

Communication-Efficient Orchestrations for URLLC Service via Hierarchical Reinforcement Learning

Wei Shi^{*†}, Milad Ganjalizadeh^{*†}, Hossein Shokri Ghadikolaei^{*}, and Marina Petrova^{†‡}

^{*}Ericsson Research, Sweden

[†]School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden

[‡]Mobile Communications and Computing, RWTH Aachen University, Germany

Email: {wei.b.shi, milad.ganjalizadeh, hossein.shokri.ghadikolaei}@ericsson.com, petrovam@kth.se

Abstract—Ultra-reliable low latency communications (URLLC) service is envisioned to enable use cases with strict reliability and latency requirements in 5G. One approach for enabling URLLC services is to leverage Reinforcement Learning (RL) to efficiently allocate wireless resources. However, with conventional RL methods, the decision variables (though being deployed at various network layers) are typically optimized in the same control loop, leading to significant practical limitations on the control loop's delay as well as excessive signaling and energy consumption. In this paper, we propose a multi-agent Hierarchical RL (HRL) framework that enables the implementation of multi-level policies with different control loop timescales. Agents with faster control loops are deployed closer to the base station, while the ones with slower control loops are at the edge or closer to the core network providing high-level guidelines for low-level actions. On a use case from the prior art, with our HRL framework, we optimized the maximum number of retransmissions and transmission power of industrial devices. Our extensive simulation results on the factory automation scenario show that the HRL framework achieves better performance as the baseline single-agent RL method, with significantly less overhead of signal transmissions and delay compared to the one-agent RL methods.

Index Terms—6G, availability, factory automation, hierarchical reinforcement learning (HRL), reliability, URLLC.

I. INTRODUCTION

Nowadays, the development of fifth generation of mobile communication systems (5G) technology has achieved a mature technical standard aiming to provide wireless communication services to multiple vertical industrial areas [1]. According to 3rd Generation Partnership Project (3GPP) [2], ultra-reliable low-latency communications (URLLC) stands as one of the three main services for 5G standards, and beyond the standardization, it has shown significant improvements in the efficiency and performance of communication systems [3]–[5]. The main requirements for URLLC (especially in the context of cyber-physical systems (CPSs)) are high reliability (e.g., 10 years without failure), high availability (e.g., 99.9999%), and low latency (below some tens of milliseconds). Machine learning (ML) has proven effective in meeting these stringent

requirements over resource-limited and faulty wireless channels [6]–[9].

A. Literature Review

Various ML-based optimization schemes have been proposed for URLLC. For example, reference [6] proposes a distributed risk-sensitive ML solution for hybrid orthogonal/non-orthogonal radio resource slicing, regulating the spectrum to satisfy the URLLC requirements. Reference [7] implements an reinforcement learning (RL) framework with the deep deterministic policy gradient algorithm [8] into an orthogonal frequency-division multiple access system, minimizing the transmission power. Reference [9] optimizes both power and hybrid automatic repeat request (HARQ) retransmission scheme, leading to further improvements in terms of reliability and availability in factory automation use cases. However, adopting these strategies in real-life applications could be impractical:

- Conventional (flat) RL methods only have one action vector that forces all decision variables to be designed at the same control loop (e.g., resource blocks and power allocation in [7] and power and HARQ retransmission scheme in [9]). However, 5G services usually require various decision variables to be tuned on different control loops. For example, we prefer to have one fixed slicing decision for many coherence intervals while we can constantly adjust the transmit power at every coherence interval.
- To make a decision in a flat RL, we need to collect the information required to determine the state including the ones corresponding to variables with a much slower control loop (slicing in our example). The interval of such data collection is determined by the faster control loop. Such unnecessary data collection results in more energy consumption, network congestion, and extra latency, which could be detrimental to the sustainable operation of URLLC service.

To address these limitations, we develop a hierarchical reinforcement learning (HRL) framework to optimize the operation of URLLC.

B. Hierarchical Reinforcement Learning

Different from one-level RL methods, HRL decomposes one RL problem into a hierarchy of sub-problems, which

This work was partly supported by Swedish Foundation for Strategic Research (SSF) under Grant iPhD-ID17-0079.

©2023 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

allows optimizing different tasks independently with different algorithms, timescales, models, and multiple agents [10]. This usually brings a reduction in exploration complexity, computation, signaling, and time required for the training and inference processes. Reference [11] proposes a hierarchical deep actor-critic method for the resource allocation problems of the 6G massive Internet of Things scenarios. Reference [12] introduces a hierarchical deep Q-networks model with one main controller and multiple sub-controllers to partition a dynamic spectrum access problem into separate sub-problems, reducing the complexity of band selection. In [13], the authors deploy a Meta-HRL framework for resource allocation in vehicular networks to enable faster learning on newly discovered sub-tasks.

C. Our Contributions

In this paper, we propose novel methods to orchestrate parameters of URLLC services. Reference [14] developed a single-agent RL method on a factory automation scenario to jointly optimize the downlink (DL) transmission power and HARQ retransmission control for the communication service performance (availability and reliability) in URLLC services. Here, we address the same problem with a novel HRL framework that supports better performance but with a more flexible structure that enables the system to allocate multiple agents and execute different operations with multi-level policies in different timescales. Our solution substantially reduces the signaling requirements for training and inference. In particular, two of the agents are located at the gNodeBs (gNBs), which significantly reduces the data exchange between gNBs and the centralized remote HRL agent. This efficiency results in time and energy savings in decision-making, thus simplifying the adaptation of our framework to the complex network requirements of real-world wireless systems and 6G.

Notations: Normal font x and X , bold font \mathbf{x} , and uppercase calligraphic font \mathcal{X} denote scalars, vectors, and sets respectively. Besides, $[X]$ denotes the set $\{1, 2, \dots, X\}$, and $|X|$ is the cardinality of set X .

II. SYSTEM MODEL AND PERFORMANCE METRICS

A. System Model

We consider a factory automation scenario where a set of $\mathcal{B} := [B]$ gNBs are present, each serving a set of $\mathcal{U}_b := [U_b]$ industrial devices, where $\mathcal{U}_b \subset \mathcal{U}$, and $\mathcal{U} := [U]$ is the set of all devices. These devices are responsible for executing various functions that facilitate automated production. In this scenario, the communication system must be capable of delivering monitoring data to gNBs and computed or emergency control commands to the actuators in a timely and reliable manner. For the channel model, we assume indoor factory with dense clutter and high base station height (InF-DH) from 3GPP in [15]. Nevertheless, our problem formulation and approach, described in Section III and Section IV, are not limited to this channel model. To enable URLLC efficiently, we consider the orchestration of a set of reliability enhancement features, such as the transmission power to industrial devices and the

maximum number of diversity transmissions (i.e., transmitting multiple instances of a packet or its segments in space, time, and/or frequency). In this paper, we focus on the orchestration of transmission power and HARQ retransmissions in DL direction. However, our framework can easily be extended for other reliability enhancement features and uplink transmissions.

From the network management perspective, in the hierarchical multi-tier architecture of cellular networks, we assume two levels of control for global and local optimizations. Although there is only one top-level controller, we assume a set of low-level controllers co-located with the gNBs, managing the transmissions towards \mathcal{U}_b together.

B. From Network to Communication Service Performance

The key element in characterizing service performance is defining service failures accurately. In [16], 3GPP defines survival time denoted as T_s , as the duration for which an application can continue to function without receiving an expected packet. Therefore, a communication service failure occurs if no packets have been received by the reception entity for the duration of survival time.

We can define the network layer state variable $Y_{b,u}(t)$ for industrial device u associated to gNB b at time t , where $Y_{b,u}(t)$ is considered 0 if the last packet fails to reach the communication interface within a specified delay bound due to decoding issues at lower layers, excessive retransmissions, or queuing delays, and 1 otherwise. We can define the network state variable $Y_{b,u}(t)$ for the u th industrial device at time t , where $Y_{b,u}(t)$ is considered 0 if the last packet fails to reach the communication interface within a specified delay bound due to decoding issues at lower layers, excessive retransmissions, or queuing delays, and 1 if the packet is received successfully and timely. Since sporadic packet losses¹ within T_s do not impact the end-to-end service performance, the application layer state variable can be defined as [14]:

$$Z_{b,u}(t) := \max_{t-T_s \leq \tau \leq t} Y_{b,u}(\tau). \quad (1)$$

The application state variable $Z_{b,u}$ enables us to define and formulate our two reliability key performance indicators (KPIs), communication service availability and communication service reliability.

Communication Service Availability: It refers to the ability of an end-to-end communication service to perform its intended function without failure at a given point in time and is commonly expressed as a probability or as a percentage of time that the system is operational within a specified time period [2]. Considering the failure definition and $Z_{b,u}(t)$ in (1), the communication service availability, $\alpha_{b,u}$, is [14]

$$\alpha_{b,u} := \lim_{t \rightarrow \infty} \Pr(Z_{b,u}(t) = 1) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T Z_{b,u}(t) dt. \quad (2)$$

¹Here, packet loss refers to all packets that fail to reach their intended recipient within their delay bound.

However, the communication service availability for u th device can be approximated in a short time, Δt_k , via [14]

$$\bar{\alpha}_{b,u}(\Delta t_k) := \frac{1}{\Delta t_k} \int_{t_k - \Delta t_k}^{t_k} Z_{b,u}(t) dt. \quad (3)$$

Communication Service Reliability: It refers to the ability of an end-to-end communication service to operate without failures over a specific period, given certain environmental and operational conditions [16]. It can be expressed as the meantime that the service is operational, that is $Z_{b,u}(t) = 1$. Therefore, reliability $\rho_{b,u}$ is formulated as [14]

$$\rho_{b,u} := \lim_{T \rightarrow \infty} \frac{1}{F_{b,u}(T)} \int_0^T Z_{b,u}(t) dt, \quad (4)$$

where $F_{b,u}(T)$ denotes the number of crossings from $Z_{b,u}(t) = 1$ to $Z_{b,u}(t) = 0$ within $[0, T]$. Since communication service reliability's unit is time, we can alternatively approximate it via the crossing rate, $\psi_{b,u}$, representing the crossings of $Z_{b,u}(t)$ from 1 to 0 during Δt_k ; defined as $\psi_{b,u} := \lim_{T \rightarrow \infty} F_{b,u}(T)/T$. Note that $\psi_{b,u}$ is inversely proportional to $\rho_{b,u}$ in (4). Then, the crossing rate can be approximated by [17]

$$\bar{\psi}_{b,u}(\Delta t_k) := \frac{F_{b,u}(\Delta t_k)}{\Delta t_k}. \quad (5)$$

In the following section, we introduce our HRL solution based on the formulation of the KPIs.

III. OPTIMIZATION WITH HRL FRAMEWORK

In this paper, the objective is to maximize the communication service availability, in (2), and communication service reliability, in (4), of a CPS by optimizing the configuration of transmission power levels and the number of retransmissions. Nevertheless, our framework can easily be extended for more decision variables. We propose to solve this problem using a HRL framework, where a high-level agent collaborates with low-level agents to manage the bi-level control of the communication system. The high-level agent is responsible for inter-agent coordination and, therefore, we assign the task of mitigating inter-cell interference globally to it by placing the transmission power under its control. Hence, we model the problem as a twin timescale Markov decision process and then apply the soft actor-critic (SAC) algorithm, as presented in [18], to solve it. Additionally, we use the branching technique described in [14] to enhance the performance of the algorithm. For simplicity, we assume that the timescales of the low-level agents are identical and denote it with Δt_k , for an iteration starting from t_k , and represent the high-level agent's timescale with $\Delta t_k^h := t_k - t_{k-1}$, where $\Delta t_k^h = c\Delta t_k$, $\forall c, k, \kappa \in \mathbb{N}$.

A. State Space

The state represents the set of various measurements from the environment that affects the performance of our main KPIs, namely communication service availability and reliability. The state of u th device associated with the b th low-level agent, $s_{b,u}(\Delta t_k)$, is measured within $[t_k - \Delta t_k, t_k]$, and consists of

various factors that can be classified as direct and indirect, based on their effects on the two KPIs. As mentioned, individual availability $\alpha_{b,u}$ and reliability $\rho_{b,u}$ are the functions of the probability of packet loss and average operational time. Therefore, we add packet loss rate and mean downtime of the network layer as the direct factors to the state space. Apart from these two variables, we also consider the various factors that are not included in the KPI functions but importantly affect the communication quality indirectly, as signal to noise and interference ratio (SINR), packet transmission delays, the status of the radio link control (RLC) layer buffers, path gain, the number of HARQ transmissions, and the number of used resource blocks. To enable a more concrete description of the communication environment, we include mean, median values, 95th, and 5th percentile of SINR, path gain, and RLC buffer status. However, for the rest of the factors, we incorporate the mean of the samples in the state. Thus, the state of b th low-level agent for its k th iteration is defined as $s_b(\Delta t_k) := \{s_{b,u}(\Delta t_k) | \forall u \in \mathcal{U}_b\}$.

As for the high-level agent, we define the global state consisting of all elements in $s_{b,u}$, but measured in a much larger time scale, Δt_k^h , and including all devices. Then, the high-level agent's state for κ th iteration can be defined as $s(\Delta t_k^h) := \{s_{b,u}(\Delta t_k^h) | \forall u \in \mathcal{U}_b, \forall b \in \mathcal{B}\}$.

B. Action and HRL Policy

The action space consists of a series of decision parameters for the agents to interact with environments. Similar to reference [14], we consider the quantized transmission power levels and the number of retransmissions in action set for each device. However, we decompose the joint action into a multi-level policy that enables the different HRL agents to learn one or a combination of different network functions with different timescales based on network requirements. Specifically for our scenario, we set up a two-level policy that includes a global optimization of transmission powers (via a high-level agent), and local optimization of the maximum number of retransmissions (via several low-level agents). Hence, the high-level action in κ th iteration is defined as $\mathbf{a}_\kappa^h := (p_\kappa^1, \dots, p_\kappa^u, \dots, p_\kappa^U)$, where $p_\kappa^u \in \{p_{\min}, p_1, p_2, \dots, p_{\max}\}$, and the action for b th low-level agent in k th iteration, $\mathbf{a}_{b,k}$, is a vector where each element represents the configured maximum numbers of retransmissions for devices in \mathcal{U}_b . Compared to flat (or single-level) RL, HRL decomposes the action space into layers, resulting in higher scalability for handling complex orchestrations in cellular systems.

C. Reward Functions

The RL agents follow an explicit objective to maximize the sum of discounted rewards. Considering the estimations on communication service availability and crossing rate in (3) and (5), respectively, we introduce two different reward functions. The first one, inspired by [14], targets the maximization of the average of reliability KPIs, and is defined as

$$r_b(\Delta t_k) := \frac{1}{\omega U_b} \sum_{u=1}^{U_b} \left(\omega \bar{\alpha}_{b,u}(\Delta t_k) - (1-\omega) \bar{\psi}_{b,u}(\Delta t_k) \right), \quad (6)$$

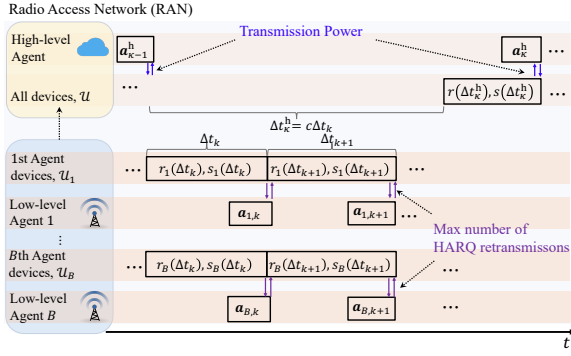


Fig. 1: Learning procedure of the two-level policy HRL framework

where $0 < \omega < 1$ decides the importance of $\bar{\alpha}_{b,u}(\Delta t_k)$ and $\bar{\psi}_{b,u}(\Delta t)$, and its placement in the denominator helps to bound the reward function by 1. Consequently, the high-level reward is defined as $r(\Delta t_k^h) := \frac{1}{B} \sum_{b=1}^B r_b(\Delta t_k^h)$, where $r_b(\cdot)$ is derived as in (6), but within $(t_{k-1}, t_k]$. For the second function, we adopt a risk-sensitive reward that replaces the average value of users with the extremum within \mathcal{U}_b , where agents can be more aggressive in exploration. Enabling agent m to maximize the availability and minimize the crossing rate, the k th iteration reward is defined as

$$r_b(\Delta t_k) := \exp\left(\frac{\eta}{\omega} (r'_b(\Delta t_k) - \omega)\right), \quad (7)$$

where

$$r'_b(\Delta t_k) := \omega \min_{u \in \mathcal{U}_b} (\bar{\alpha}_{b,u}(\Delta t_k)) - (1-\omega) \max_{u \in \mathcal{U}_b} (\bar{\psi}_{b,u}(\Delta t_k)), \quad (8)$$

and η represents a fixed coefficient predefined to adjust the reward reduction with the application's sensitivity to reliability KPIs. The high-level reward, $r(\Delta t_k^h)$, can then be calculated as (7), but within $(t_{k-1}, t_k]$ where the min and max functions in (8) are calculated for $\forall u \in \mathcal{U}$.

D. Learning Procedure

The learning procedure of our two-level policy HRL framework is as shown in Fig.1, where two network operations, global power control and maximum HARQ retransmission number are decided by a high-level agent and low-level agents, parallelly following two timescales (Δt_k , Δt_k^h). At every time step Δt_k , the low-level agents collect the states of the users they control, also the calculated rewards. Then the generated actions for HARQ retransmission number are sent to the communication system. The high-level agent issues power values to all the devices every $c\Delta t_k$. Similar to the design in [9], all the agents are model-free and can individually implement any RL algorithms according to the task requirements. For the sake of fair comparison to the single-agent RL in [14], we deploy the same SAC algorithm [18] to all the HRL agents, and adopt the same branching technique that enables the continuous actions to describe our discrete actions in our factory automation scenario.

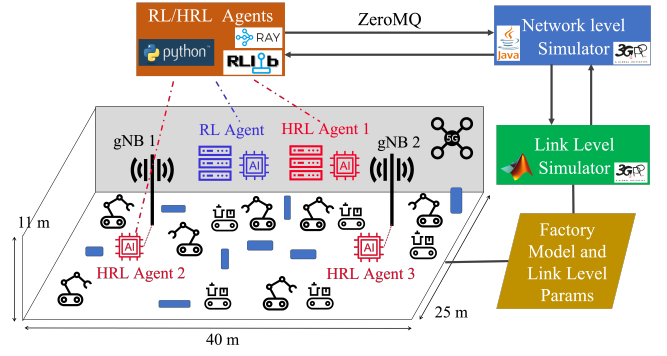


Fig. 2: Framework of the simulators with scenario setup

IV. SIMULATION METHODOLOGY AND RESULT ANALYSIS

In this section, we present the simulation methodology and evaluate the performance of our HRL framework.

A. Simulation Configuration and Methodology

Fig. 2 shows our simulation architecture, which consists of a link-level Matlab simulator, a network-level Java simulator, and the RL agents implemented using an open source library, RLlib. With path gain matrices, 3D channel data [15], and nodes allocation provided by the link-level simulator, the network-level simulator is able to simulate the network with multiple gNBs and users in physical, medium access control, and other higher layers of the 5G-NR network. As for the RL/HRL deployment, our network-level simulator supports interacting with external agents using the pipelines based on ZeroMQ protocol.

We considered a $40 \times 25 \times 11 \text{ m}^3$ factory with two 10m height gNBs providing communication service to 10 industrial devices. Moreover, we assumed that the devices are in a high interference condition, and they move with the speed of 30 km/h, while staying in close proximity to the original position. We considered periodic control traffic with periodicity 2 ms, and a delay bound 2.5 ms. On the transmitter side, packets are queued in the RLC buffers and then sent via transport blocks based on the selected modulation and coding scheme. On the receiver side, the successfully decoded packets that were received after the delay bound were discarded by packet data convergence protocol. In the end, the reliability and availability are calculated in the application layer based on survival time, T_s . The network and learning parameters of our simulations are presented in TABLE I. For more details on our simulation setup, you can refer to [19, §VI.A, §VI.B].

The goal of our HRL solution is to find the optimal solution for power control and HARQ retransmissions to maximize communication service availability and reliability. We considered the following baselines in our evaluations:

- **MaxRetPwr**: Similar to [14], all resource blocks are configured with 0.02 W, and the maximum number of transmissions is set to 2.
- **RLAvg/RLRiskSen**: There exist only one RL agent (in a remote server) interacting with the two gNBs (the blue agent in Fig.2). The step period is set to 0.1 s. The

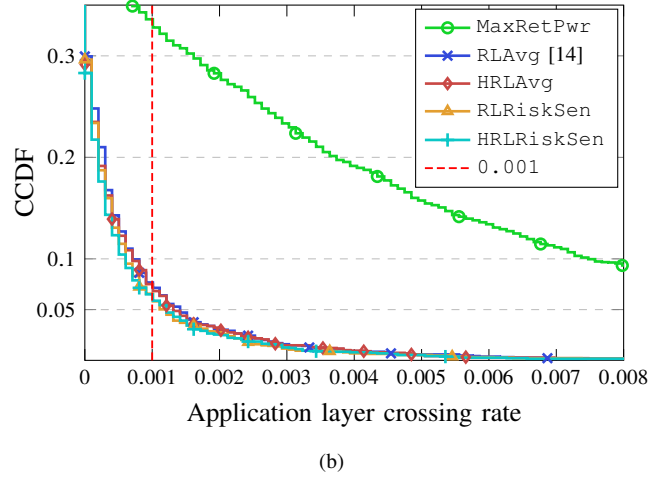
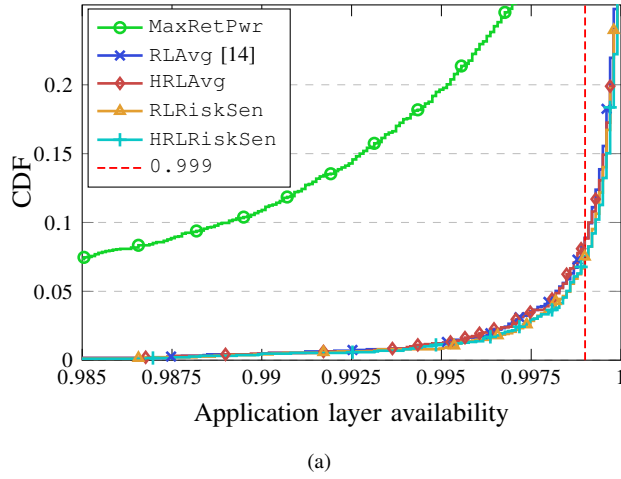


Fig. 3: Comparison of single-agent RL solution and HRL framework in terms of (a) communication service availability, and (b) crossing rate of user equipments, sampled from all the devices from the 10 s simulations.

TABLE I: Simulation Parameters.

Network Parameters	
Parameter	Value
Deployment	2 gNBs, 1 cell each
gNBs antenna height	8 m
Devices' height	1.5 m
Carrier frequency	2.6 GHz
Bandwidth	20 MHz
TTI length/Subcarrier spacing	0.5 ms/30 KHz
DL transmit power (p_{min}/p_{max})	0.2 W/0.5 W
Number of gNB/Devices' antennas	2/2
DL URLLC delay bound	2.5 ms
DL URLLC Survival time (T_s)	5 ms
Simulation time	10 s/Episode
Learning Parameters	
Parameter	Value
Neural network hidden layers	128×128
Activation function	ReLU
Loss function	MSE
Optimizer	mini-batch SGD
Discount factor	0.1
Batch size	200
Learning rate	0.0003
RL step period (Δt)	0.1 s
HRL step period low/high-level ($\Delta t_k / \Delta t_k^h$)	0.1/0.5 s, $\forall k, \kappa \in \mathbb{N}$

agent selects from the two possible levels of 0.008 W, and 0.02 W, set for all resource blocks allocated to a specific device, and the maximum number of transmissions (which can be either 1 or 2). While RLavg used the average reward from [14] and presented in (6), RLRiskSen was implemented with the risk-sensitive reward in (7).

Furthermore, for the HRL solution, we implemented one top-level agent and two low-level agents (the red agents in Fig. 2). The former agent's objective was to optimize the DL transmission powers globally, while the numbers of HARQ retransmissions were optimized per gNB by the latter agents. For the sake of fair comparison, the action space was set similar to RLavg and RLRiskSen. Besides, we considered two setups as HRLavg, where we incorporated the average reward in (6), and HRLRiskSen, where we incorporated the risk-sensitive reward, in (7).

B. Result and Analysis

Fig. 3a and Fig. 3b present the cumulative distribution function (CDF) of the communication service availability and the complementary CDF (CCDF) of the device crossing rate, determined by (3) and (5), respectively. In these figures, each data point represents the device availability (in Fig. 3a) or crossing rate (in Fig. 3b) in one simulation round. Assuming an availability requirement of 0.999, Fig. 3a shows that both HRLRiskSen and RLRiskSen achieve a similar violation probability of 0.083, while MaxRetPwr can only reach violation probability of 0.35. Similarly, assuming a crossing rate requirement of 0.001, Fig. 3b indicates that HRLRiskSen and RLRiskSen can obtain violation probability of 0.058, closely followed by HRLavg, and they all significantly outperform MaxRetPwr. Furthermore, our HRL framework shows a lower violation probability for availability > 0.999 and crossing rate < 0.001 , outperforming the RL.

It is surprising that our HRL framework achieves better performance than the single-agent RL method, which learns the states from all the devices within one agent. In comparison, the HRL agents have three different actions, where two of them only learn partial states of the devices from one gNB. Since we run both RL and HRL training for a fixed number of iterations, such improvement can be contributed by the reduction in the action space (as a result of action decomposition), leading to faster convergence.

The error bar plot in Fig. 4 demonstrates the mean (shown by square) and 5th percentile (shown by line) device availability. The average and 5th percentile availability of single-agent RL and HRL in both rewards achieve similar performance that is much higher than MaxRetPwr. Although the baseline simulation adopts the maximum power and 2 HARQ retransmissions, its poor performance indicates that the lack of freedom in operations could reduce interference management capability. Same as in Fig. 3, the risk-sensitive reward shows improvement in performance than the average reward for the HRL and RL.

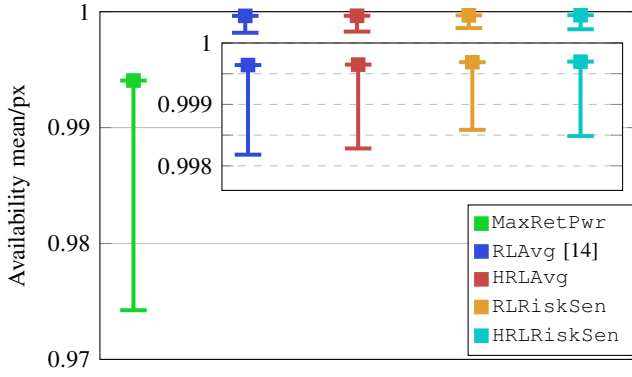


Fig. 4: Mean (square) and 5th percentile (line) availability for the simulations

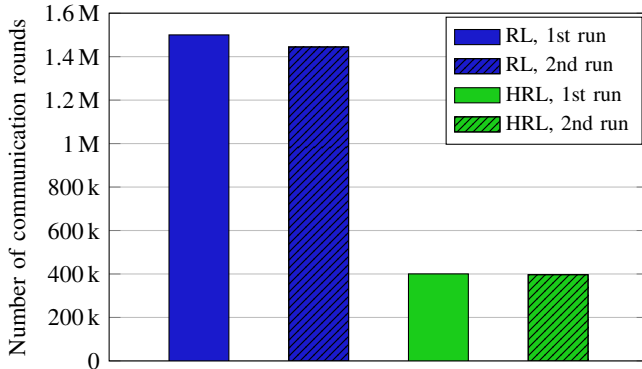


Fig. 5: Number of signal exchanges in the training process.

Fig. 5 presents the number of signal exchanges (communication rounds) the learning frameworks require to converge (i.e., all reward values tend to stabilize and converge, and they no longer increase with further exploration). One signal exchange in Fig. 5 represents a signal transmission between the remote agent and one of the gNBs. According to the allocation of agents in Fig. 2, there are two learning agents assembled locally with the gNBs, the transmission overhead of which can be ignored compared to the remote server. In the case of the single-agent RL solution, the agent transmits one message of action to each gNB and receives a reply message with states and rewards from them. Therefore, there are four signal transmissions at every step. Similarly, for HRL framework, there are four signal transmissions at every high-level step. We performed two RL simulations and two HRL simulations and logged the total number of exchanged signals till convergence. As Fig. 5 confirms, compared to HRL, RL simulations required over triple the number of communication rounds to converge. Such reduction in signal exchanges can significantly improve the energy efficiency of large-scale communication systems with many remote devices and simultaneous operations. Beyond signal exchanges, decomposing the RL action allows the HRL's average learning time per iteration to be 33% less than that of the single-agent RL (i.e., 6.288 ms vs 8.386 ms). This translates into massive gains in terms of latency and energy saving to run the training, especially in dynamic environments where we may need to retrain the models regularly.

V. CONCLUSIONS

In this paper, we propose a HRL framework and implement that into a simulated factory automation model to optimize the operations of power control and HARQ retransmissions in 5G communication, aiming to achieve the optimal availability and reliability according to the standard of URLLC. We design and compare five simulations that separately deploy the single-agent RL strategy, our HRL framework with average and risk-sensitive reward functions, and one with fixed operations as the baseline. Our HRL framework achieves better performance on availability and reliability to the ideal single-agent RL solution and significantly outperforms the baseline simulation. Besides, our HRL framework enables better flexibility that allows the operations to be executed in different timescales. Furthermore, due to the flexible allocation of agents, the HRL solution can save signal consumption by at least triple less than that of the RL in the scenario of our factory model.

REFERENCES

- [1] H. Basilier *et al.*, "Applied network slicing scenarios in 5G," *Ericsson Technol. Rev.*, vol. 2021, no. 2, pp. 2–11, 2021.
- [2] *Service requirements for the 5G system*, 3GPP, TS 22.261 v19.2.0, 2023.
- [3] Z. Li *et al.*, "5G URLLC: Design challenges and system concepts," in *IEEE Int. Symp. Wireless Commun. Syst. (ISWCS)*, 2018.
- [4] P. Popovski *et al.*, "Wireless access in ultra-reliable low-latency communication (URLLC)," *IEEE Trans. Commun.*, vol. 67, no. 8, pp. 5783–5801, 2019.
- [5] H. Ren *et al.*, "Joint pilot and payload power allocation for massive-MIMO-enabled URLLC IIoT networks," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 5, pp. 816–830, 2020.
- [6] A. Azari, M. Ozger, and C. Cavdar, "Risk-aware resource allocation for urllc: Challenges and strategies with machine learning," *IEEE Commun. Mag.*, vol. 57, no. 3, pp. 42–48, 2019.
- [7] A. T. Z. Kasgari *et al.*, "Experienced deep reinforcement learning with generative adversarial networks (GANs) for model-free ultra reliable low latency communication," *IEEE Trans. Commun.*, vol. 69, no. 2, pp. 884–899, 2021.
- [8] T. P. Lillicrap *et al.*, "Continuous control with deep reinforcement learning," 2015, *arXiv:1509.02971* [cs.LG].
- [9] M. Ganjalizadeh *et al.*, "An RL-based joint diversity and power control optimization for reliable factory automation," in *IEEE Global Commun. Conf. (GLOBECOM)*, 2021.
- [10] S. Pateria *et al.*, "Hierarchical reinforcement learning: A comprehensive survey," *ACM Comput. Surveys*, vol. 54, no. 5, pp. 1–35, 2021.
- [11] H. A. Shah, L. Zhao, and I.-M. Kim, "Joint network control and resource allocation for space-terrestrial integrated network through hierarchal deep actor-critic reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 70, no. 5, pp. 4943–4954, 2021.
- [12] S. Liu, J. Wu, and J. He, "Dynamic multichannel sensing in cognitive radio: Hierarchical reinforcement learning," *IEEE Access*, vol. 9, pp. 25 473–25 481, 2021.
- [13] Y. He *et al.*, "Meta-hierarchical reinforcement learning (MHRL)-based dynamic resource allocation for dynamic vehicular networks," *IEEE Trans. Veh. Technol.*, vol. 71, no. 4, pp. 3495–3506, 2022.
- [14] M. Ganjalizadeh *et al.*, "Saving energy and spectrum in enabling URLLC services: A scalable RL solution," *IEEE Trans. Ind. Informat.*, early access, 2023. doi: 10.1109/TII.2023.3240592.
- [15] *Study on channel model for frequencies from 0.5 to 100 GHz*, 3GPP, TR 38.901 v17.0.0, 2022.
- [16] *Service requirements for cyber-physical control applications in vertical domains*, 3GPP, TS 22.104 v19.0.0, 2021.
- [17] A. Hoyland and M. Rausand, *System reliability theory: models and statistical methods*. John Wiley & Sons, 2009.
- [18] T. Haarnoja *et al.*, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 1861–1870.
- [19] M. Ganjalizadeh *et al.*, "Device selection for the coexistence of URLLC and distributed learning services," 2022, *arXiv:2212.11805* [cs.NI].