# Resource Allocation for a Wireless Coexistence Management System Based on Reinforcement Learning

Philip Söffker, Dimitri Block, Nico Wiebusch, Uwe Meier inIT, Institute of Industrial Information Technologies

OWL University of Applied Sciences

Liebigstraße 87

32657 Lemgo, Germany

Email: {philip.soeffker, dimitri.block, nico.wiebusch, uwe.meier}@hs-owl.de

Abstract—In industrial environments an increasing amount of wireless devices are used, which utilize licence-free bands. As a consequence of this mutual interferences of wireless systems might decrease the state of coexistence. Therefore, a central coexistence management system is needed, which allocates conflict-free resources to wireless systems. To ensure a conflict-free resource utilization, it is useful to predict the prospective medium utilisation before resources are allocated. This paper presents a self learning concept, which is based on reinforcement learning. A simulative evaluation of reinforcement learning agents based on neural networks, called deep Q-networks and double deep Q-networks, was realised for exemplary and practically relevant coexistence scenarios. The evaluation of the double deep Q-network showed, that a prediction accuracy of at least 98 % can be reached in all investigated scenarios.

### I. INTRODUCTION

License-free radio frequency (RF) bands such as the 2.4-GHz-ISM band are shared between incompatible heterogeneous wireless communication systems. In industrial environments, typically standardized wireless communication systems (WCSs) within this band are wide-band high-rate IEEE 802.11 called wireless local area network (WLAN), narrow-band low-rate IEEE 802.15.4-based WirelessHART and ISA 100.11a, and IEEE 802.15.1-related PNO WSAN-FA and Bluetooth (BT). Additionally, the spectrum band is shared with many proprietary wireless technologies (WTs) which target specific application requirements. Hence, sharing the spectrum may cause interferences between these heterogeneous WTs.

Therefore, the norm IEC 62657-2 [1] for industrial WCSs recommends an active coexistence management for reliable medium utilization and mitigation of interferences. The IEC recommends the use of a (i) manual, (ii) automatic noncooperative or (iii) automatic cooperative coexistence management. The first approach is the most inefficient one, due to time-consuming complex configuration effort. The automatic approaches (ii) and (iii) enable efficient self-reconfiguration without manual intervention and radio-specific expertise. An automatic cooperative coexistence management (iii) requires a control channel, i.e. a logical common communication connection between each coexisting wireless system to enable deterministic medium access. In case of a single legacy coexisting wireless system without such connection, the noncooperative approach (ii) is recommended. Non-cooperative coexistence management approaches may also utilize cooperative WCSs but are able to react on non-cooperative

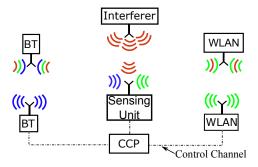


Fig. 1. Structure of a central coexistence management with two WCSs and one interferer

WCS which cause independent interferences. Such coexistence managements require approaches, which mitigate temporary interferences but also predict future medium utilization. Hence, a self-optimizing resource allocation behavior of the remaining cooperative WCSs is required.

Additionally, heterogeneity in industrial environments leads to high complexity. Hence, coexistence management requires self-learning approaches which models the dynamic heterogeneity of the industrial environment.

In particular, the actual non-cooperative coexistence managements are able to manage the resource allocation of connected WCS. The control channel can be used for a bidirectional communication of arbitrary WCSs and a dedicated coexistence management entity, which is called central coordination point (CCP). So the CCP can allocate resources to the different WCSs. This coexistence management principle is shown in Fig. 1 with BT and WLAN as heterogeneous WCSs. There are also WCSs which are not connected to the CCP. So the coexistence management cannot control these WCSs. Therefore, these uncontrollable WCSs are called interferers. A sensing unit observes the current RF spectral emissions as resource information for the CCP. The CCP uses these current informations to allocate free resources to the WCSs, without predicting the possible utilization of these allocated resources.

In this paper we propose a resource allocation concept for industrial non-cooperative coexistence management with a reinforcement learning (RL) approach. This approach learns to predict the future medium utilization. Additionally, self-optimization improves the resource allocation behavior.

In general, RL [2] targets self-optimization and prediction problems of agents which interact with the surrounding envi-

ronment. The agents observe the environment, take decisions based on the observations, and are rewarded therefore by the environment. A RL approach enables self-optimization without requirements of certain problem-domain knowledge.

For resource allocation within industrial coexistence management systems, a RL agent observes RF spectral emissions from utilized as well as from interfering WCSs. Then, the agent has to take decisions for resource allocation, which involves for example spectral, temporal and transmission power adjustments for the utilized WCSs. The WCSs apply and evaluate the adjustments based on various quality indicators such as link quality indicator (LQI) and packet loss ratio (PLR). Based on the quality indicators the RL agent get rewarded for its decision.

The following section II presents the related work. Section III will explain the concept of reinforcement learning based resource allocation for a central coexistence management system. This leads to section IV, where the presented concept will be simulated for exemplary and practically relevant coexistence scenarios. The results of this simulation will be presented in section V. Finally, section VI concludes the paper.

### II. RELATED WORK

The requirement for a self-optimizing resource utilization approach based on RL was already proposed by Ren et al. [3] in 2010. They use a special Q-learning algorithm [2] to find and predict non-utilized time intervals called whitespaces within licensed RF bands. Moreover, the distributed WCSs do this prediction autonomously. Additionally, there is no consultation between the individual distributed WCSs and no management entity. So each WCS acts opportunistically.

Liu et al. [4] propose a RL coexistence management approach for a time-slotted medium access of LTE-U and WLAN systems in license-free RF bands. Therefor, they also use a Q-learning algorithm to allocate dynamically free time slots to the LTE-U and WLAN systems. They assume, that there are no other WCSs with different WTs, which use the RF band at the same time. So it is an exclusively occupied RF band. Even if there are other WCSs with different WTs it is assumed, that these technologies are known and a influence on them is possible. So it is a cooperative coexistence management, which cannot handle non-cooperative WCS. In [4] the individual WCSs communicate directly with each other to negotiate for time slots. They do not use a central management entity.

A neural net based RL approach is used in [5] for a distributed medium access. In that approach one medium utilization solution is trained at a single central unit for all distributed WCS. The approach is limited to orthogonal resource utilization, which can not be assumed for heterogeneous WTs. This single trained solution is transferred to all distributed WCSs and rarely updated. Afterwards every WCS uses the trained solution to access the medium independently. So there is no consultation among the several distributed WCSs or between WCSs and central unit to manage the spectrum access.

All showed approaches act opportunistic and predict only for themselves or their WT which resources will be occupied in the future. Hence, there is no RL approach, which proactively predicts the medium utilization for heterogeneous WCSs in a central non-cooperative coexistence management.

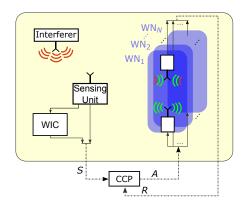


Fig. 2. Architecture of the RL-based coexistence management environment, and interactions between CCP and coexistence management environment

#### III. CONCEPT

The allocation of conflict-free resources is a fundamental part of a coexistence management system. However, before a resource allocation can be executed, a prediction of the future medium utilization is required. Such a central prediction and resource allocation is achieved in this paper with a RL approach. Hence, the coexistence management system can be described by a RL structure. As mentioned before, such a RL structure basically consist of two parts. The first part is the environment. This environment is a coexistence management environment, e.g. a shop floor. The second part is the agent. This agent is the CCP of the coexistence management system.

# A. Coexistence Management Environment

Shop floor environments are often equipped with metallic objects like machines. These cause reflections, absorptions and dispersions. Additionally, the environment contains many WCSs. These WCSs are usually organized in wireless networks (WNs). Thereby, N WNs use M different WTs, with  $M \leq N$ . Despite the heterogeneity of the WNs, a control channel is required to link the different WNs with the CCP. Such a control channel can be realized, like in [6], with Ethernet, whereby simple network management protocol (SNMP) is used. Furthermore, there are emitting interferers in the environment. Such an environment is illustrated in Fig. 2, with the CCP, WNs and interferers.

The actual state of the coexistence management environment is captured for the CCP as observation S. This observation consists of two separate observation elements. The first element captures the RF spectral emissions from WNs and interferers with the aid of a sensing unit. This element consists of discrete-time samples which are captured for a fixed time interval. Afterwards, these samples are transformed with a fast Fourier transform (FFT) and the magnitude of the spectrum is computed to  $|S_{\text{FFT}}|$ . Hence, the magnitude of the spectrum is like a snapshot of the coexistence management environment. The second observation element is an analysis of the captured RF spectral emissions for the same fixed time interval. For example this analysis classifies the captured spectrum on occupied frequency channels of a-priori known WTs. Hence, non-cooperative interfering WCSs can be classified. This is called wireless interference classification (WIC). It is possible to classify different WTs simultaneously, which is helpful in crowded wireless environments. Such WICs are neuro-fuzzy signal classifier (NFSC) [7] or convolutional neural network

(CNN) approaches [8], [9]. So the observation can be written as:

$$\mathbf{S} = \begin{pmatrix} |S_{\text{FFT}}| \\ S_{\text{WIC}} \end{pmatrix} \tag{1}$$

Based on the observation, the agent performs actions  $\boldsymbol{A}$ . For a wireless coexistence management system these actions are resource allocations. They are allocated by the CCP. Each WN gets its own dedicated resource a allocated:

$$\mathbf{A} = \begin{pmatrix} a_1 & a_2 & a_3 & \dots & a_N \end{pmatrix}^{\mathrm{T}} \tag{2}$$

Some WNs use a static channel selection such as one based on WLAN. These WNs get for example an allocation of a single resource, which is a frequency channel. However, other WNs use frequency hopping such as one based on BT. These WNs get for example an allocation of multiple frequency channels. So a bunch of channels are allocated, whereby the WN can select its own specific channel among them. Thus there are two types of resources, which have to be handled by the CCP.

The WNs use these allocated resources which have to be evaluated with a reward R. The reward is derived from the quality of data transmission on the allocated resources. This quality of data transmission can be expressed as a metric like the quality-of-coexistence (QoC) parameter [10]. The QoC parameter aggregates transmission-related characteristics, such as transmission time, update time and PLR. Each WN evaluates the QoC for itself with the scalar value  $QoC_{WN_i}$ . If other interfering systems use these allocated resource too, then the value of the QoC will decrease. This feedback is used, to validate the resource allocation of the CCP to each WN as reward:

$$\mathbf{R} = (R_1 \quad R_2 \quad R_3 \quad \dots \quad R_N)^{\mathrm{T}} \ \forall R_i = QoC_{\mathrm{WN}_i}$$
 (3)

Hence, this QoC-parameter is a feedback for each resource allocation.

#### B. Central Coordination Point

The CCP is the RL-agent. It learns to allocate conflict free resources to the WNs. However, at the beginning the CCP has no a-priori knowledge about the coexistence management environment. Therefore, the CCP has to optimize its resource allocation and learn to predict the occupancy of resources.

The CCP interacts with the environment in multiple steps. At the beginning of each step the CCP checks the initial observation *S*. Based on this observation, the CCP takes decisions for actions *A* for all WNs. The WNs perform their data transmissions on the resources, which were allocated by the actions. Then each data transmission is individually rewarded with the QoC parameter, which is an evaluation of the actions. Because of these actions the CCP needs a new observation for evaluating the new environmental situation. This new observation is the last part of each step and is also the basis for the next step. Each step follows this mentioned order.

The CCP interacts in two phases with the environment. The first phase is the training phase, wherein the CCP is trained for a particular environment. The second phase is the operational phase, wherein the CCP has completed its training but still has the ability for minor optimizations. These two phases are divided into episodes. Each episode has a predefined number of steps, e.g. 20 steps form an episode. After each episode the

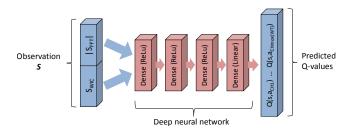


Fig. 3. Architecture of a DQN or DDQN With observation as input and predicted quality of WT specific resources as output

environment can be reset, which are like arbitrary processing time gaps. Meanwhile the interferer may change their resource, which is a challenge for the CCP. The CCP keeps his learned knowledge, despite the reset of the environment. Hence, the CCP learns a policy to maximize its reward. Thus it learns when to allocate which conflict-free resource to what WN, by predicting the occupation of resources due to interferers. This process of autonomously learning the behavior of interferers and the consequent opportunity of a prediction and allocation of free resources, is the advantage of a RL based central coexistence management system.

The CCP applies learning approaches. A well known technique of learning is Q-learning [2]. It uses Q-values, which addresses in this context the predicted quality of WT specific resources, a learning rate  $\alpha$ , which describes how fast newly learned knowledge overrides old knowledge, and a discount factor  $\gamma$ , which weights the relevance of immediate and future rewards. So the CCP learns by comparing predicted Q-values with the reward of the actual executed actions. This Q-learning can be extended with neural networks. Such an extension is helpful for large observation spaces and is called deep Qnetwork (DQN) [11]. These large observation spaces are also given with a large FFT length. Another problem of wireless environments is, that they are noisy. For such cases van Hasselt [12] proposed the double Q-learning approach, which uses two independent Q-functions. So even if one Q-function is biased it is not correlated to the other one, which reduces the prediction error. This double Q-learning was also combined with neural networks, which leads to double deep Q-networks (DDQNs) [13]. The CCP either uses a DQN or a DDQN. Each of them uses a deep neural network (DNN), as pictured in Fig. 3. Additionally experience replay is used for both CCP types and a target network [13] is used at the DDQN.

# IV. SIMULATION & RESULTS

The similation can be divided into coexistence management environment and CCP. For the environment the signal processing framework GNU Radio<sup>1</sup> is used. It enables the simulation of multiple WCSs, interferers and noise. For the proof of concept a lean realization is addressed which could be scaled later on. Therefore, the environment only processes synchronous streams without asynchronous events.

The lean realization is limited to a single WN, i.e. M=1 This WN consist of two WCSs which utilize unidirectional communications. The WN uses a static channel selection with four possible channels. These four channels are the action space. Then, the transmitting WCS applies a phase-shift keying (PSK) modulation. Afterwards the receiving WCS

<sup>&</sup>lt;sup>1</sup>gnuradio.org version 3.7.11 (27 Feb., 2017)

demodulates the transmitted data. Additionally for simplicity, it calculates the reward from the bit error rate (BER) for each step with  $R=1-{\rm BER}$ .

The data transmission is efficiently disturbed by an interferer. For a worst-case disturbance, this interferer inverses the PSK modulation of the WN. Hence, the power spectral densities (PSDs) of the interferer and WN are almost indistinguishable. This interferer will be a problem for the WN, if both use the same frequency channel. For the interference two coexistence scenarios are used: (i) static interferer, where the interferer occupies a channel for the duration of an episode like WLAN, and (ii) sequential hopping interferer, where the interferer sequentially changes the channel after each step like WirelessHART. In both scenarios the interferer chooses for each episode the initial channel randomly.

The sensing unit capture the complete band of all four channels for the duration of each step. It contains 1024 I/Q-samples, which is used for the generic observation element  $|S_{\rm FFT}|$ . The additional specific WIC observation element is omitted, because of the lean realization.

For the simulation of the DQN and DDQN CCP OpenAI Gym<sup>2</sup> is used. Both networks use a DNN with four dense layers with the output size of 256, 64, 32 and 4, respectively. Additionally, they use experience replay and the DDQN uses a target network. Further hyperparameters of the networks are listed in Table I.

Each simulation experiment consist of 250 episodes and is repeated 15 times. The first 100 episodes are the training phase. It takes place with the help of an  $\mathcal{E}$ -greedy exploration approach. The remaining 150 episodes are the operational phase. Each episode is separated into 20 steps. So, random frequency channel choice results in an average accumulated reward of 15.

TABLE I HYPERPARAMETERS OF THE CCP

Hyperparameters	Values
Learning rate $\alpha$	0,0001
Discount factor $\gamma$	0,96
Minibatch size	32
Initial exploration (training phase)	1
Final exploration (operational phase)	0,01
Update frequency of target network	20 episodes

The results of the simulation are shown in Fig. 4. For a better comparison of the operational phase a mean reward is estimated. This mean reward is shown in Fig. 4 as purple colored line. Thereby, the DDQN CCP reaches in all scenarios a mean reward of at least 19,6 which corresponds to a prediction accuracy of circa 98%, see Fig. 4(b) and 4(d). That implies, nearly all 20 data transmissions in each episode are received correctly, because of an appropriate resource allocation by the CCP. The DQN CCP in the first scenario has a mean reward of 18,08 which is not noteworthy in comparison to the other DQN and DDQNs.

For the training phase the exponential learning behavior is estimated as shown in Fig. 4 with the green solid line. The rise time of the exponential behavior is shown in the same figure as green dashed line. The trainings phase of the DQN in scenario (ii) has a rise time of In the second scenario the DDQN CCP results in a rise time of 31.8 episodes. It is 8.9

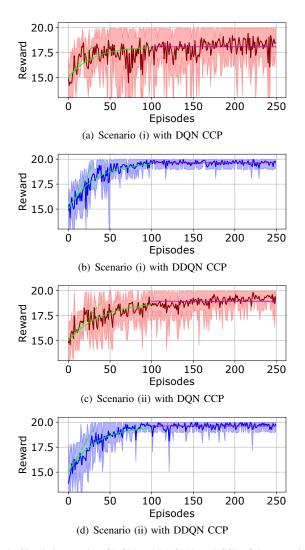


Fig. 4. Simulation results of DQN and DDQN based CCP of the scenario (i) and (ii); the solid lines are the mean reward, the shaded areas represent the values between the 10% and 90% percentile

episodes faster than the DQN. Within the first scenario, the DDQN also shows a comparable rise time of 30.3 episodes.

## V. Conclusion

In this paper we suggested a new concept for a selfoptimizing central coexistence management system. This concept used a deep RL approach, which learns to predict future medium utilizations and concludes this knowledge to allocate unoccupied frequency channels to WNs. The RLbased concept was evaluated with a simulation as proof of concept. This simulation considered a wireless environment with practically relevant coexistence scenarios. Additionally, the CCP applied DQN and DDQN based RL-agents. The DDQN shows a faster learning process during training and a higher prediction accuracy during operational phase in both scenarios. Hence, the advantage of such a RL-based central coexistence management system is the autonomously learning of interferer behaviors, without any a-priori knowledge about the wireless environment, and the consequential prediction and allocation of unoccupied resources.

 $<sup>^2</sup>$ gym.openai.com version 0.7.4 (5th Mar., 2017)

In the future software-defined radios have to be integrated to interact with real wireless environments.

#### ACKNOWLEDGEMENTS

Parts of this research were funded by KoMe (IGF 18350 BG/3 over DFAM) and HiFlecs (16KIS0266 over BMBF).

### REFERENCES

- IEC, "Industrial communication networks" wireless communication networks" part 2: Coexistence management," 2013.
   R. Sutton and A. Barto, "Reinforcement learning: An introduction: Second edition," 24.03.2018. [Online]. Available: http://incompleteideas. net/book/the-book-2nd.html
- [3] Y. Ren, P. Dmochowski, and P. Komisarczuk, "Analysis and implementation of reinforcement learning on a gnu radio cognitive radio platform," in Cognitive Radio Oriented Wireless Networks & Communications (CROWNCOM), 2010 Proceedings of the Fifth International Conference on, 2010.
- [4] Y.-Y. Liu and S.-J. Yoo, "Dynamic resource allocation using reinforcement learning for lte-u and wifi in the unlicensed spectrum," in 2017 Ninth International Conference on Ubiquitous and Future Networks (ICUFN), 2017, pp. 471–475.
- [5] O. Naparstek and K. Cohen, "Deep multi-user reinforcement learning for dynamic spectrum access in multichannel wireless networks," in GLOBECOM 2017 - 2017 IEEE Global Communications Conference, 2017, pp. 1-7.
- [6] N. Wiebusch, D. Block, and U. Meier, "A centralized cooperative snmp-based coexistence management approach for industrial wireless systems," in 2017 IEEE 13th International Workshop on Factory Communication Systems (WFCS), 2017, pp. 1–4.
  [7] D. Block, D. Tows, and U. Meier, "Implementation of efficient real-time
- industrial wireless interference identification algorithms with fuzzified neural networks," in 2016 24th European Signal Processing Conference (EUSIPCO), 2016, pp. 1738-1742.
- [8] M. Schmidt, D. Block, and U. Meier, "Wireless interference identification with convolutional neural networks," in 2017 IEEE 15th International Conference on Industrial Informatics (INDIN), 2017, pp. 180-185.
- [9] S. Grunau, D. Block, and U. Meier, "Multi-label wireless interference identification with convolutional neural networks," 2018 IEEE 16th International Conference 2018 in Press, 2018.
- [10] N. Wiebusch, P. Šoffker, D. Block, and U. Meier, "A multidimensional resource allocation concept for wireless coexistence management," in 2017 22nd IEEE International Conference on Emerging Technologies and Factory Automation, 2017, pp. 1-4.
  [11] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G.
- Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," Nature, vol. 518, no. 7540, pp. 529-533,
- [12] H. Hado V., "Double q-learning," in Advances in Neural Information Processing Systems 23, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc, 2010, pp. 2613-2621.
- [13] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in 2016 Proceedings of the Thirtieth AAAI, 2016, pp. 2094–2100.