

A Reinforcement Learning Framework for PQoS in a Teleoperated Driving Scenario

Federico Mason*, Matteo Drago*, Tommaso Zugno^{†*}, Marco Giordani*, Mate Boban[†], Michele Zorzi*

*Department of Information Engineering, University of Padova, Italy. Email: {name.surname}@dei.unipd.it

[†]Huawei Technologies, Munich Research Center, Germany. Email: name.surname@huawei.com

Abstract—In recent years, autonomous networks have been designed with Predictive Quality of Service (PQoS) in mind, as a means for applications operating in the industrial and/or automotive sectors to predict unanticipated Quality of Service (QoS) changes and react accordingly. In this context, Reinforcement Learning (RL) has come out as a promising approach to perform accurate predictions, and optimize the efficiency and adaptability of wireless networks. Along these lines, in this paper we propose the design of a new entity, integrated at the RAN level that implements PQoS functionalities with the support of an RL framework. Specifically, we focus on the design of the reward function of the learning agent, able to convert QoS estimates into appropriate countermeasures if QoS requirements are not satisfied. We demonstrate via ns-3 simulations that our approach achieves better results in terms of QoS and Quality of Experience (QoE) performance of end users in a teleoperated driving scenario.

Index Terms—Predictive Quality of Service (PQoS), teleoperated driving, reinforcement learning (RL), RAN, ns-3.

I. INTRODUCTION

We are witnessing a rapid evolution towards autonomous systems, able to safely operate without human intervention [1]. In this scenario, the dynamic nature of the environment in which autonomous machines are deployed may result in the Quality of Service (QoS) to change and degrade unexpectedly, with potentially catastrophic consequences if communication failures are not promptly notified. To solve this issue, the research community has been investigating Predictive Quality of Service (PQoS) [2] as a means to provide autonomous systems with advance notifications about QoS changes. Unlike *reactive* mechanisms, which respond to QoS deterioration only after it occurs, a *predictive* approach gives time to the applications to foresee unanticipated events, and adapt their operations accordingly. In the automotive sector, Vehicle-To-Everything (V2X) is associated to strict QoS requirements in terms of ultralow latency, as well as ultrahigh throughput, reliability, and security [3]. Hence, PQoS may assist V2X networks to preemptively evaluate the potential risks of QoS degradation, and guarantee reliable driving.

At first, network prediction was developed based on linear regression and/or filter-based models (e.g., Kalman filters) [4], [5]. However, linear regression came out as a good estimator when the underlying relationship between the system's variables and the response was known to be linear, even though most real systems are non-linear. At the same time, filter-based prediction mechanisms have been shown to be very sensitive to

imperfections of the environment, as well as random dynamics, and mainly focus on real-time or short-time estimation of the target state. In turn, autonomous applications may also require predictions of the system's future behavior [6], which may span several minutes or even hours [7]. Moreover, despite their quick convergence in dynamic processes, these methods tend to require a priori knowledge of the process to learn, which is not always available in autonomous systems [8]. More recently, machine learning (ML) technologies have emerged as a powerful approach to predict and optimize wireless networks [9], without requiring explicitly pre-programmed a priori rules, which makes these techniques particularly interesting for enabling PQoS.

Based on the above introduction, this paper addresses the problem of developing and testing PQoS mechanisms for V2X networks. To this aim, we propose the design of a new entity called "RAN-AI" (Sec. II), connected to the different components of the Radio Access Network (RAN) that: (i) collects data from different sources; (ii) makes QoS predictions based on the acquired data, recognizes and notifies upcoming QoS variations; and (iii) defines and applies network countermeasures in case QoS requirements are not satisfied. The RAN-AI entity integrates a Reinforcement Learning (RL) framework, able to identify the optimal (set of) countermeasures for PQoS (Sec. III). In particular, we focus on the design of the learning algorithm and the reward function of the learning agent. The performance of the RL model for PQoS is validated via ns-3 simulations in a teleoperated driving scenario, considering realistic setup and input parameters (Sec. IV). We demonstrate that our RL framework guarantees the optimal trade-off between QoS and Quality of Experience (QoE) for end users (and satisfies QoS requirements of V2X applications) compared to other baseline solutions. Finally, we formalize our conclusions and next research steps (Sec. V).

II. OUR PQoS FRAMEWORK

We consider a teleoperated driving scenario, in which connected vehicles are controlled by a remote driver (either human or software). In view of the highly dynamic nature of the V2X environment, PQoS shall be able to predict QoS degradation and gracefully change the operational mode of the system to ensure that end users satisfy their strict latency/reliability constraints. For example, PQoS may assist the control center in adapting the vehicles' speed and trajectory based on the condition of the radio environment [2]. PQoS procedures

can be implemented at both the Core Network (CN) and the RAN. In this work we focus on the latter, and describe the functionalities of a new entity called “RAN-AI” that we designed to facilitate PQoS in V2X systems.

A. RAN-AI Entity for PQoS

RAN operations in V2X, from modulation and coding to resource allocation and scheduling, can affect QoS performance, if not properly configured. Along these lines, we propose the design of a RAN-AI entity, installed in a remote/edge server, or at the gNB, that executes methods to improve the network performance if the agreed QoS is not satisfied. In our implementation, the RAN-AI entity is connected to the different components of the RAN (Centralized Unit (CU), Distributed Units (DUs) and Remote Radio Units (RRUs)), as well as to the core network by means of dedicated interfaces, and hosts an RL agent for PQoS. Notably, the RAN-AI entity is in charge of the following tasks:

- (i) Get application statistics, and collect full-stack RAN measurements from the gNB and the end users, that may be used as input parameters of the RL agent (Sec. II-B).
- (ii) Share such measurements with the RL agent, which determines the optimal set of countermeasures to maximize the overall performance (Sec. II-C).
- (iii) Inform the end users about QoS changes, and suggest the countermeasure(s) to adopt, based on the agent’s decision(s), so that they can adjust their behavior accordingly.

B. PQoS Inputs

For the RAN-AI entity to operate properly, the availability of measurements (or “inputs”) from diverse sources is the key. The RAN-AI may have access to the following data from both the environment and the RAN itself.

- *Context information.* This incorporates information about (i) the operational scenario, (ii) road elements, (iii) the network deployment, and (iv) the time of the day.
- *Trajectory of the users.* Driving applications may provide to the network the planned trajectory of end users. Moreover, the RAN-AI entity shall acquire data (based on trajectory statistics) about the gNB at which end users are/will be attached in future locations, as an indication of the cell load.
- *Traffic information.* The RAN-AI entity may gather traffic conditions from external control centers, which may also provide traffic predictions based on historical data.
- *Network metrics.* The RAN-AI entity can get access to measurements gathered at the Physical (PHY) and Medium Access Control (MAC) layers (e.g., L1 measurements such as Reference Signal Received Power (RSRP), Reference Signal Received Quality (RSRQ), Signal to Interference plus Noise Ratio (SINR), Physical Resource Blocks (PRBs) utilization, and Modulation and Coding Scheme (MCS) index), Radio Link Control (RLC) and Packet Data Convergence Protocol (PDCP) layers (e.g., statistics of the data traffic exchanged across the users).

- *Higher layers metrics.* The RAN-AI may be informed by the end users’ applications about the experienced end-to-end (E2E) performance (e.g., mean, standard deviation, minimum and maximum value of delay and throughput).

C. PQoS Countermeasures

If QoS requirements are not satisfied, the role of the RAN-AI is to convert PQoS inputs into appropriate network countermeasures (or “actions”). Besides operating directly on the driving patterns, the RAN-AI may undertake lower-layer actions (e.g., changing scheduling decisions, adapting the radio resource allocation as a function of the propagation conditions, modulating traffic requests based on the available network capacity, or modifying the system numerology to provide more resilient communication channels).

Another action that can alleviate the burden on the channel is to reduce the size of the data generated at the application layer before transmission. In our scenario, end users transmit Light Detection and Ranging (LiDAR) data in the form of point clouds, whose size depends on the *application mode*, i.e., the level of compression of the observations [10]. In the context of this work, the RAN-AI entity implements an RL agent whose goal is to identify the optimal application mode (i.e., the action) for the end users when generating the data (and thus the corresponding size of the packets to send). Finally, the RAN-AI communicates to each end user the corresponding best action that the agent has identified. This involves the transmission of a notification packet, which may be exposed to both transmission delays and errors.

III. REINFORCEMENT LEARNING MODEL

Reinforcement Learning (RL) is a powerful paradigm that models the target scenario as a Markov Decision Process (MDP), where time is discretized into slots and, in each slot t , an agent observes the system state, takes a new action, and receives a reward or a penalty accordingly. The goal of the agent is then to determine the optimal policy π^* , i.e., the map between states and actions $\pi : \mathcal{S} \rightarrow \mathcal{A}$ leading to maximizing the cumulative reward $G_t = \sum_{\tau=t}^{\infty} \lambda^{\tau-t} R(\tau)$, where R is the reward and $\lambda \in [0, 1)$ is the future reward discount factor.

In this work, the RAN-AI is installed at the gNB and implements an RL agent which controls all the vehicles connected to that cell. Hence, the agent’s actions correspond to the application modes described in Sec. II-C, and the reward is obtained by evaluating the vehicles’ performance. This approach ensures that the size of the state/action spaces is invariant with respect to the number of users, and improves the learning efficiency since, during the training phase, the system can exploit the data gathered by all the vehicles.

A. A Learning Framework for PQoS

In our system, at each step t , the agent observes the state $s_t \in \mathcal{S}$ of each vehicle, where s_t is a vector containing the RAN-AI entity’s inputs. Given the state s_t , the agent computes the q -value $Q(s_t, a)$ associated with each action $a \in \mathcal{A}$, where $Q(s_t, a)$ is the estimate of the cumulative reward G_t that can be obtained playing action a in state s_t , and then following the

optimal policy. During the training, the agent will eventually converge towards the action with the highest q-value, which ensures the best cumulative reward.

We adopt a Deep Reinforcement Learning (DRL) approach, and approximate the agent’s policy by means of a Neural Network (NN), which makes it possible to handle continuous state spaces and overcome the “curse of dimensionality” phenomenon [11]. In particular, we consider a Feed-Forward NN, with S inputs and A output neurons, and implement the Rectified Linear Unit (ReLU) activation function across the different layers [12]. The input size (S) of the NN coincides with the number of input parameters of the RAN-AI entity, while the output size (A) corresponds to the number of possible agent’s actions, i.e., the different application modes. Hence, our architecture is trained to approximate the function $Q(\cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, that establishes the quality of each state-action pair. In particular, the training of the agent is performed according to the Double Q-learning (DQL) algorithm described in [13], which is an extended version of the classical Q-learning algorithm introduced in [14]. The main details of our NN architecture are given in Table I.

B. Reward Function

The design of the reward function is a particularly critical task in RL. If the reward does not fully represent the system’s requirements, the agent may learn an undesirable behavior, and lean towards suboptimal actions. In our scenario, analyzing the performance of the system is not straightforward, as it depends on two different factors:

- The Quality of Service (QoS): The agent’s decision should ensure that the vehicle satisfies the agreed QoS in terms of different Key Performance Indicators (KPIs), in particular in terms of maximum end-to-end delay δ_m and Packet Reception Ratio (PRR) (among the most representative KPIs for teleoperated driving scenarios [15]).
- The Quality of Experience (QoE): The agent’s decision should ensure that the quality of the transmitted data is good enough to perform teleoperated driving operations (e.g., object detection [16]). This is measured based on the symmetric point-to-point Chamfer Distance CD_{sym} between the transmitted data \hat{D} and original data D acquired by the LiDAR, expressed as [17]

$$CD_{\text{sym}} = \sum_{\forall \mathbf{d} \in D} \min_{\forall \hat{\mathbf{d}} \in \hat{D}} \|\mathbf{d} - \hat{\mathbf{d}}\|_2^2 + \sum_{\forall \hat{\mathbf{d}} \in \hat{D}} \min_{\forall \mathbf{d} \in D} \|\mathbf{d} - \hat{\mathbf{d}}\|_2^2. \quad (1)$$

To make the agent capture both the above aspects, we design the reward as a piece-wise function, which returns 0 whenever the QoS requirements are not met; otherwise, a positive value that depends on both the QoS and the QoE of the end users. Let $\hat{P}RR_t$, $\hat{\delta}_t$, and δ_m be the PRR and average delay of the vehicle at time t , and the maximum delay tolerated by the system (based on the use case of interest), respectively. At time slot t , if the QoS requirements are met, i.e., $\hat{\delta}_t < \delta_m$ and $\hat{P}RR_t = 1$, the agent reward $R(t)$ is given by

$$R(t) = (1 - \alpha) \frac{\delta_m - \hat{\delta}_t}{\delta_m} + \alpha \frac{CD_{\text{sym},m} - \hat{C}D_{\text{sym},t}}{CD_{\text{sym},m}}, \quad (2)$$

TABLE I: Simulation parameters.

Parameter	Description	Value
f_c	Carrier frequency	3.5 GHz
B	Total bandwidth	50 MHz
P_{TX}	Transmission power	23 dBm
T	RAN-AI update periodicity	100 ms
τ_s	Simulation time	80 s
N_u	Number of vehicles	{1, 5}
λ	Discount factor	0.95
ζ	Learning rate	10^{-4}
ϵ	Weight decay	10^{-3}
α	QoS/QoE weight	{0.5, 1}
δ_m	Max. tolerated delay	50 ms
$CD_{\text{sym},m}$	Max. tolerated Chamfer Distance	45
Layer size (inputs \times outputs)		$8 \times 12 \rightarrow 12 \times 6 \rightarrow 6 \times 3$

and is equal to 0 otherwise. Specifically, in (2), α is a positive value in $[0, 1]$, while $CD_{\text{sym},t}$ and $CD_{\text{sym},m}$ are the Chamfer Distance measured at time t and the maximum Chamfer Distance that can be tolerated, respectively. The balancing between QoS and QoE is determined by α , which is a tuning parameter to be set according to the target scenario.

IV. PERFORMANCE EVALUATION

In Sec. IV-A we describe our simulation setup, while in Sec. IV-B we validate through numerical simulations the performance of our RAN-AI implementation for PQoS, compared to other baseline methods.

A. Simulation Setup and Parameters

Our results are based on ns-3 simulations, thereby enabling full-stack end-to-end analyses. To do so, we extended the ns-3-mmwave module [18] to incorporate a new RAN-AI class and its functionalities. Simulation parameters are reported below and summarized in Table I.

a) *Scenario*: Our PQoS framework was validated in a test scenario with one gNB covering a portion of the city of Bologna, Italy, and $N_u = \{1, 5\}$ vehicles moving according to realistic mobility traces generated using Simulation of Urban MObility (SUMO) [19]. The system is operating at 3.5 GHz (corresponding to NR band n78) with a bandwidth of 50 MHz.

b) *V2X application*: Each vehicle runs an uplink application streaming LiDAR data, modeled according to the Kitti multi-modal dataset [20]. Moreover, it receives in the downlink commands from the remote driver for teleoperated driving operations. Each raw LiDAR perception generates a point cloud of around 120 000 points at 10 Hz, with an average file size of 3 200 KB. Based on [17], compression is accomplished using Draco [21] (a software designed by Google to compress 3D-like data) combined with the semantic segmentation functionalities of RangeNet++ [22] (a NN able to assign class labels to data points). Our compression pipeline consists of the following steps:

- We first infer semantic segmentation of LiDAR data with RangeNet++, so as to identify the most valuable objects in the scene. Three semantic levels are defined:

TABLE II: Application modes and RL reward parameters.

Application Mode	RangeNet++	Draco		Avg. file size [KB]	Chamfer Distance CD_{sym}
		Compression	Quantization		
0 (baseline)	NO	NO		3200	0
1450	NO	5	14	200	0.000044
1451	LEVEL 1	5	14	104	5.476881
1452	LEVEL 2	5	14	17	35.634660

- LEVEL 0: The raw LiDAR acquisition is considered.
- LEVEL 1: The points associated to the road elements are removed from the cloud, thus reducing the file size.
- LEVEL 2: The points associated to buildings, vegetation, and traffic signs are also removed; the resulting data consists only of dynamic elements like pedestrians and vehicles, i.e., the most critical elements in autonomous and remote driving scenarios.
- The resulting point cloud is then compressed with Draco, which defines 15 quantization levels and 11 compression levels that trade off efficiency against speed.

We consider four application modes, as reported in Table II.

c) RL agent architecture: The state of our RL agent includes only a subset of the RAN-AI input parameters (Sec. II-B), to reduce the state space dimension, and ensure a faster learning convergence. Specifically, we focused on RAN-level metrics usually available at the gNB, averaged over a reporting period of 100 ms, i.e., the value of the MCS, the number of Orthogonal Frequency Division Multiplexing (OFDM) symbols used to transmit, the value of the SINR, the mean/max/min/std of the packets' delay, and the PRR at the PDCP layer. We also limited the agent's action space to three actions, corresponding to application modes {1450, 1451, 1452},¹ so that $S = 8$ and $A = 3$.

With respect to DQL, we set the discount factor of the RL algorithm to $\lambda = 0.95$ and we adopt a batch-learning approach, where the learning transitions are grouped in batches of size $B_{\text{size}} = 10$. We implement the Adaptive moment estimator (Adam) algorithm to update the weights of the NN, considering $\zeta = 10^{-5}$ as the maximum learning rate, and $\epsilon = 10^{-3}$ as the weight decay [23]. The agent's training follows two subsequent phases, organized into episodes of 800 steps each, where the step duration is set to 100 ms so that each episode lasts 80 seconds. First, we perform an *offline* training phase, where the agent's actions are kept fixed for the whole duration of the episode. Then, we perform an *online* training phase, where the agent's actions change in time according to an ϵ -greedy exploration policy. The duration of each training phase depends on the number of vehicles in the scenario, and is set to 2500 when $N_u = 1$ and 500 when $N_u = 5$.

In terms of the reward function in Eq. (2), we set $\delta_m = 50$ ms, as specified by the 5GAA for teleoperated driving applications [24], and $CD_{\text{sym},m} = 45$, while the Chamfer Distance

¹From preliminary experiments, we obtained that application mode 0 (where data are not compressed nor segmented) was never selected by the agent during training, so it was not included in the action space to promote faster convergence; this strategy will represent the benchmark solution in our tests.

CD_{sym} associated to each application mode is reported in Table II. Finally, we set $\alpha \in \{0.5, 1\}$ to consider different QoS/QoE weights.

d) Performance evaluation: To validate our framework, we compared the following action policies:

- *DQL* (proposed), where at each step the agent implements our proposed RL framework described in Sec. III;
- *Constant* (benchmark), where at the beginning of the simulation the end user maintains one application mode among $\{0, 1450, 1451, 1452\}$ for the whole simulation.

The two strategies have been tested separately, and will be compared in terms of (i) the reward gained by the agent, (ii) the probability to satisfy the QoS requirements, (iii) the user's QoS, expressed in terms of delay at the application layer, and (iv) the user's QoE, expressed in terms of CD_{sym} .

B. Numerical Results

First, in Fig. 1 we evaluate the DQL statistics of the training phase vs. N_u and α . We observe that, during the first phase of the training period, the agent mostly makes random decisions because of the ϵ -greedy policy, which ensures that the action and state spaces are fully explored. As the training progresses, the agent starts acting greedily and prioritizes the actions that maximize the long-term reward. In this case, the agent's decisions depend on the observed state and, consequently, the action probability distribution may change completely from an episode to another, which justifies the variability in Fig. 1. We can see that, when $N_u = 1$, the agent prefers application mode 1450, since it offers the best trade-off between QoS and QoE compared to other actions. Instead, when $N_u = 5$, the system is more congested, and the agent tends to penalize more the actions leading to a violation of the QoS requirements; as such, action 1451 (where the original LiDAR data is subject to both compression and segmentation before transmission) is preferred most of the time.

Second, we analyze the QoE and QoS performance of the DQL agent against other benchmarks in different system configurations, during a test phase of 100 episodes. In Fig. 2, we illustrate the distribution of the Chamfer Distance (an indication of the QoE of the system), considering $N_u = 1$ and for $\alpha \in \{0.5, 1.0\}$. We observe that, when the RAN-AI entity implements strategies $\{0, 1450, 1451, 1452\}$, the Chamfer Distance remains constant throughout all the episodes. Instead, with DQL, the RAN-AI entity tries to adapt to the conditions of the environment. In particular, when $\alpha = 1$, the agent has an incentive to improve the QoE's reward component, and prioritizes action 1450, so that CD_{sym} is 0 for more than 45%

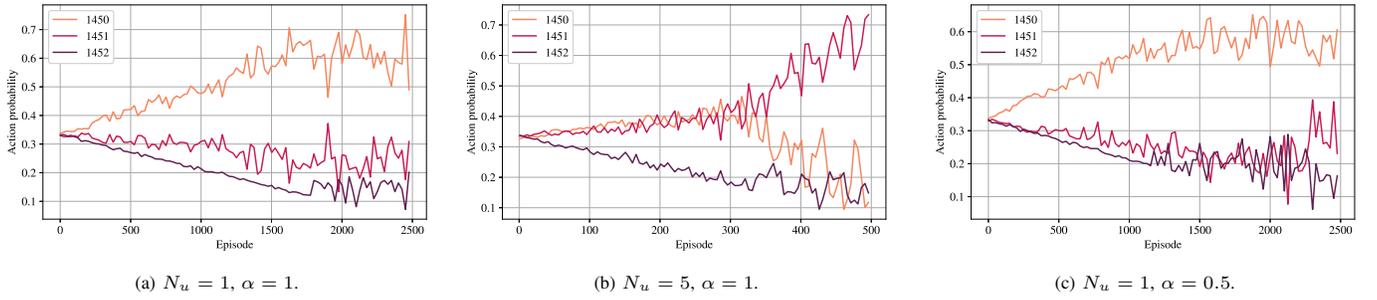


Fig. 1: Action probability during the training phase, vs. the number of vehicles N_u and α .

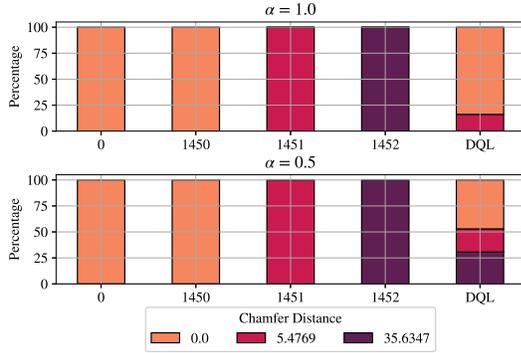


Fig. 2: Chamfer Distance (CD_{sym}) distribution, for $N_u = 1$.

of the time. On the other hand, for $\alpha = 0.5$, the reward function discourages an overly aggressive behavior, which may lead to a violation of the QoS requirements, and prefers action 1452, which however leads to some QoE degradation. Similarly, Fig. 3 plots the distribution of the QoS for the different strategies, where a QoS of 0 implies that the delay/PRR requirements of the teleoperated driving application are not satisfied. We observe that DQL improves as α decreases, since the agent receives a higher reward when the delay is minimized. In particular, when $\alpha = 1$ ($\alpha = 0.5$), the QoS requirements are satisfied around 50% (75%) of the time. Notably, while the Constant 1452 action ensures better QoS (since the data are compressed and segmented before transmission, which can reduce the size of the packets to send), it is characterized by a very high CD_{sym} (Fig. 2) and, consequently, a bad reward distribution. Similarly, while the Constant 1450 action promotes better QoE, it results in a QoS degradation of up to 40% compared to DQL when $\alpha = 0.5$.

The above considerations are validated in Fig. 4, which reports the distribution of the mean packet delay at the application layer, i.e., the distribution of the average delay of all the packets delivered at the application layer during the same time slot. Specifically, we use the box plot representation, where the black line in the middle of each box is the median, the edges represent the 25th and 75th percentiles, while the whiskers identify the outliers of the distribution. We observe that all the policies can maintain the median delay below the required threshold ($\delta_m = 50$ ms) for both $N_u = 1$ and $N_u = 5$. However, while for Constant 1450 and 1451 the top whiskers of the distribution approach a value of 100 ms, which indicates that such strategies may be unable to meet the QoS requirements when the working conditions are critical,

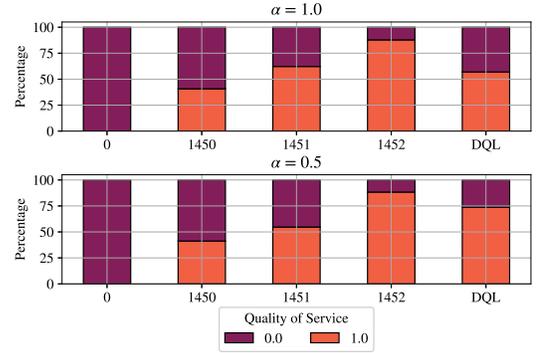


Fig. 3: QoS distribution, when for $N_u = 1$.

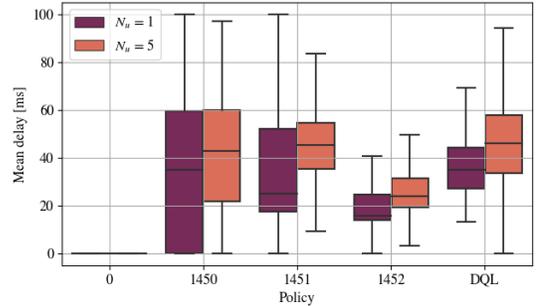


Fig. 4: Mean delay (application layer), for $\alpha = 1$.

for DQL the whiskers are equal to 70 and 95 ms for $N_u = 1$ and $N_u = 5$, respectively. Notice that the only strategy with a better QoS is Constant 1452, resulting in a delay distribution squeezed towards zero, thus ensuring that QoS requirements are rarely violated, at the expense of a poor QoE.

To evaluate the trade-off between QoE and QoS of the different PQoS policies, in Fig. 5 we plot the agent’s reward (normalized in $[-1, +1]$) for different values of α and N_u . Specifically (i) the white dot represents the median, (ii) the thick black bar in the center represents the interquartile range, (iii) the thin black line represents the rest of the distribution, except for “outliers.” Wider sections of the violin plot indicate values that occur more frequently. From Fig. 5a we can see that the DQL solution achieves the best trade-off between QoS and QoE, meaning that the adaptive behavior of DQL is desirable for PQoS. We notice that the Constant 0 and 1452 policies display the lowest reward due to their poor QoS and QoE performance, respectively. For $N_u = 1$, DQL and the Constant 1450 and 1451 policies achieve similar rewards. On the other hand, when $N_u = 5$, i.e., considering more congested

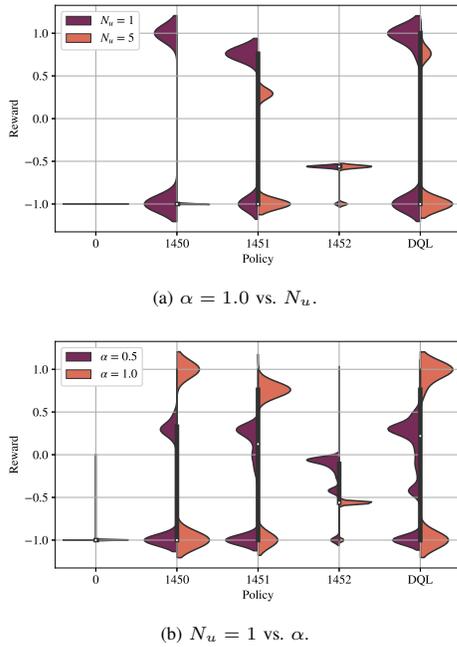


Fig. 5: Distribution of the agent's reward, normalized in $[-1, +1]$.

networks, the advantages of DQL are more evident. Notably, DQL is the only strategy that can provide a reward higher than 0.5, while for Constant 1450 and 1451 the reward peaks at about -1 and 0.4 , respectively. By tuning parameter α , it is possible to further adjust the policy learned by the DQL agent, promoting more or less conservative communication settings. From Fig. 5b, we observe that with $\alpha = 0.5$ the DQL agent is able to reduce the probability of negative rewards, compared to $\alpha = 1$. In this case, however, the maximum reward is 0.3 (vs. 1 when $\alpha = 1$), through still higher than any competitor.

V. CONCLUSIONS AND FUTURE WORKS

In this paper we analyzed the potential of PQoS to predict and optimize autonomous networks. Specifically, we proposed the design of a new “RAN-AI” entity, installed in the gNB and interacting with a custom RL agent, to identify the optimal set of PQoS actions/countermeasures to satisfy QoS requirements of end users. We performed simulations in ns-3 in a teleoperated driving scenario, and demonstrated that the adaptive behavior of our RL model can achieve the best trade-off between QoS and QoE performance, compared to other baseline solutions that do not support machine learning. Also, we make the case that, by properly tuning the parameters of the reward function, it is possible to adjust the policy learned by the learning agent, promoting more or less conservative communication settings.

Among our future research activities, we will extend our current PQoS framework to incorporate advanced functionalities, including the support for multi-cell scenarios, as well as the definition of new ML technologies and countermeasures for PQoS (e.g., based on federated learning). In particular, we will investigate whether a fully distributed architecture, in which end machines make autonomous decisions, can promote more efficient PQoS.

REFERENCES

- [1] H. Shariatmadari, R. Ratasuk, S. Iraj, A. Laya, T. Taleb, R. Jäntti, and A. Ghosh, “Machine-type communications: current status and future perspectives toward 5G systems,” *IEEE Communications Magazine*, vol. 53, no. 9, pp. 10–17, Sep. 2015.
- [2] M. Boban, M. Giordani, and M. Zorzi, “Predictive Quality of Service (PQoS): The Next Frontier for Fully Autonomous Systems,” *IEEE Network*, 2021.
- [3] 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.
- [4] P. Del Moral, “Nonlinear filtering: Interacting particle resolution,” *Comptes Rendus de l'Académie des Sciences-Series I-Mathematics*, vol. 325, no. 6, pp. 653–658, Sep. 1997.
- [5] E. A. Wan and R. V. D. Merwe, “The unscented Kalman filter for nonlinear estimation,” in *IEEE Adaptive Systems for Signal Processing, Communications, and Control Symposium*, Oct. 2000, pp. 153–158.
- [6] F. Mason, M. Giordani, F. Chiariotti, A. Zanella, and M. Zorzi, “An adaptive broadcasting strategy for efficient dynamic mapping in vehicular networks,” *IEEE Transactions on Wireless Communications*, vol. 19, no. 8, pp. 5605–5620, Aug. 2020.
- [7] 5GAA Automotive Association, “Making 5G Proactive and Predictive for the Automotive Industry,” *White Paper*, Aug. 2019.
- [8] A. Carron, M. Todescato, R. Carli, L. Schenato, and G. Pillonetto, “Machine learning meets Kalman filtering,” in *IEEE 55th Conference on Decision and Control (CDC)*, 2016.
- [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. Nardo, D. Peressoni, P. Testolina, M. Giordani, and A. Zanella, “Point Cloud Compression for Autonomous Driving: A Performance Comparison,” in *IEEE Wireless Communications and Networking Conference (WCNC) [To Appear]*, 2022.
- [11] R. Bellman, “Dynamic programming,” *Science*, vol. 153, no. 3731, pp. 34–37, Jul. 1966.
- [12] A. F. Agarap, “Deep learning using rectified linear units (ReLU),” *arXiv preprint arXiv:1803.08375*, 2018.
- [13] H. Van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double Q-learning,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 30, no. 1, 2016.
- [14] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3-4, pp. 279–292, May 1992.
- [15] 3GPP, “Service requirements for enhanced V2X scenarios (Release 15),” *TS 22.186*, Sep. 2018.
- [16] V. Rossi, P. Testolina, M. Giordani, and M. Zorzi, “On the Role of Sensor Fusion for Object Detection in Future Vehicular Networks,” in *Joint European Conference on Networks and Communications 6G Summit (EuCNC/6G Summit)*, 2021.
- [17] A. Varischio, F. Mandruzzato, M. Bullo, M. Giordani, P. Testolina, and M. Zorzi, “Hybrid Point Cloud Semantic Compression for Automotive Sensors: A Performance Evaluation,” *IEEE International Conference on Communications (ICC)*, 2021.
- [18] M. Mezzavilla, M. Zhang, M. Polese, R. Ford, S. Dutta, S. Rangan, and M. Zorzi, “End-to-End Simulation of 5G mmWave Networks,” *IEEE Communications Surveys and Tutorials*, vol. 20, no. 3, pp. 2237–2263, Apr. 2018.
- [19] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, “Recent development and applications of SUMO - Simulation of Urban MObility,” *International Journal On Advances in Systems and Measurements*, vol. 5, no. 3&4, pp. 128–138, Dec. 2012.
- [20] Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago, “The KITTI Vision Benchmark Suite.” [Online]. Available: <http://www.cvlibs.net/datasets/kitti/>
- [21] Google, “Draco 3D Data Compression,” 2017, [Online]. Available: <https://github.com/google/draco>.
- [22] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, “RangeNet++: Fast and accurate LiDAR semantic segmentation,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- [23] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [24] 5GAA, “C-V2X Use Cases Volume II: Examples and Service Level Requirements,” *White Paper*, Oct. 2020.