from fully overlapping with the receiver filters. Hence, it can reduce CCI even if the pulses are not aligned in time, which results in higher throughput for TP-RP links. For a large network, the assignment of the FOs for each link should be chosen such that the capacity of the entire network should be optimized. Also, the parameters of the filters need to be optimized as they also affect the SI, as can be seen from (5).

### A. Prototype Filters

The filters used in (1) and (3) determine the inter-symbol interference characteristics. In this study, we consider three filters for analysis: 1) RRC filter with a roll-off factor  $\alpha$ , 2) Gaussian filter with a time-frequency dispersion parameter  $\rho$ , and 3) IOTA filter with a dispersion parameter  $\rho$ . While RRC filters give a set of orthogonal functions, which is a commonly used for single-carrier schemes, Gaussian pulses yield a non-orthogonal base where the basis functions are localized in both time and frequency optimally. IOTA filter is derived from the Gaussian filter. A Gaussian filter is modified such that when it is orthogonal to its shifted versions in time and frequency. Therefore, its time-frequency characteristics are similar to Gaussian filter [14].

## III. COMBINING Q-LEARNING WITH POT

In this study, we consider Q-learning to address the FO assignment problem. Q-learning is a model-free RL algorithm whose learned decision policy is determined by state-action value function  $Q$  [15], which estimates long-term discounted rewards for each state-action pair. We assume that the pairs in the network do not communicate with each other since we consider an uncoordinated network. We also assume that once an aggressor TP-RP pair enters the network, they stay active. Let  $S_u$  be the number of aggressors adjacent to the  $u$ th link. The parameter  $S_u$  is estimated through a simple counting method, i.e.,  $S_u$  increases by 1 if the signal-to-interference-plus-noise ratio (SINR) of the  $u$ th link drops by more than 3 dB or decreases by 1 if the SINR increases by 3 dB. Each pair actively monitors  $S_u$ .

For the sake of illustrating the proposed scheme, consider a network with  $U = 3$  links. When the first link is established, there is no aggressor. After the first aggressor enters the network, it will fully overlap with the victim at the first link. In this case, TPs at the victim link increases its counter by 1 (if the SINR drops by more than 3 dB). The aggressor senses transmission at this desired frequency and it also increases its count by 1. At this point, the first aggressor knows that it is the second to enter the network and it takes the FO provided

by the Q-learning algorithm for  $S_{u=2} = 1$ . When the second aggressor enters, it can again fully overlap with the victim at ( $u = 1$ )th link. Now, the number of aggressors from the perspective of the second aggressor is  $S_{u=3} = 1$ . Hence, it takes the same FO as the first aggressor, which results in both aggressors fully overlapping. Therefore, both aggressors increase their counters, i.e.,  $S_{u=2} = 2$  and  $S_{u=3} = 2$  and infer that there are two aggressors in the network. Since the victim has fully overlapped with both aggressors, its count will also be  $S_{u=1} = 2$ . The trained Q-table is stored in each TP-RP pair. Therefore, each link looks for the FO values that avoid full overlapping for the new count. Through this sequence, the aggressors will also keep in memory the order in which they enter the network to ensure they do not pick the same offset provided by the proposed algorithm.

The Q-learning algorithm is trained separately for different values of  $S_u$  in this study, adding another dimension to the Q-table. For every  $S_u$  value, all the aggressors choose from the set of possible actions and observe the corresponding rewards. The algorithm goes through a state-action value iteration process dictated by (6) and computes the optimal value of each state-action pair.

The states in our Q-learning algorithm are the current FOs for all the aggressors. There are two possible actions that an agent can take in every iteration during training - 1) *Change FO*: The RL algorithm can change the FO by a fraction of the frequency carrier spacing and 2) *Change dispersion parameter*: choose the filter dispersion parameter to reduce the interference. In this paper, we focus on changing the FO, while the use of the dispersion parameter in the action space is left as a future study. In each time step, the agent updates the  $Q(s, a)$  value, where  $s$  and  $a$  are one of the possible states and actions respectively, by recursively discounting future rewards and weighing them by a positive learning rate  $\beta$ :

$$Q_{\text{new}}(s_t, a_t) \leftarrow (1 - \beta)Q_{\text{old}}(s_t, a_t) + \beta[r + \gamma \max_{a' \in \mathcal{A}} Q_{\text{old}}(s_{t+1}, a')], \quad (6)$$

where  $s_t$  and  $s_{t+1}$  are the current and next state, respectively,  $a_t$  is the action taken at state  $s_t$ ,  $r$  is the reward for taking action  $a_t$  at state  $s_t$ ,  $\gamma \in [0, 1)$  is the discount parameter, and  $\mathcal{A}$  is the set of possible actions. After the training process, the algorithm converges to optimal Q-values for each state-action pair [15],  $Q^*(s, a)$  and the optimal policy can be obtained by acting greedily in every state  $s$  as

$$\pi^* = \arg \max_{a \in \mathcal{A}} Q^*(s, a). \quad (7)$$

In this study, the reward  $r$  is defined as the improvement in capacity due to the change in FO and/or the filter parameter, which can be expressed as

$$r = \lambda_1(C_t - C_{t-1}), \quad (8)$$

where  $C_t$  is the current sum-capacity of the entire channel calculated during simulations,  $C_{t-1}$  is the sum capacity in the previous iteration, and  $\lambda_1$  is a weight parameter that can be tuned. The change in the sum capacity after an iteration may be very small and  $\lambda_1$  allows us to amplify it for faster convergence. The Q-value at every state is the current capacity of

the desired user, given the current FO and the filter parameters for the aggressors. After training,  $u$ th TP-RP pair can look up the ideal FO and filter parameter based on  $S_u$ .

#### IV. MULTI-USER EFFICIENCY IN POT

The ME is a measure of the signal quality over the total interference in the network. In [16], *asymptotic* ME for the desired signal  $u$  is defined as  $\eta_u \triangleq \lim_{\sigma \rightarrow \infty} \frac{e_u(\sigma)}{G_u^2}$ , where  $e_u$  is the effective energy of the desired signal at the  $u$ th user and  $\sigma$  is the standard deviation of the noise in the channel. For a conventional detector (i.e., single-user matched filter [16]),  $\eta_u$  is simplified to  $\eta_u = \max^2 \left\{ 0, 1 - \frac{\sum_i G_i \varrho_{i,u}}{G_u} \right\}$ , where  $\varrho_{i,u}$  is the correlation factor between the desired user and the  $i$ th aggressor. For our case, the partially overlapping between the victim and aggressor links due to intentional FOs leads to correlation between them. So this intentional FOs can be considered analogous to the correlation factor. From (4), the energies of the desired signal of user  $u$ , the SI, and the CCI can be calculated as

$$E_S = G_u^2 A_{mkmk}^2, \quad (9)$$

$$E_{SI} = G_u^2 \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} A_{nlmk}^2, \quad (10)$$

$$E_{OI} = \sum_i G_i^2 \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} A_{nlmk}^2. \quad (11)$$

The effective symbol energy for the desired user after taking the interference into account is then calculated as

$$e_u = G_u^2 A_{mkmk}^2 - G_u^2 \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} A_{nlmk}^2 - \sum_i G_i^2 \sum_{l=-K+1}^{K-1} \sum_{n=0}^{N-1} A_{nlmk}^2. \quad (12)$$

As a result, the ME is obtained as

$$\eta_u = \frac{e_u}{G_u^2} = \max^2 \left\{ 0, 1 - \frac{\sqrt{E_{SI} + E_{OI}}}{G_u A_{mkmk}} \right\}. \quad (13)$$

#### V. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed learning method with computer simulations. We uniformly distribute the TPs and RPs in an area of 1 km<sup>2</sup> and pair them. Considering IoT applications, we assume a modestly small resource allocation for data with  $N = K = 12$  for all links in the network. We use quadrature phase shift keying (QPSK) modulation throughout the simulations. To calculate path loss, we use the free-space path loss model without loss of generality of the proposed scheme. The carrier frequency is 800 MHz with the bandwidth of 200 kHz for each TP-RP link. In current simulations, we keep the filter parameter to be constant ( $\rho = \alpha = 0.2$ ). The multi-path channel is modeled based on extended pedestrian A (EPA) specified in Long-Term Evolution (LTE) standards, unless otherwise stated.

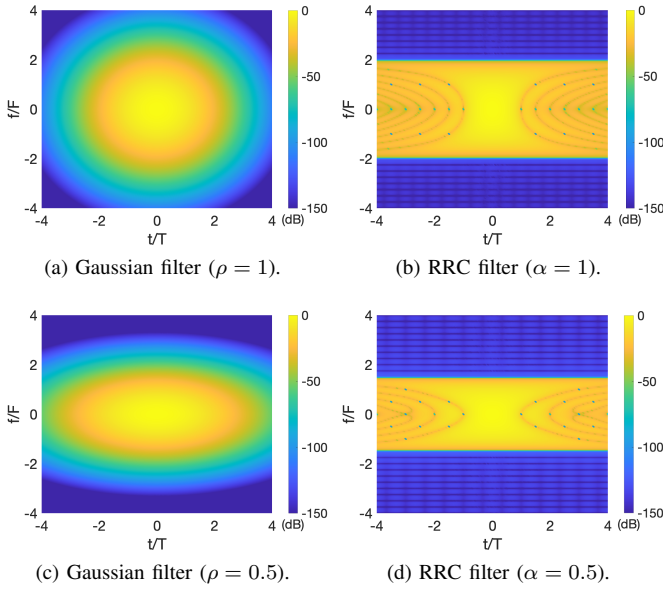


Figure 1. The ambiguity functions of RRC and Gaussian filters for different dispersion parameters.

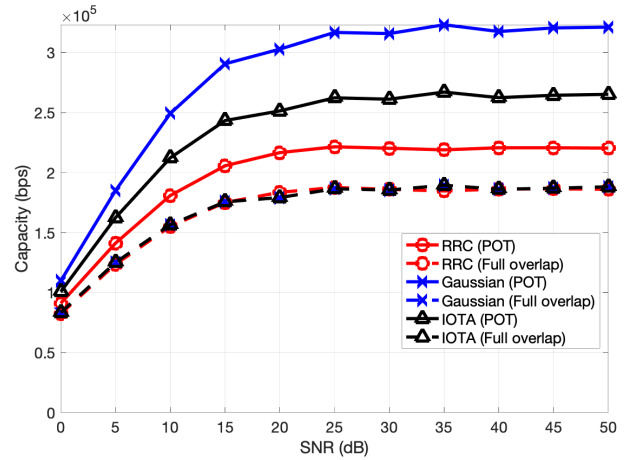
### A. Effect of Filter Parameters on Interference

In this section, we analyze the impact of the dispersion parameter on the systems performance. In Figure 1, we compare the ambiguity plots for RRC and Gaussian filters for different dispersion parameters. For Gaussian filter, when the dispersion parameter  $\rho = 1$ , the energy is distributed equally in the time and frequency domain (i.e., isotropic) as in Figure 1(a). However, when we set  $\rho = 0.5$ , the filter starts to *squeeze* in the frequency domain and expand in the time domain as in Figure 1(c). Thus, in the presence of CCI, a lower  $\rho$  allows other users to exploit available responses in frequency. However, this leads to the problem of SI at the desired link as it causes non-orthogonal pulses. We observe the same behaviour for RRC filter in Figure 1(b) and Figure 1(d) for different  $\alpha$  parameters. Dealing with this SI requires an equalizer being used at the receiver, and is out of the scope of this paper. The behavior of IOTA filter is similar to the Gaussian filter with more energy spreading in time and frequency. For the ambiguity function for IOTA filter, we refer the reader to [14].

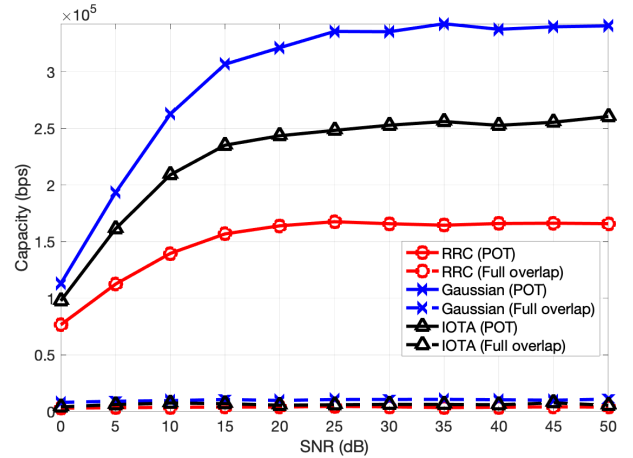
### B. Capacity Analysis

In Figure 2, we analyze the capacity at the desired link (i.e., average capacity over different instances for the same link) for different filters against when ten aggressors are present. Figure 2(a) and Figure 2(b) show the capacity curves for AWGN and EPA channel models, respectively. Among the filters, Gaussian filter provide the best capacity performance in both channel condition. As compared to fully overlapping, POT provides a gain by approximately 80% with Gaussian filters without noise being present. Another important observation is that the capacity reduces without POT under EPA channel model, while the drop is not substantial for POT.

In Figure 3, we investigate the capacity in the channel for a victim link against the number of aggressors present in the network at 10 dB signal-to-noise ratio (SNR). As



(a) Capacity under AWGN channel.



(b) Capacity under EPA channel model.

Figure 2. Capacity analysis for different channels.

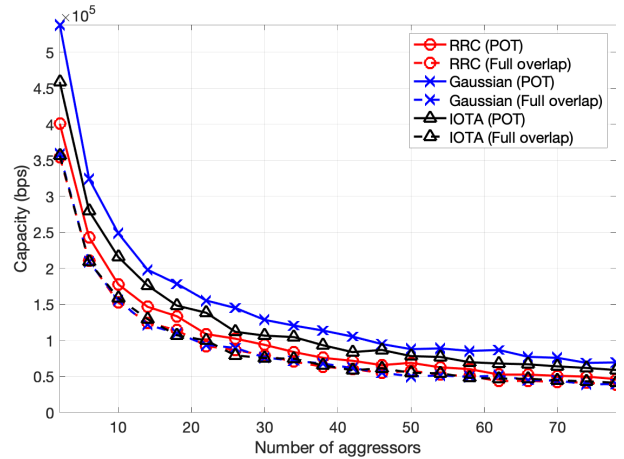


Figure 3. Capacity versus number of aggressors with POT and full overlap.

expected, with the increasing interference from the number of aggressors, the capacity reduces gradually. Gaussian filter offers a small improvement over other filters when POT are used. The capacity when full overlapping occurs is the same for any type of filter used. In the case of severe CCI (e.g. 1000 femtocells/km<sup>2</sup> [17]), the curves for full overlapping and

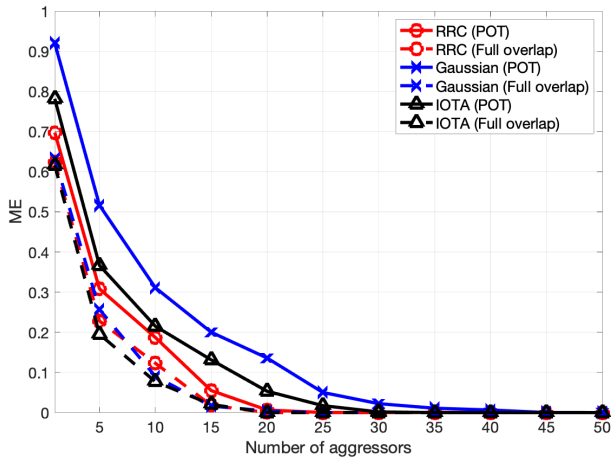


Figure 4. Multi-user efficiency with POT and full overlap.

partial overlapping converge.

### C. Multi-user Efficiency and Outage Probability

Based on (13), the ME for different filters at different interference conditions are plotted in Figure 4. Among all the filters, the Gaussian filter gives a higher ME for POT. For a full overlapping scenario, all filters provide similar ME. We present the outage probabilities in Figure 5. The outage occurs when the SINR of the desired signal falls below a pre-defined threshold denoted by  $\Gamma$ . Figure 5 shows the outage probabilities for  $\Gamma = -6$  (specified in LTE). The Gaussian filter outperforms other filters again and shows the least probability of outage against any number of aggressors.

## VI. CONCLUSION AND FUTURE WORK

Our preliminary results show the benefits of using POT in uncoordinated networks. We show that users to the order of  $10^2$  can be accommodated at a cost (to performance) significantly lower than the one with fully overlapping. We observe that a non-orthogonal filter (Gaussian) has the highest potential to fully exploit benefits offered by POT. Gaussian filters yield superior capacity against noise and interference, while also providing higher ME. For outage probability, using RRC filters provides similar performance but Gaussian slightly edges out ahead. The primary disadvantage of the proposed scheme is that the algorithm becomes more computationally expensive with the increasing number of aggressors, requiring a longer training period. As an extension of the current study, a detailed analysis of the computational complexity versus the accuracy of training is needed. Another extension is an *online* Q-learning algorithm which accounts for cases where a TP-RP pair has a wrong estimate of the number of aggressors. Incorporating an equalizer which addresses the inherent SI for non-orthogonal filters also needs to be investigated to improve the link performance further.

## REFERENCES

[1] Nokia, “Ultra Dense Network (UDN) White Paper,” Jun. 2016, available at <https://onestore.nokia.com/asset/200295>.

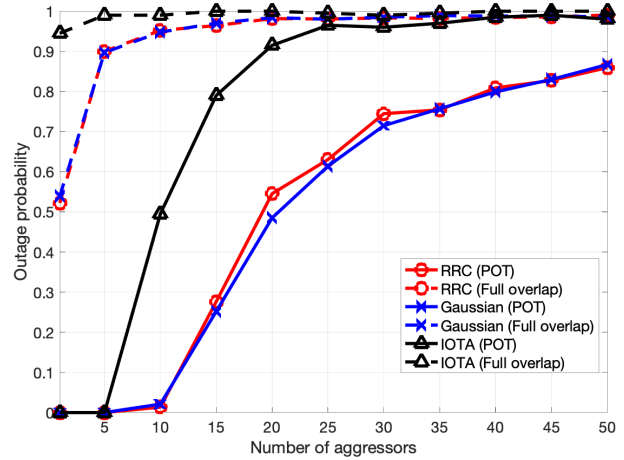


Figure 5. Outage probability with POT and full overlap ( $\Gamma = -6$  dB).

[2] F. Tang, Z. M. Fadlullah, N. Kato, F. Ono, and R. Miura, “AC-POCA: Anticoordination game based partially overlapping channels assignment in combined UAV and D2D-based networks,” *IEEE Trans. Veh. Technol.*, vol. 67, no. 2, pp. 1672–1683, Feb. 2018.

[3] Y. Ding, Y. Huang, G. Zeng, and L. Xiao, “Using partially overlapping channels to improve throughput in wireless mesh networks,” *IEEE Trans. Mobile Comput.*, vol. 11, no. 11, pp. 1720–1733, Nov. 2012.

[4] X. Zhao, L. Li, S. Geng, H. Zhang, and Y. Ma, “A link-based variable probability learning approach for partially overlapping channels assignment on multi-radio multi-channel wireless mesh information-centric IoT networks,” *IEEE Access*, vol. 7, pp. 45 137–45 145, 2019.

[5] A. Şahin, E. Bala, I. Güvenç, R. Yang, and H. Arslan, “Partially overlapping tones for uncoordinated networks,” *IEEE Trans. Commun.*, vol. 62, no. 9, pp. 3363–3375, Sep 2014.

[6] M. H. Yilmaz and H. Arslan, “Game theoretical partially overlapping filtered multi-tones in cognitive heterogeneous networks,” in *Proc. IEEE Military Commun. Conf. (MILCOM)*, Oct. 2014, pp. 411–415.

[7] Y. Wei, F. R. Yu, M. Song, and Z. Han, “User scheduling and resource allocation in HetNets with hybrid energy supply: An actor-critic reinforcement learning approach,” *IEEE Trans. Wireless Commun.*, vol. 17, no. 1, pp. 680–692, 2017.

[8] H. Ye, G. Y. Li, and B. F. Juang, “Deep reinforcement learning based resource allocation for V2V communications,” *IEEE Trans. Veh. Technol.*, vol. 68, no. 4, pp. 3163–3173, Apr. 2019.

[9] R. Amiri, H. Mehrpouyan, L. Fridman, R. K. Mallik, A. Nallanathan, and D. Matolak, “A machine learning approach for power allocation in HetNets considering QoS,” in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2018, pp. 1–7.

[10] M. Simsek, M. Bennis, and A. Czylik, “Dynamic inter-cell interference coordination in HetNets: A reinforcement learning approach,” in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2012, pp. 5446–5450.

[11] M. Bennis and D. Niyato, “A Q-learning based approach to interference avoidance in self-organized femtocell networks,” in *Proc. IEEE Global Commun. Conf. (GLOBECOM) Workshops*, Dec. 2010, pp. 706–710.

[12] F. B. Mismar, B. L. Evans, and A. Alkhateeb, “Deep reinforcement learning for 5G networks: Joint beamforming, power control, and interference coordination,” *IEEE Trans. Commun.*, pp. 1–1, 2019.

[13] F. Schaich and T. Wild, “Waveform contenders for 5G – OFDM vs. FBMC vs. UFMC,” in *Proc. IEEE Int. Sym. Commun. Control Signal Process. (ISCCSP)*, 2014, pp. 457–460.

[14] A. Şahin, I. Güvenç, and H. Arslan, “A survey on multicarrier communications: Prototype filters, lattice structures, and implementation aspects,” *IEEE Commun. Surveys Tuts.*, vol. 16, no. 3, pp. 1312–1338, Third 2014.

[15] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3–4, pp. 279–292, 1992.

[16] S. Verdu *et al.*, *Multuser Detection*. Cambridge University Press, 1998.

[17] Y. Teng, M. Liu, F. R. Yu, V. C. M. Leung, M. Song, and Y. Zhang, “Resource allocation for ultra-dense networks: A survey, some research issues and challenges,” *IEEE Commun. Surveys Tuts.*, vol. 21, no. 3, pp. 2134–2168, Third quarter 2019.