# Distributed Resource Allocation for URLLC in IIoT Scenarios: A Multi-Armed Bandit Approach

Francesco Pase*, Marco Giordani*, Giampaolo Cuozzo°, Sara Cavallero°,
Joseph Eichinger†, Roberto Verdone°, Michele Zorzi*

*WiLab and University of Padova, Italy. Email: {name.surname}@dei.unipd.it
°WiLab and University of Bologna, Italy. Email: {name.surname}@unibo.it
†Huawei Technologies, Munich Research Center, Germany. Email: joseph.eichinger@huawei.com

*Abstract*—This paper addresses the problem of enabling inter-machine Ultra-Reliable Low-Latency Communication (URLLC) in future 6G Industrial Internet of Things (IIoT) networks. As far as the Radio Access Network (RAN) is concerned, centralized pre-configured resource allocation requires scheduling grants to be disseminated to the User Equipments (UEs) before uplink transmissions, which is not efficient for URLLC, especially in case of flexible/unpredictable traffic. To alleviate this burden, we study a distributed, user-centric scheme based on machine learning in which UEs autonomously select their uplink radio resources without the need to wait for scheduling grants or preconfiguration of connections. Using simulation, we demonstrate that a Multi-Armed Bandit (MAB) approach represents a desirable solution to allocate resources with URLLC in mind in an IIoT environment, in case of both periodic and aperiodic traffic, even considering highly populated networks and aggressive traffic.

*Index Terms*—6G, URLLC, Industrial IoT (IIoT), resource allocation, machine learning, Multi-Armed Bandit (MAB).

## I. INTRODUCTION

With early 5th generation (5G) deployments already rolled out, the research community is discussing use cases, requirements, and enabling technologies towards 6th generation (6G) systems [1]. Among other services, 6G will introduce new communication interfaces and innovative architectures to support the Industrial Internet of Things (IIoT) in 2030 and beyond, where the 6G network connects sensors and machines in factories, plants, mines, to enable analytics, diagnostics, monitoring, asset tracking, as well as process, regulatory, supervisory, and safety control [2]. In this context, the need for robots to complete cooperative operations that require high precision and coordination in real time comes with its own set of requirements, e.g., in terms of reliability (up to 99.99999%) and latency (below 1 ms, or even 0.1 ms, in the radio part), making it crucial to support Ultra-Reliable Low-Latency Communication (URLLC) in the industrial domain [3]. The factory of the future will further operate to support high-density deployments of machines and end users.

In this context, the time introduced by the Radio Access Network (RAN) operations, from routing and scheduling to resource allocation and modulation, represents one of the most impactful latency components. Specifically, a *centralized pre-configured* scheduling protocol usually requires the prior exchange of scheduling requests (grants) to (from) the Next Generation Node Base (gNB), which is not compatible with URLLC in IIoT scenarios [4], [5]. To partially address this issue, 3GPP NR supports *semi-persistent* and *grant-free* communication in the uplink (UL) [6], in which the network pre-allocates radio resources, thereby eliminating the need for User Equipments (UEs) to wait for network grants before transmission. However, reserving resources to dedicated UEs can be inefficient if traffic demands are aperiodic [7], and it is not possible to anticipate when resources will be needed [8].

Another solution is to design a *user-centric* architecture (as foreseen in 6G [9]) in which end machines make autonomous decisions, "disaggregated" from the network [10]. Along these lines, in this paper we explore the feasibility of a decentralized/distributed scheduling algorithm that, exploiting machine learning (ML) technologies, allows UEs to optimize their UL transmission strategies by autonomously selecting the available physical resources. This framework is able to learn from the application, and could work well even considering architectures for IIoT scenarios in which communication is on the sidelink, with no or limited support from the gNB [11].

Despite this potential, however, distributed scheduling may create collisions during communication, raising the question of whether this approach is compatible with URLLC applications. To this aim, we apply the Multi-Armed Bandit (MAB) theory [12] to evaluate how autonomous machines should select transmission resources based on previous scheduling decisions and the effect they produced on the network in terms of reliability. While the MAB approach is well known, most related work focused on downlink (DL) [13], cellular [14], or IoT [15] networks. In turn, we consider a UL scenario modeled according to the "Motion Control" 5G-ACIA geometry (in which a remote server sends commands to control the moving parts of machines), thus ensuring that our results are representative of a typical IIoT environment. Other notable papers consider vehicular scenarios [16], [17], where the target is to enable URLLC for vehicle-to-vehicle communications via Deep Reinforcement Learning (DRL). However, we argue that for IIoT use cases, state-of-the-art MAB algorithms may better exploit the strong correlation typical of the industrial environment while, at the same time, reducing the computational complexity and training time to converge to optimal solutions, compared to more sophisticated
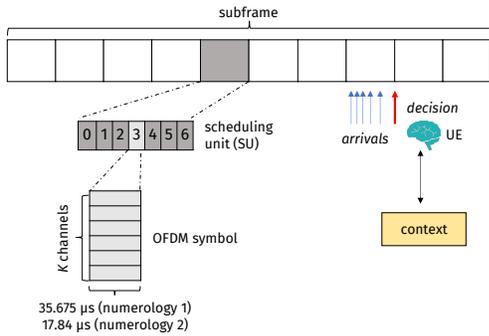
Fig. 1: Transmission structure.

DRL alternatives.

We perform simulations with both periodic and aperiodic traffic, and as a function of the UEs' density and spatial distribution, the traffic periodicity (thereby modeling aggressive or conservative applications), and the transmit power, thus considering a low-power performance regime. From our results, we conclude that the Thompson Sampling agent [18] is a promising candidate method to minimize the collision probability even in the presence of unscheduled transmissions.

The rest of the paper is organized as follows. In Sec. II we present the distributed resource allocation problem, in Sec. III we introduce possible ML methods based on MAB to solve it, and in Sec. IV we describe our simulation setup and discuss our main results. Finally, Sec. V concludes our work with suggestions for future research.

## II. PROBLEM FORMULATION AND SYSTEM MODEL

We consider an Orthogonal Frequency Division Multiplexing (OFDM) system in which devices, also denoted as agents in machine learning parlance, are located in a factory environment, and have to autonomously choose the orthogonal channel to be used for UL transmissions. The time domain is discretized into intervals of duration equal to the OFDM symbol (with a Scheduling Unit (SU) consisting of 7 OFDM symbols), whose duration depends on the adopted NR numerology. The frequency domain is also discretized into $K$ orthogonal channels, whose size depends on the available bandwidth $B$ and the subcarrier spacing $\Delta f$.

At the beginning of each SU, the agents make their scheduling decisions, that is the channel to be used for transmission, as shown in Fig. 1. Unlike in a centralized pre-configured resource allocation approach, in which radio resources are scheduled by the gNB via scheduling grants, we study the feasibility of a decentralized algorithm based on ML in which each agent autonomously optimizes its channel selection policy relying only on the gNB feedback, without prior ad hoc message exchange with the gNB itself. The rationale behind this scheme is to exploit the underlying correlations typical of the IIoT traffic to avoid the transmission of centralized scheduling grants, thus reducing the end-to-end latency and promoting URLLC.

If multiple agents use the same physical channel during a specific SU, we assume that their packets are lost due to a *collision* event. At the end of each SU, the gNB broadcasts a message indicating in which channel(s) data were successfully received. This message is used by the pool of agents to optimize their subsequent decision strategies, and achieve coordination without communication.

We formalize the problem using the MAB framework, which is used to model many sequential decision processes in computer science and engineering [12]. In this particular multi-agent scenario, there are $N$ agents, i.e., the $N$ UEs, interacting with the same environment. Whenever an agent $n \in \{1, \ldots, N\}$ generates a new packet during SU $t$, it schedules its transmission at the beginning of SU $t + 1$, choosing one among the $K$ available channels, which will be used for transmission for the whole SU duration. According to the MAB notation, we refer to the action of using channel $k \in \mathcal{K} = \{1, \ldots, K\}$ as "playing the arm" $k$. At the end of SU $t + 1$, the message received from the gNB is converted into a reward $r_{n,t}$, indicating whether or not the transmission was successful, i.e., $r_{n,t} = 1$ or $r_{n,t} = 0$, respectively: maximizing the reward implies transmitting the data successfully in low latency, as there is no need to exchange scheduling grants between the UEs and the gNB, leading to the URLLC objective. In our model, we assume that the reward behind each action is sampled from a Bernoulli distribution with unknown parameter $\mu_n(k_{n,t})$, which depends on the action taken by the agent, and captures the probability of the other agents transmitting at the same time. Thus, in each SU $t$, agent $n$ samples an action $k \in \mathcal{K}$ according to its policy $\pi_n : \mathcal{H}^{t-1} \to \Delta_K$, which is, in general, a map from history $H_n(t-1) = \{(k_{1,n}, r_{1,n}), \ldots, (k_{t-1,n}, r_{t-1,n})\} \in \mathcal{H}^{t-1}$ to a probability distribution over the action set $\mathcal{K}$, where $\Delta_K$ denotes the $K$-simplex. The history vector $H_n(t-1)$ is used by the agent to optimize its policy $\pi_n$, so as to maximize the expected cumulative reward $R(\pi_n, T) = \mathbb{E}_{\pi_n} \left[ \sum_{t=1}^{T} \mu_n(k_{n,t}) \right]$.

## III. MULTI-ARMED BANDIT (MAB) AGENTS

To solve the problem in Sec. II and maximize the reward, many algorithms have been proposed in the literature over the past years [12]. In this paper, we study the performance of different MAB agents to solve the problem of distributed resource allocation, in the specific context of URLLC for IIoT.

*a) Random Agent (RA):* It implements the simplest decision scheme, and is used as a lower bound. Nonetheless, it represents well the case of 5G NR grant-free scheduling, where the access decision is random, and re-transmissions are optimized to achieve reliability [19]. In particular, in each SU, the RA selects uniformly, at random, one of the $K$ arms, and no learning is involved.

*b) UCB Agent (UCB-A):* It implements the Upper Confidence Bound (UCB) algorithm [20], i.e., the agent plays, in each SU $t$, the arm $k_t$ such that

$$k_t = \operatorname{argmax}_{k \in \mathcal{K}} \left[ Q_t(k) + c \sqrt{\frac{\log t}{n_t(k)}} \right], \quad (1)$$

where $Q_t(k)$ is the empirical average at step $t$ of the experienced rewards for arm $k$, $n_t(k)$ is the number of times arm $k$ has been played until time step $t$, and $c$ is an exploration

parameter to be optimized. In Eq. (1), $Q_t(k)$ represents the *exploitation* part, as it is related to the past experience, while $\sqrt{\log t / n_t(k)}$ quantifies the uncertainty around the empirical average, and decreases as we collect more samples, i.e., as $n_t(k)$ increases. The larger this second term for an action $k$, i.e., the uncertainty of its performance, the higher the probability of choosing that arm, meaning that we need more samples to have a good estimate of its related reward. This principle is also known as *"optimism in the face of uncertainty."*

*c) Thompson Sampling Agent (TS-A):* The agent adopts a *Bayesian inference* approach to identify the most promising arms. In particular, TS-A builds a distribution for each reward, thus modeling not only its mean, but the whole statistics [18]. Given that our problem includes a binary reward $\{0, 1\}$ behind each arm, it is quite natural to model the rewards according to a Bernoulli distribution, which is parameterized by the success probability vector $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_K)$, where $\mu_k$ represents the average unknown reward behind arm $k \in \mathcal{K}$. Following the Bayesian framework, parameter $\mu_k$ of arm $k$ is modeled as a Beta$(\alpha_k, \beta_k)$ random variable, where $\alpha_k$ counts the number of successful transmissions after playing arm $k$, and $\beta_k$ represents the number of collisions. Therefore, the mean of $\mu_k$ is equal to $\alpha_k / (\alpha_k + \beta_k)$. The Beta distribution parameters are initialized to $\{\alpha_k = 1, \beta_k = 1\}$ for all $k \in \{1, \ldots, K\}$.

As the TS-A collects more data, $\alpha_k$ and $\beta_k$ are updated accordingly, inducing biased probabilities for the different arms. These informed distributions are also called *posterior probabilities*, in Bayesian parlance. Whenever the agent makes a decision, i.e., it chooses a physical channel based on the probability of that channel not being accessed by other agents in that time interval, it samples a vector $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_K)$, and plays the arm $k^*$ such that $k^* = \operatorname{argmax}_k \{\mu_k\}$. This algorithm is known as the Thompson Sampling (TS) algorithm [18].

*d) Neural Agent (NA):* The NA is equipped with a small-size Neural Network (NN) used to represent its decision policy. In particular, the agent receives, as an input, context information $s_t \in \mathcal{S}$ from the environment, thus the problem is formulated as a contextual MAB, i.e., the average reward depends on the played arm $k_{n,t}$, and on the state $s_{n,t}$ [20]. The NN input represents the feedback on the results of the last transmission attempt, broadcast by the gNB. As such, the input data is a vector of $K + 1$ entries: the first $K$ values are the results of the transmission attempts in the $K$ orthogonal channels, whereas the last value indicates whether it is a first-time transmission or a re-transmission. Again, the $0/1$ reward given to failed/successful transmission, respectively, is used by the NA to optimize the NN parameters, and maximize the given rewards. The model is an adaptation of that in [21].

*Remark.* The UCB and TS algorithms exhibit good theoretical properties in terms of convergence time to optimal strategies, as long as some critical assumptions are satisfied [20]:

1) The rewards behind each action need to exhibit a sub-Gaussian distribution. Any distribution with limited support has this property, which is also verified in our setting.
2) The reward samples after playing action $k$ are i.i.d. This assumption is more critical in real scenarios, and in

particular in our problem. In fact, each agent interacts with many other devices, and so the rewards depend on the actions of the other agents, which are continuously learning and changing their decision schemes. This leads to highly non-stationary environments, meaning that the reward distribution may change over time. However, empirical results show that state-of-the-art MAB algorithms can still be applied even though the stationarity assumption for the rewards is not satisfied [22].

In Sec. IV-B we compare the performance of the MAB agents presented above, and provide guidelines towards the best schemes to satisfy URLLC requirements for IIoT.

## IV. PERFORMANCE EVALUATION

In this section, after introducing our simulation setup, we evaluate the performance of the proposed distributed resource allocation scheme implementing one of the MAB agents described in Sec. III, in different IIoT scenarios.
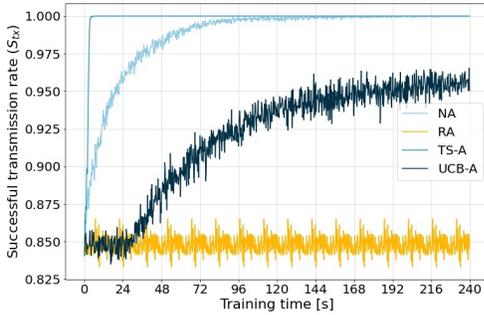
### A. Simulation Setup

End machines transmit at frequency $f_c = 3.5$ GHz and with a bandwidth of $B = 20$ MHz. The subcarrier spacing is set to $\Delta f = 30$ KHz (i.e., 3GPP NR numerology 1), which results in $K = 55$ orthogonal channels, and an OFDM symbols duration of $T_{\text{OFDM}} \simeq 35.675 \, \mu$s [23]. With an SU of 7 OFDM symbols, we get an SU duration of $T_{SU} \simeq 0.25$ ms. We assume that, whenever a packet is to be sent, it can be transmitted within one SU. If two or more UEs select the same UL channel for transmission in the same SU, we consider those packets to be lost (due to a collision event). Assuming that the gNB feedback (informing about the collision) is received within the current SU, the retransmission can be scheduled in the subsequent SU.
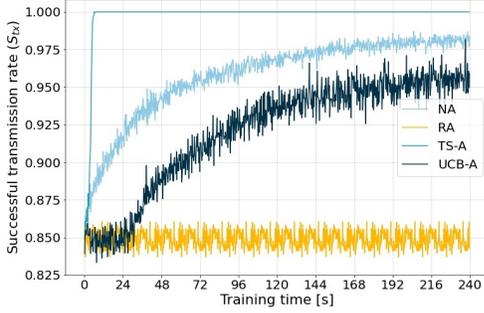
The factory floor is characterized according to the 5G-ACIA "Motion Control" scenario, as described in [11]. Hence, the geometry is modeled as a parallelepiped of length $\ell = 15$ m, width $w = 15$ m, and height $h = 3$ m, and machines are randomly and uniformly distributed inside the factory. The gNB is located at the center of the ceiling, and communicates with power $P_{\text{TX,DL}} = 30$ dBm. The transmit power of the UEs is set to $P_{\text{TX,UL}} \in \{8, 10, 23\}$ dBm. Also, we consider omnidirectional transmissions, therefore the antenna gain is fixed to $G = 1$ for both the UEs and the gNB. The channel model is based on the 3GPP Indoor Factory (InF) scenario [24], where UEs are assumed to communicate in Non-Line-of-Sight (NLOS) if the joining line between the UE's and the gNB's centers intersects one or more machines.

In our simulations, the traffic can be either *periodic* or *quasi-periodic*. In the first case, packets are generated at constant periodicity $\tau$. In the second case, the application still generates packets with periodicity $\tau$, upon which a random component $t_{\text{off}}$ of $\{-2, -1, 0, +1, +2\}$ OFDM symbols is added.

The performance of the different MAB agents' policies is assessed in terms of *successful transmission rate* $S_{TX}$, which indicates the ratio between the successfully received packets and the total number of attempts within one SU, averaged over $1\,000$ steps, as a function of the traffic periodicity $\tau$, the

(a) Periodic traffic.



(b) Quasi-periodic traffic

Fig. 2: $S_{TX}$ vs. the training time, for different MAB agents, with periodic and quasi-periodic traffic, $\tau = 1.5$, and $N = 50$.

number of UEs $N$, and the UL transmission power $P_{\text{TX,UL}}$. Notice that $S_{TX}$ is inversely proportional to the number of re-transmissions and, as such, represents well the theoretical rewards $r_{n,t}$ of the MAB agents.

### B. Numerical Results

**Impact of the training.** In Fig. 2 we analyzed the training curve of the agents with periodic and quasi-periodic traffic, with a periodicity $\tau = 1.5$ ms, and considering $N = 50$ UEs in the system, for a total training time of $T = 240$ s. For the periodic case, we observe from Fig. 2a that TS-A is the best performing agent. In particular, the TS agents are able to learn their optimal strategy, achieving zero collisions (i.e., $S_{TX} = 1$, our target for URLLC) in a very short training time ($< 10$ s). NA achieves a similar performance to that of TS-A, though after a longer training process. This is due to the fact that NA needs more interactions with the system to optimize the network parameters, thus slowing down the training phase. For UCB-A, the exploration parameter $c$ in Eq. (1) was set to 2, as it showed the most stable configurations in our experiments. Still, it results in an even slower convergence compared to NA, due to the fact that it struggles to achieve coordination. Also, UCB-A presents significant oscillations over time, due to the impact of collisions and retransmissions. As expected, RA (our baseline) performs poorly, and there is no improvement over time, as feedback signals are not exploited by the algorithm to adjust the access scheme.

For the quasi-periodic case, we observe from Fig. 2b that TS-A presents again the best performance despite the more complex scenario, converging to zero collisions within 15 s. Now, NA no longer achieves perfect convergence
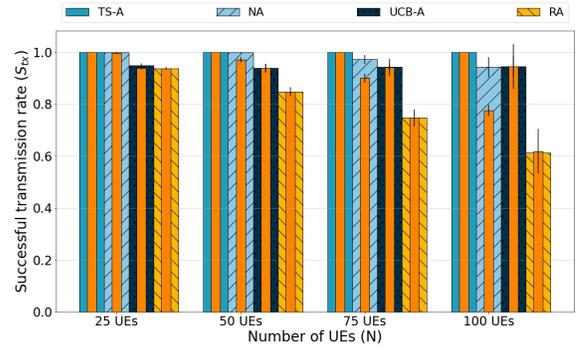


Fig. 3: $S_{TX} \pm$ one standard deviation vs. $N$ for different MAB agents, after a training time of 60 s, with $\tau = 1.5$ ms, with periodic (wide bars) and quasi-periodic (narrow bars) traffic.

within the training time, suggesting that it cannot work well in non-stationary multi-agent scenarios, or deal with non-deterministic traffic requests. However, we believe that, with a better tuned training process, and with more relevant context information as input, the final performance would reasonably improve. Finally, UCB-A and RA perform similarly to the case of periodic traffic.

**Impact of the number of users.** In Fig. 3 we evaluate the performance of the MAB agents as a function of $N \in \{25, 50, 75, 100\}$. In particular, we studied the statistics of the successful transmission rate $S_{TX}$ after 60 s of training, where again the total training time is set to $T = 240$ s. First, we observe that TS-A converges to the optimal scheme (i.e., $S_{TX} = 1$) within 60 s in all configurations, thus achieving coordination without communication even in dense ($N = 100$) networks. Second, NA outperforms UCB-A with periodic traffic, but suffers with quasi-periodic traffic: notably, $S_{TX}$ decreases by 10% in the quasi-periodic case, for $N = 100$. This is due to the fact that NA implements and exploits an NN to optimize its decisions, thus the learning phase can take more time in the most complex scenarios. Interestingly, compared to other agents, UCB-A's performance is less sensitive to $N$, and eventually outperforms NA's approach in the most crowded scenarios. On the downside, it exhibits wider oscillations, i.e., higher standard deviation in Fig. 3, an indication of a less stable behavior of the agent in non-stationary environments.

**Impact of the traffic periodicity.** Fig. 4 explores the effect of the traffic periodicity $\tau$ on the successful transmission rate $S_{TX}$. As expected, the more aggressive the traffic, the more difficult for the agents to achieve convergence, which is also highlighted by the increased standard deviation in all MAB configurations. Again, TS-A is the best agent, and can converge to the optimal scheme regardless of the value of $\tau$. Eventually, NA is also able to achieve zero collisions (i.e., $S_{TX} = 1$) when $\tau = 5$ ms in case of periodic traffic. Even the RA approach (our baseline) achieves a successful transmission rate of around 0.9 as $\tau$ grows, i.e., considering less bandwidth-hungry applications, thanks to the lower collision probability as the contention on the channel becomes less intense. Notably, UCB-A is the only method that improves the average accuracy as $\tau$ decreases: the shorter traffic periodicity implies more transmission attempts within the training time, which in turn
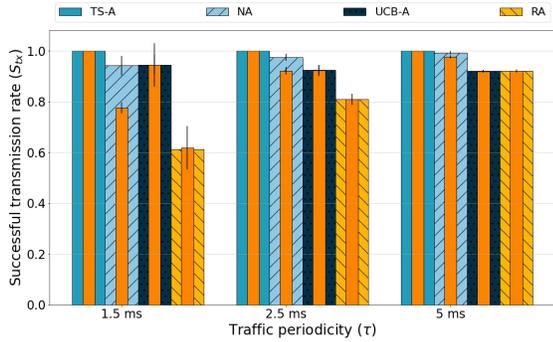
Fig. 4: $S_{TX} \pm$ one standard deviation vs. $\tau$ for different MAB agents, after a training time of 60 s, with $N = 100$, with periodic (wide bars) and quasi-periodic (narrow bars) traffic.
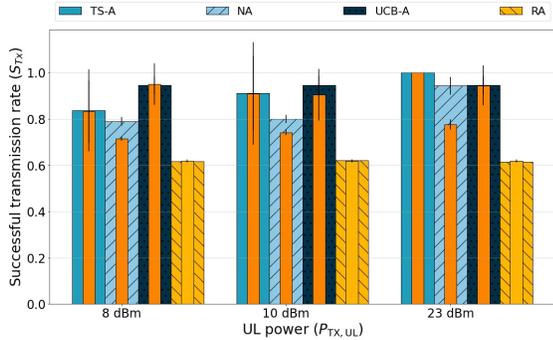


Fig. 5: $S_{TX} \pm$ one standard deviation vs. $P_{TX,UL}$ for different MAB agents, after a training time of 60 s, with $N = 100$ and $\tau = 1.5$ ms, with periodic (wide bars) and quasi-periodic (narrow bars) traffic.

provides more data to the agent to optimize its decisions. However, oscillations become significant when $\tau = 1.5$ ms.

**Impact of the UL transmission power.** IIoT devices, such as industrial sensors, may be subject to battery lifetime constraints. In light of this, we studied the impact of the UL transmission power $P_{TX,UL} \in \{8, 10, 23\}$ dBm on the MAB convergence. While decreasing $P_{TX,UL}$ promotes energy savings and mitigates interference, it may also lead to communication outage when the Signal to Interference plus Noise Ratio (SINR) goes below a pre-defined sensitivity threshold, set to $-5$ dB in our simulations. In Fig. 5, with $P_{TX,UL} = 23$ dBm, the outage probability is very small, leading to $S_{TX} \approx 1$ in most configurations (if convergence is achieved). As $P_{TX,UL}$ starts decreasing, outage events, besides collisions, lead to additional packet losses, and to a more complex environment. Unlike TS-A and NA, UCB-A is less sensitive to this effect. The reasons are twofold. On one side, NA converges slowly, and is more exposed to retransmissions. At the same time, TS-A converges quickly to a specific solution, meaning that unpredictable outage events may break the environment statistics underlying the TS algorithm, and lead to unexpected negative feedback from the gNB. On the contrary, UCB-A initially explores more, and can better adapt to new configurations in more dynamic scenarios. When $P_{TX,UL} = 8$ dBm, UCB-A is the best performing agent, and achieves $+16\%$ $S_{TX}$ compared to TS-A.

**TS-A performance.** In view of the above results, we further analyzed TS-A's convergence time to the optimal solution
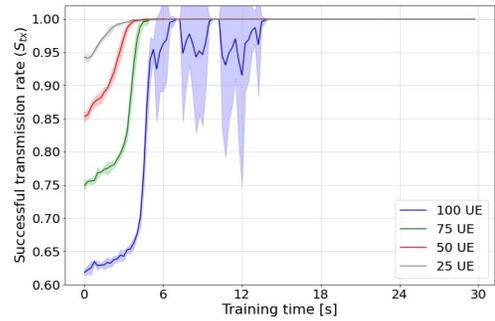


Fig. 6: $S_{TX}$ vs. the training time and as a function of $N$, for TS-A with periodic traffic, and $\tau = 1.5$ ms. The curves report mean $\pm$ standard deviation over the simulation runs.
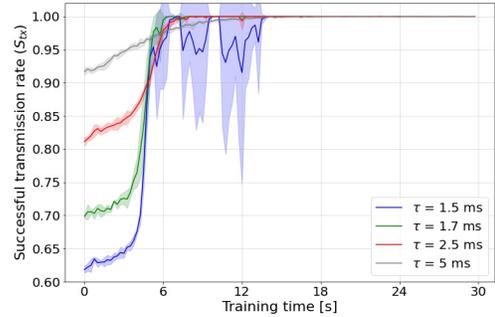


Fig. 7: $S_{TX}$ vs. the training time and as a function of $\tau$, for TS-A with periodic traffic, and $N = 100$. The curves report mean $\pm$ standard deviation over the simulation runs.

(where no collisions are experienced) as a function of (i) the number of users $N$, and (ii) the traffic periodicity $\tau$. In Fig. 6, we observe that, as $N$ increases, the TS algorithm takes more time to converge to the best solution, as expected. Notably, the curve with $N = 100$ presents the highest variance, due to the fact that many users are learning an individual policy, leading to a highly non-stationary environment.

In Fig. 7, we see that when $\tau$ decreases the convergence time grows accordingly, even though the gap among different configurations is relatively small (convergence is achieved after $\sim 8$ s). This is due to the fact that, on the one hand, when the traffic periodicity is short, the problem becomes more complex, as more packets have to be allocated. On the other hand, the agents receive more feedback signals within the same time interval, thus leveraging more data for the training.

### C. Final remarks

Our initial experiments confirm that there exists a MAB configuration for which distributed resource allocation can achieve zero collisions in low latency, i.e., without gNB scheduling grants, thus supporting URLLC.

In particular, TS-A is the best performing approach in case of both dense systems and aggressive aperiodic traffic (where conventional semi-persistent/grant-free NR schedulers may fail). Consequently, our experiments suggest that the Bayesian formulation, together with the exploration strategy of TS, are good starting points to build distributed resource allocation in real IIoT environments, reducing the latency introduced by centralized protocols. Interestingly, UCB-A works well in

complex scenarios, or when UEs communicate with limited power, thus supporting energy efficiency at the expense of some collisions. Moreover, the superior performance in terms of $S_{TX}$ of the MAB schemes against RA shows that machine learning can dramatically reduce, if not completely eliminate, the burden of re-transmissions introduced by 5G-NR-like grant-free access scheduling schemes [19].

However, distributed resource allocation requires longer training time before convergence, which in real IIoT systems may not be negligible. Still, the training could be run offline, which does not affect the real-time performance of the system (it can be executed when the machine is turned off, e.g., during the calibration of the electro-mechanical processes, or before the service is activated); once active, the service can run rapidly and without significant computational overhead. Moreover, our analysis evaluates the training time when the system starts the optimization process from scratch: faster adaptation can be achieved if the system faces limited changes with respect to the initial training scenario, e.g., some components join or leave the system. Nevertheless, the trained model still requires retraining when data distributions have deviated significantly from those of the original training set, which involves additional overhead [25]. This motivates further explorations in the case of more dynamic systems, that will be carried out as part of our future work.

## V. CONCLUSIONS AND FUTURE WORK

In this paper we studied the design of user-centric (rather than gNB-centric) distributed (rather than centralized) resource allocation in IIoT scenarios. This approach does not involve scheduling grants to be disseminated before UL transmissions, and is thus positioned to better support URLLC compared to conventional scheduling methods. We explored different state-of-the-art MAB agents, for the first time applied to the context of URLLC for IIoT, and identified TS-A as the best performing implementation, achieving zero collisions in our experiments. TS-A scales well with the number of users in the system compared to other MAB methods, and still achieves perfect accuracy even considering aperiodic traffic. Notably, UCB-A showed superior performance when the UEs communicate with low power, despite some collision events.

This work opens up new interesting research directions. For example, we will evaluate whether federated learning, which optimizes the scheduling policies based on the interaction among the UEs, would result in faster convergence than MAB.

## REFERENCES

[1] M. Giordani, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Toward 6G Networks: Use Cases and Technologies," *IEEE Communications Magazine*, vol. 58, no. 3, pp. 55–61, Mar. 2020.

[2] J. Lee, B. Bagheri, and H.-A. Kao, "A Cyber-Physical Systems architecture for Industry 4.0-based manufacturing systems," *Manufacturing Letters*, vol. 3, pp. 18 – 23, Jan. 2015.

[3] M. Wollschlaeger, T. Sauter, and J. Jasperneite, "The future of industrial communication: Automation networks in the era of the internet of things and industry 4.0," *IEEE Industrial Electronics Magazine*, vol. 11, no. 1, pp. 17–27, Mar. 2017.

[4] G. Cuozzo, S. Cavallero, F. Pase, M. Giordani, J. Eichinger, C. Buratti, R. Verdone, and M. Zorzi, "Enabling URLLC in 5G NR IIoT Networks: A Full-Stack End-to-End Analysis," in *Joint European Conference on Networks and Communications & 6G Summit (EuCNC/6G Summit)*, 2022.

[5] P. Yang, L. Kong, and G. Chen, "Spectrum sharing for 5G/6G URLLC: Research frontiers and standards," *IEEE Communications Standards Magazine*, vol. 5, no. 2, pp. 120–125, Apr. 2021.

[6] 3GPP, "NR; Medium Access Control (MAC) protocol specification – Release 15," 3GPP, Technical Specification (TS) 38.321, 2019.

[7] M. C. Lucas-Estañ, J. Gozalvez, and M. Sepulcre, "On the capacity of 5G NR grant-free scheduling with shared radio resources to support ultra-reliable and low-latency communications," *Sensors*, vol. 19, no. 16, p. 3575, Aug. 2019.

[8] M. Boban, M. Giordani, and M. Zorzi, "Predictive Quality of Service (PQoS): The Next Frontier for Fully Autonomous Systems," *IEEE Network*, vol. 35, no. 6, pp. 104–110, Nov/Dec 2021.

[9] M. Wang, Y. Cui, X. Wang, S. Xiao, and J. Jiang, "Machine learning for networking: Workflow, advances and opportunities," *IEEE Network*, vol. 32, no. 2, pp. 92–99, Mar. 2017.

[10] F. Pase, M. Giordani, and M. Zorzi, "On the Convergence Time of Federated Learning Over Wireless Networks Under Imperfect CSI," in *IEEE International Conference on Communications Workshops (ICC)*, 2020.

[11] 5G-ACIA, "5G for Industrial Internet of Things (IIoT): Capabilities, Features, and Potential," *ZVEI*, Nov. 2021.

[12] A. Slivkins, "Introduction to Multi-Armed Bandits," *Foundations and Trends® in Machine Learning*, vol. 12, 2019.

[13] C.-F. Liu and M. Bennis, "Data-driven predictive scheduling in ultra-reliable low-latency industrial IoT: A generative adversarial network approach," in *IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2020.

[14] H. Halabian, "Distributed resource allocation optimization in 5G virtualized networks," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 3, pp. 627–642, Feb. 2019.

[15] F. Hussain, S. A. Hassan, R. Hussain, and E. Hossain, "Machine learning for resource management in cellular and IoT networks: Potentials, current solutions, and open challenges," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 1251–1275, Jan. 2020.

[16] H. Yang, X. Xie, and M. Kadoch, "Intelligent resource management based on reinforcement learning for ultra-reliable and low-latency iov communication networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 5, pp. 4157–4169, Jan. 2019.

[17] L. Liang, H. Ye, and G. Y. Li, "Spectrum sharing in vehicular networks based on multi-agent reinforcement learning," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 10, pp. 2282–2292, Aug. 2019.

[18] D. Russo and B. Van Roy, "Learning to optimize via posterior sampling," *Mathematics of Operation research*, vol. 39, no. 4, pp. 1221–1243, Nov. 2014.

[19] Y. Liu, Y. Deng, M. Elkashlan, A. Nallanathan, and G. K. Karagiannidis, "Analyzing Grant-Free Access for URLLC Service," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 3, pp. 741–755, Mar. 2021.

[20] T. Lattimore and C. Szepesvári, *Bandit Algorithms*. Cambridge University Press, 2020.

[21] C. Boutilier, C. wei Hsu, B. Kveton, M. Mladenov, C. Szepesvari, and M. Zaheer, "Differentiable Meta-Learning of Bandit Policies," in *Advances in Neural Information Processing Systems*, 2020.

[22] R. Bonnefoi, L. Besson, C. Moy, E. Kaufmann, and J. Palicot, "Multi-Armed Bandit Learning in IoT Networks: Learning Helps Even in Non-stationary Settings," in *Cognitive Radio Oriented Wireless Networks*, 2018, pp. 173–185.

[23] M. Polese, M. Giordani, and M. Zorzi, "3GPP NR: the cellular standard for 5G networks," *5G-ITALY White Book*, 2019.

[24] 3GPP, "Study on channel model for frequencies from 0.5 to 100 GHz (Release 16)," Technical Specification (TS) 38.901, 2019.

[25] S. Kuutti, R. Bowden, Y. Jin, P. Barber, and S. Fallah, "A survey of deep learning applications to autonomous vehicle control," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 2, pp. 712–733, Jan. 2020.