



UNIVERSITY OF CATANIA  
DEPARTMENT OF ELECTRICAL, ELECTRONICS  
AND COMPUTER ENGINEERING  
PH.D. COURSE IN SYSTEMS, ENERGY, COMPUTER AND  
TELECOMMUNICATION ENGINEERING

---

OPTIMIZATION OF THE 5G NR INTERFACE  
BASED ON INNOVATIVE TECHNIQUES

Doctoral Dissertation of:  
**Luciano Miuccio**

Tutor:  
**Prof. Daniela Panno**

Coordinator:  
**Prof. Paolo Arena**

XXXV Cycle





*To all the people I love.*



---

---

## Acknowledgments

---

I would like to take this opportunity to express my gratitude to my supervisors, family, fiancée, colleagues, and friends who gave me encouragement and support through all these years.

First and foremost, I would like to express my deepest gratitude to my supervisor, Prof. Daniela Panno for all the support and encouragement she gave me. During this difficult pandemic period, she has always motivated me and allowed to perform extraordinary experiences. I also want to express my gratitude to Prof. Salvatore Riolo for his unconditional trust, help, friendship and discussions during this intense journey.

Moreover, I also like to thank Prof. Mehdi Bennis with the Centre for Wireless Communications, Oulu, Finland, for all the generous support and guidance he provided during my internship abroad. He shared his knowledge and expertise with me, being always helpful and available.

A big thank goes to my fiancée, Marialuisa, who has always been on my side with constant support, encouragement and understanding. Moreover, I sincerely appreciate my parents, sister, grandparents, uncles/aunties, and all Marialuisa's family, for the endless love and support through all this time.

I want also to acknowledge the support and help of all my friends. Finally, thanks to all my colleagues in Catania and Oulu for creating unforgettable memories together.



---

---

## Abstract

---

**T**HE fifth generation (5G) of cellular networks aims at providing connectivity for a large number of applications. To achieve this goal, 5G has been designed considering three scenarios with vastly heterogeneous requirements: Enhanced Mobile Broadband (eMBB), massive Machine-Type Communications (mMTC), and Ultra-Reliable Low-latency Communications (URLLC). On the basis of the stringent Key Performance Indicators (KPIs), optimizing separately the services belonging to each one of these scenarios is already challenging. Furthermore, simultaneously addressing the KPIs of the different scenarios is even more difficult. Considering the usage of common techniques cannot permit to overcome the fundamental trade-off among the achievable data rate, latency, reliability, and spectral efficiency. Therefore, to efficiently deal with these challenges, several optimization aspects exploiting innovative techniques should be introduced.

In the view of 5G and future 6G wireless networks, the aim of this Dissertation is to provide the capability to efficiently support different service classes and their diverse Quality of Service (QoS) requirements. To do so, we optimize several aspects of the Radio Resource Management (RRM) level and Medium Access Control (MAC) procedures, by exploiting innovative techniques. The analysis provided considers different access techniques, like Time Division Multiple Access (TDMA), Orthogonal Frequency Division Multiple Access (OFDMA) and the novel Sparse Code Multiple Access (SCMA), the support of advanced wave-

## Abstract

---

form, such as filtered OFDM, and Universal Filtered Multi-Carrier (UFMC), and the exploitation of a multi-numerology frame structure with different sub-carrier spacing. In addition, innovative technique of Access Class Barring (ACB) and predictive estimation of the traffic load are considered. Finally, the introduction of Deep Learning (DL) techniques and Deep Reinforcement Learning (DRL) algorithms permitted to overcome limitations set by mathematical and statistical approaches. The research activity can be grouped into three main areas.

The first one considers the coexistence among a wide variety of services inside the same OFDM grid with the aim of providing flexibility and improvement of performances. Initially, we propose a QoS-aware and Channel-aware two-levels RRM framework, that appropriately allocates the band spectrum to the numerologies and properly assigns the Physical Resource Blocks (PRBs) to each numerology. Moreover, we implement a new Simulation Environment to identify the best (waveform, Guard-Band (GB) size) pair which reduces the Inter-Numerology Interference (INI) phenomenon, and maximizes the spectral efficiency while taking into account QoS requirements.

In the second research activity, we focus on the optimization of the RRM for the mMTC services. We propose a new framework that includes a joint control of the dynamic resource allocation between the Physical Random Access Channel (PRACH) and the Physical Uplink Shared Channel (PUSCH), and a new random access procedure based on an adaptive ACB scheme. To further increase the spectral efficiency, we adopt the Sparse Code Multiple Access (SCMA) technique for the transmission in the PUSCH resources. Then, instead of improving the succeeded access attempts in the PRACH, we present the innovative transmission idea to exploit the unused PUSCH resources to serve an additional part of MTC devices. Moreover, we propose an accurate current access attempts estimation method, based on Deep Neural Network (DNN), which accepts as input only the information really available at the next generation NodeB (gNB). Finally, to improve the transmission performance of SCMA in practical networks we design an end-to-end SCMA en/decoding scheme robust to the channel noise.

Finally, in the third activity, we consider the problem of automatically learning MAC protocols with good generalization proprieties across sev-



---

eral transmission environments. These protocols take into account both the control plane and the data plane point of view, and are learned between several User Equipments (UEs) cast as Reinforcement Learning (RL) agents and one Base Station (BS) cast as an expert.



---

---

# Contents

---

<b>Acknowledgments</b>	<b>III</b>
<b>Abstract</b>	<b>VII</b>
<b>List of Publications</b>	<b>XIII</b>
<b>List of Acronyms</b>	<b>XV</b>
<b>List of Figures</b>	<b>XIX</b>
<b>List of Tables</b>	<b>XXV</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 Research Questions . . . . .	7
1.3 Research Contributions and Thesis outline . . . . .	7
<b>2 A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems</b>	<b>11</b>
2.1 Overview . . . . .	11
2.2 5G Downlink Frame Structure and Assumptions . . . . .	14
2.3 System Model and Problem Formulation . . . . .	15
2.4 Benchmark Schemes . . . . .	20
2.4.1 1st level allocation algorithms . . . . .	21
2.4.2 2nd level scheduling algorithms . . . . .	21

## Contents

---

2.5	The Proposed Control Framework . . . . .	24
2.5.1	Phase 1.1 . . . . .	24
2.5.2	Phase 1.2 . . . . .	32
2.5.3	Time complexity analysis . . . . .	38
2.6	Performance Evaluation . . . . .	40
2.6.1	Stationary analysis with homogeneous GBR requirements . . . . .	42
2.6.2	Stationary analysis with heterogeneous GBR requirements . . . . .	46
2.6.3	Dynamic analysis with homogeneous GBR requirements . . . . .	47
2.7	Related Works . . . . .	50
2.7.1	New 5G NR waveforms . . . . .	51
2.8	THE VIENNA 5G Link Level (LL) SIMULATOR . . . . .	53
2.9	System Model and QoS-aware RRM framework . . . . .	54
2.10	The Proposed Simulation Environment for the Multi-Numerology scenario . . . . .	56
2.10.1	Adaptation of the 5G Vienna LL Simulator scheduling for the multi-numerology scenario . . . . .	58
2.11	Case Study Analysis . . . . .	60
2.11.1	Simulation Configuration and Parameters . . . . .	61
2.11.2	Performance Analysis . . . . .	62
2.12	Conclusion . . . . .	63
<b>3</b>	<b>Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios</b>	<b>65</b>
3.1	Overview . . . . .	65
3.2	Background . . . . .	70
3.2.1	LTE Uplink . . . . .	70
3.2.2	eMTC and NB-IoT . . . . .	73
3.2.3	5G Uplink frame structure . . . . .	74
3.3	Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA . . . . .	75
3.3.1	System Model . . . . .	75
3.3.2	Problem Formulation . . . . .	83

---

3.3.3 Overview of the proposed Control Framework . . .	88
3.3.4 Enhanced Dynamic Clustering Dimensioning Algorithm . . . . .	90
3.3.5 Performance Evaluation . . . . .	100
3.3.6 Appendix . . . . .	113
3.3.7 Proof of Equation (3.77) . . . . .	113
3.4 Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH . . . . .	116
3.4.1 System Model and Background . . . . .	116
3.4.2 Our proposal . . . . .	118
3.4.3 Performance Evaluation . . . . .	123
3.5 A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond . . . . .	126
3.5.1 Motivations and related work . . . . .	126
3.5.2 Contributions . . . . .	128
3.5.3 System Model for the PRACH traffic load prediction	129
3.5.4 The proposed DNN-based traffic load estimation method . . . . .	130
3.5.5 Comparison and assessment of estimation methods in a stand-alone RA cycle . . . . .	136
3.5.6 Comparison and assessment of the estimation methods in a long-term analysis . . . . .	140
3.5.7 Appendix . . . . .	151
3.5.8 Calculation of $\bar{P}_S$ , $\bar{P}_C$ and the pmf of $P_C$ . . . . .	151
3.6 A Wasserstein GAN autoencoder for SCMA networks . .	153
3.6.1 System model and Problem Formulation . . . . .	153
3.6.2 Our solution . . . . .	155
3.6.3 Performance evaluation . . . . .	160
<b>4 Learning Generalized Wireless MAC Communication Protocols via Abstraction</b>	<b>163</b>
4.1 Overview . . . . .	163
4.2 System Model . . . . .	165
4.3 UEs-BS Interaction as an MARL Problem . . . . .	167
4.4 Policy learning via abstraction . . . . .	169

## Contents

---

4.4.1 Abstract Formulation . . . . .	169
4.5 Performance Evaluation . . . . .	173
4.5.1 Setting . . . . .	173
4.5.2 Benchmarks . . . . .	174
4.5.3 Results and Discussion . . . . .	174
<b>5 Conclusions and Perspectives</b>	<b>179</b>
5.1 Conclusions . . . . .	179
5.2 Works in progress and Future works . . . . .	180
<b>Bibliography</b>	<b>183</b>

---

---

## List of Publications

---

### Journal Papers

---

- **L. Miuccio**, D. Panno and S. Riolo, "A Wasserstein GAN Autoencoder for SCMA Networks," in *IEEE Wireless Communications Letters*, vol. 11, no. 6, pp. 1298-1302, June 2022.
- **L. Miuccio**, D. Panno, P. Pisacane and S. Riolo, "A QoS-aware and channel-aware Radio Resource Management framework for multi-numerology systems," *Computer Communications*, vol. 191, pp. 299-314, 2022.
- **L. Miuccio**, D. Panno and S. Riolo, "A DNN-based estimate of the PRACH traffic load for massive IoT scenarios in 5G networks and beyond," *Computer Networks*, vol. 201, p. 108608, 2021.
- S. Riolo, D. Panno and **L. Miuccio**, "Modeling and Analysis of Tagged Preamble Transmissions in Random Access Procedure for mMTC Scenarios," *IEEE Transactions on Wireless Communications.*, vol. 20, no. 7, pp. 4296-4312, July 2021.<sup>1</sup>
- **L. Miuccio**, D. Panno and S. Riolo, "A New Contention-Based PUSCH Resource Allocation in 5G NR for mMTC Scenarios," in *IEEE Communications Letters*, vol. 25, no. 3, pp. 802-806, March 2021.

---

<sup>1</sup>Paper not included in this dissertation.

## List of Publications

---

- **L. Miuccio**, D. Panno and S. Riolo, "Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for Massive MTC in 5G NR Based on SCMA," in *IEEE Internet of Things Journal*, vol. 7, no. 6, pp. 5042-5063, June 2020.

## Conference Papers

---

- **L. Miuccio**, S. Riolo, S. Samarakoon, D. Panno, and M. Bennis, "Learning Generalized Wireless MAC Communication Protocols via Abstraction," accepted for publication in IEEE GLOBECOM 2022.
- **L. Miuccio**, D. Panno, P. Pollara and S. Riolo, "Implementation of a Simulation Environment for Multi-Numerology Scenarios in 5G Vienna Link Level Simulator," *2021 IEEE International Conference on Computing (ICOCO)*, 2021, pp. 134-139.
- **L. Miuccio**, D. Panno, P. Pisacane and S. Riolo, "Channel-Aware and QoS-Aware Downlink Resource Allocation for Multi-Numerology Based 5G NR Systems," *2021 19th Mediterranean Communication and Computer Networking Conference (MedComNet)*, pp. 1-8.
- **L. Miuccio**, D. Panno and S. Riolo, "Joint Congestion Control and Resource Allocation for Massive MTC in 5G Networks Based on SCMA," *2019 15th International Conference on Telecommunications (ConTEL)*, Graz, Austria, 2019, pp. 1-8.<sup>1</sup>
- **L. Miuccio**, D. Panno and S. Riolo, "Dynamic Uplink Resource Dimensioning for Massive MTC in 5G Networks Based on SCMA," *European Wireless 2019; 25th European Wireless Conference*, Aarhus, Denmark, 2019, pp. 1-6.<sup>1</sup>

---

<sup>1</sup> Paper not included in this dissertation.



---

---

## List of Acronyms

---

<b>2G</b>	second-generation
<b>3GPP</b>	Third Generation Partnership Project
<b>3G</b>	third-generation
<b>4G</b>	fourth-generation
<b>5G</b>	fifth-generation
<b>6G</b>	sixth-generation
<b>ACB</b>	Access Class Barring
<b>AE</b>	Autoencoder
<b>AI</b>	Artificial Intelligence
<b>AWGN</b>	Additive White Gaussian Noise
<b>B5G</b>	Beyond 5G
<b>BE</b>	Best Effort
<b>BER</b>	Bit Error Rate
<b>BPSK</b>	Binary Phase Shift Keying
<b>BS</b>	Base Station
<b>CDF</b>	Cumulative Density Function
<b>CDMA</b>	Code Division Multiple Access
<b>CP</b>	Cyclic Prefix
<b>CQI</b>	Channel Quality Indicator
<b>CTA</b>	Contention Tree Algorithm
<b>DAW</b>	Default Averaging Window
<b>DNN</b>	Deep Neural Network
<b>DRL</b>	Deep Reinforcement Learning

## List of Acronyms

---

<b>DRX</b>	Discontinuous Reception
<b>DTX</b>	Discontinuous Transmission
<b>eMBB</b>	enhanced Mobile Broadband
<b>eMTC</b>	enhanced MTC
<b>FBMC</b>	Filter-Band Multi-Carrier
<b>FER</b>	Frame Error Rate
<b>FFT</b>	Fast Fourier transform
<b>GB</b>	Guard Band
<b>gNB</b>	next Generation NodeB
<b>GBR</b>	Guaranteed Bit Rate
<b>IFFT</b>	Inverse Fast Fourier Transform
<b>INI</b>	Inter-Numerology Interference
<b>ITU</b>	International Telecommunication Union
<b>KPI</b>	Key Performance Indicator
<b>LDPC</b>	Low-Density Parity-Check
<b>LDS-CDMA</b>	Low-Density Spreading-based CDMA
<b>LDS-OFDM</b>	Low-Density Spreading-based OFDM
<b>LL</b>	Link Level
<b>LTE</b>	Long Term Evolution
<b>mIOT</b>	massive Internet of Things
<b>mMTC</b>	massive Machine Type Communications
<b>MAC</b>	Medium Access Control
<b>MARL</b>	Multi-Agent Reinforcement Learning
<b>MCS</b>	Modulation and Code Scheme
<b>MDP</b>	Markov Decision Process
<b>MPOMPD</b>	Multi-Agent Partially Observable MDP
<b>NB-IoT</b>	NarrowBand IoT
<b>NOMA</b>	Non-Orthogonal Multiple Access
<b>NR</b>	New Radio
<b>NSA</b>	Non-StandAlone
<b>OA</b>	Observation Abstraction
<b>OFDM</b>	Orthogonal Frequency Division Multiplexing
<b>OFDMA</b>	Orthogonal Frequency Division Multiple Access
<b>OOB</b>	Out-Of-Band
<b>OMA</b>	Orthogonal Multiple Access
<b>PAC</b>	Pearson Auto-Correlation Coefficient

---

<b>PAPR</b>	Peak-To-Average Power Ratio
<b>PDF</b>	Probability Density Function
<b>PDSCH</b>	Physical Downlink Shared Channel
<b>pmf</b>	probability mass function
<b>PRACH</b>	Physical Random Access Channel
<b>PRB</b>	Physical Resource Block
<b>PUSCH</b>	Physical Uplink Shared Channel
<b>QAM</b>	Quadrature Amplitude Modulation
<b>QoS</b>	Quality of Service
<b>RA</b>	Random Access
<b>RE</b>	Resource Element
<b>ReLU</b>	Rectified Linear Unit
<b>RL</b>	Reinforcement Learning
<b>RRM</b>	Radio Resource Management
<b>RX</b>	Reception
<b>SA</b>	StandAlone
<b>SB</b>	Sub-Band
<b>SER</b>	Symbol Error Rate
<b>SIC</b>	Successive Interference Cancellation
<b>SINR</b>	Signal to Interference plus Noise Ratio
<b>SNR</b>	Signal to Noise-Ratio
<b>SCMA</b>	Sparse Code Multiple Access
<b>SCS</b>	Sub-Carrier Spacings
<b>TBLER</b>	Transport Block Error Rate
<b>TDMA</b>	Time Division Multiple Access
<b>TTI</b>	Transmission Time Interval
<b>TX</b>	Transmission
<b>UE</b>	User Equipment
<b>UFMC</b>	Universal Filtered Multi-Carrier
<b>URLLC</b>	Ultra-Reliable and Low Latency Communications
<b>WGAN</b>	Wasserstein Generative Adversarial Networks
<b>ZP</b>	Zero-Postfix



---

## List of Figures

---

1.1	5G usage scenarios . . . . .	3
2.1	Multi-Numerology Resource Grid. . . . .	15
2.2	1 Sub-Band (SB), i.e., granularity of the minimum amount of resources that can be assigned to each numerology. . . .	17
2.3	1st level allocation algorithms at a glance. . . . .	22
2.4	Two control levels framework for radio resource allocation and scheduling. . . . .	25
2.5	Radio Resource Management Flow diagram, adopting CARAM-new as 1st level allocation algorithm. . . . .	26
2.6	Flow diagram of the preliminary control framework proposed in [1], adopting the previous version of CARAM as 1st level allocation algorithm. . . . .	27
2.7	Steps of the SB allocation: (a) Step 3; (b) Step 4 first round ; (c) Step 4 second round; (d) final SB allocation. . .	37
2.8	Average runtime together with the standard deviation (shaded area) for the Phase 1.1 of the new-CARAM under different $N_{UE}$ values. . . . .	39
2.9	Average runtime together with the standard deviation (shaded area) for the Phase 1.2 of the new-CARAM under different $N_{SB}$ values. . . . .	39
2.10	UEs Arrangements. . . . .	43

## List of Figures

---

2.11 Stationary analysis with homogeneous GBR requirements and UE ratio 1:1:1. Medium load scenario with UEs arrangement 1. . . . .	45
2.12 Stationary analysis with homogeneous GBR requirements and UE ratio 1:1:1. High load scenario with UEs arrangement 2. . . . .	45
2.13 Stationary analysis with homogeneous GBR requirements. (a)-(b) UE ratio 1:1:4 and medium load scenario with UEs arrangement 1. (c)-(d) UE ratio 4:1:1 and high load scenario with UEs arrangement 2. . . . .	46
2.14 Stationary analysis with heterogeneous GBR requirements.	46
2.15 Dynamic analysis of a single simulation test with homogeneous GBR requirements. Medium-high load scenario. . . . .	47
2.16 Dynamic analysis of a single simulation test with homogeneous GBR requirements. Overload scenario. . . . .	48
2.17 Out of band emission for the considered waveforms. . . . .	52
2.18 Flow chart of the proposed Simulation Phase . . . . .	57
2.19 Radio resource allotment for the different phases of the Simulation Environment . . . . .	58
2.20 Adaptation procedure for the GB insertion in the proposed Simulation Environment . . . . .	60
2.21 Throughput of the Best Effort Numerology vs different waveforms and guard band sizes. . . . .	62
3.1 Uplink frame structure. . . . .	69
3.2 A typical RA cycle of 5ms. . . . .	71
3.3 Conventional 4-step RA procedure, when the communication is successful on the second attempt, and signaling messages for data packet transmission in LTE/LTE-A and eMTC technologies. . . . .	72
3.4 Example of multiple access and bit-to-codeword mapping of an SCMA encoder with $Q = 4$ , $S = 2$ , $K_{max} = 3$ , $L_{SB} = 6$ , $I = 4$ . . . . .	76
3.5 SCMA Codebooks. Example of the first dimension of the codebook in RE1. . . . .	77
3.6 The proposed basic two-step RA procedure. . . . .	81

3.7 Two-step RA procedure and data packet transmission. . . .	82
3.8 Pearson Auto-Correlation Coefficient (PAC) of the random process $\{N(\omega, n)\}$ . . . . .	84
3.9 Number of succeeded MTC devices transmission vs number of MTC devices which carry out the RA procedure for different $T_{pr}$ values when $\theta_{max} = 160$ bits. . . . .	86
3.10 Flow diagram of the MTC device behavior which adopts the proposed joint control scheme. . . . .	91
3.11 Examples of Cluster, each one has its own $T_{pr}$ dimension. . .	92
3.12 Enhanced Dynamic Cluster Dimensioning Algorithm (eD-CDA) phases. . . . .	95
3.13 Representation of $V_{min}$ vs $\tilde{\sigma}_{r_i}^j$ when $\theta_{max} = 160$ bits. . . .	100
3.14 Representation of the incremental ratio of $V_{min}$ with respect to $\tilde{\sigma}_{r_i}^j$ vs $\tilde{\sigma}_{r_i}^j$ . . . . .	101
3.15 Flow diagram of the MTC device behavior which adopts the adapted $ACB_{p^*}$ . . . . .	102
3.16 Connectivity Capability Loss in terms of number of failed communications, when $N_{MTC} = 50000$ . . . . .	104
3.17 Connectivity Capability Loss in terms of number of failed communications, under different $N_{MTC}$ values. . . . .	104
3.18 Mean value of the service delay for succeeded communication, under different $N_{MTC}$ values. . . . .	105
3.19 Adopted energy model. . . . .	107
3.20 Energy blocks and energy consumption events. . . . .	107
3.21 Analytical representation of the energy consumed by adopting the proposed control framework and the $ACB_{p^*}$ scheme. . . . .	110
3.22 Mean value of the energy consumed per MTC device for a single succeeded communication. . . . .	111
3.23 CDF of the MTC device energy consumption for one succeeded communication. . . . .	111
3.24 Number of successful transmissions in an RA cycle vs the number of MTC devices ( $M$ ) which carry out the RA procedure for different $T_{pr}$ values. . . . .	119
3.25 Throughput Gain under different systems and $M$ . . . . .	126
3.26 Energy Consumption Indicator under different systems and $M$ . . . . .	126

## List of Figures

---

3.27 Example of Step 1 of the RA procedure. . . . .	131
3.28 Architecture of $DNN_K$ . . . . .	132
3.29 Average number of succeeded and collided preambles vs $M$ for different $L$ values. . . . .	134
3.30 RMSE in the estimation of $M$ vs $M$ values with the $T_{pr} = 1$ dimensioning . . . . .	140
3.31 Density scatter plots in the estimation of $M$ vs $M$ values with the $T_{pr} = 1$ dimensioning. . . . .	140
3.32 RMSE in the estimation of $M$ vs $M$ values . . . . .	141
3.33 Density scatter plots in the estimation of $M$ vs $M$ values with the $T_{pr} = 4$ dimensioning. . . . .	142
3.34 Expected number of successful communications vs the number of MTC devices ( $M$ ) which attempts access in a given RA cycle, assuming different PRACH dimensionings. . . . .	144
3.35 The temporal distributions of actual or estimated $M$ values with different control schemes for the simulated Beta arrival distribution in the space of 10 s. . . . .	145
3.36 Confusion Matrices with the simulated Beta arrival distribution. . . . .	146
3.37 The temporal distributions of actual or estimated $M$ values with different control schemes for the real arrival distribution in the space of 10 s. . . . .	147
3.38 Confusion Matrices for the real arrival distribution. . . . .	147
3.39 Desired evolution of the WGAN-SCMA training procedure from the beginning (a) to the objective (d), with $M = 4$ and $d_f = 3$ . Each sub-figure (a)-(d) represents the complex plane related to dimension $k$ . . . . .	156
3.40 Overview of the proposed WGA-SCMA architecture. . . . .	158
3.41 SER comparison. . . . .	162
4.1 High-level depiction of the system model. . . . .	166
4.2 Proposed training procedure in the abstracted observation space. . . . .	170
4.3 The proposed AE-based abstraction framework trading-off compression with value. . . . .	172



4.4 Average total number of successfully delivered dPDUs by the  $N = 2$  agents. . . . . 174

4.5 Average total number of delivered dPDUs under different solutions. Training procedure with  $TBLER = 10^{-4}$  but the performance is evaluated with different values of TBLER. 176

4.6 Average total number of delivered dPDUs under different number of agents and Poisson arrival rate ( $\lambda$ ).  $P = 10$  for each agent. . . . . 177



---

---

## List of Tables

---

1.1	Requirements for IMT-2020 5G . . . . .	3
2.1	System Parameters . . . . .	41
2.2	Parameters of the case study . . . . .	59
2.3	Number of SBs per numerology versus (waveform, GB size) pair. . . . .	61
3.1	Simulation Parameters . . . . .	92
3.2	Simulation Test . . . . .	112
3.3	Simulation Parameters . . . . .	124
3.4	RMSE and $R^2$ values achieved in the estimation of $K_C$ for the considered methods . . . . .	137
3.5	RMSE values achieved in the estimation of $M$ for the considered methods . . . . .	137
3.6	Simulation Parameters . . . . .	144
4.1	Training algorithm Parameters . . . . .	175



---

# CHAPTER *1*

---

## Introduction

---

The objective of this Chapter is to provide a brief *excursus* of the main concepts investigated in this Dissertation, whose aim is to study and propose innovative techniques to optimize the 5G New Radio (NR) interface facing services with very heterogeneous requirements.

### 1.1 Background and Motivation

---

Looking forward to 2030, the key verticals like connected industries, intelligent transport systems and smart cities will bring our society to become totally digitized and data-driven. In this regard, it is expected a fully mobile and connected society, characterized by a huge growth in connectivity and traffic volume. Some typical trends include explosive growth of data traffic, great increase of connected devices and continuous emergence of new and stringent services. Today's statistics show that over 1 billion mobile users around the globe are intensely using the social networking media, streaming, and gaming services on a daily basis. At this regards, 5G technology has to support the proliferating traffic demand, providing a wide range of connected devices and services. Unlike earlier generations, 5G networks do not include just the enhance-

## Chapter 1. Introduction

---

ment for the legacy mobile broad-band services but rather it targets to address much more diversified usage scenarios. In particular, the International Telecommunication Union Radio communication Sector (ITU-R) defined three categories [2]:

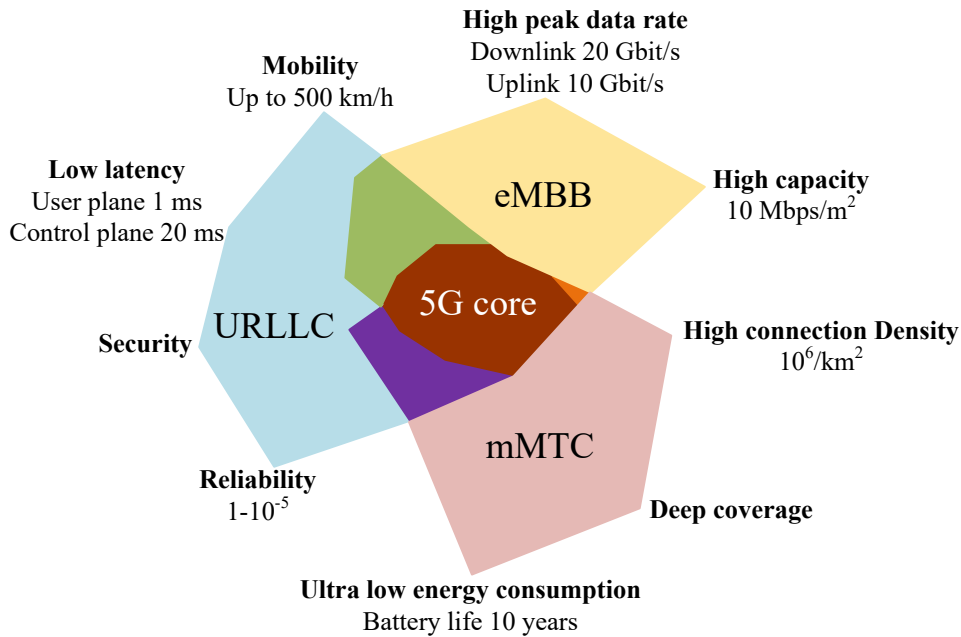
- **eMBB - Enhanced Mobile Broadband:** as an extension to 4G broadband services, the aim of this category is to have an ultra high-speed connection for both indoors and outdoors.
- **mMTC - Massive Machine Type Communications:** 5G will need to suit a whole raft of connected devices with heterogeneous quality of service requirements. The objective of this category is to have a unifying connectivity fabric with enough flexibility to support the exponential increase in the density of connected devices.
- **URLLC - Ultra-Reliable and Low Latency Communications:** this use case is related to services that are delay sensitive, and thus require stringent requirements for latency and reliability to ensure increased reactivity.

The service requirements for each of the above categories are remarkably distinct in terms of reliability, throughput, latency, among others. The first category is data rate hungry (e.g., ultra high-definition (8K) videos at 120 fps and virtual/augmented reality wireless streaming). On the other hand, mMTC envisions the massive Internet of Things (mIoT) paradigm, which requires low power consumption and very low throughput, while URLLC services need to be extremely reliable with a target latency below 1 ms. The requirements for radio access technologies are depicted in Fig. 1.1 and listed in Table 1.1.

According to these objectives, 5G must be able to support users with throughput and peak throughput of 10 and 20 times higher than what available in legacy 4G networks, respectively. The density of the maximum connection will be 10 times more and the energy saving is of importance. In order to cater this vision and to satisfy the proliferating traffic demand, 5G technology is envisioned to support 1000 times increase in capacity and 100 times more connected devices than today's 4G networks.

However, there is a fundamental trade-off among the achievable capacity, latency, and reliability, respectively, over the same spectrum. For

## 1.1. Background and Motivation



**Figure 1.1:** 5G usage scenarios

**Table 1.1:** Requirements for IMT-2020 5G

Capability	5G requirement	Usage scenario
Downlink peak data rate	20 Gbit/s	eMBB
Uplink peak data rate	10 Gbit/s	eMBB
User experienced downlink data rate	100 Mbit/s	eMBB
User experienced uplink data rate	50 Mbit/s	eMBB
User plane Latency	$\leq 4$ ms $\leq 1$ ms	eMBB URLLC
Control plane Latency	$\leq 20$ ms	eMBB/URLLC
Mobility	500 km/h	eMBB/URLLC
Connection density	$10^6/\text{km}^2$	mMTC
Energy efficiency	Equal to 4G	eMBB
Battery life	10 years	mMTC
Area traffic capacity	10 Mbps/m <sup>2</sup>	eMBB
Peak downlink spectrum efficiency	30 bit/s/Hz	eMBB
Reliability	$1-10^{-5}$	URLLC

instance, achieving ultra-reliable wireless transmissions typically require a large radio latency performance though. As the current wireless communication technologies do not have the capabilities to handle the future wireless networks requirements, advanced technologies and intelligent radio resource management techniques have to be developed for 5G and beyond wireless networks. Specifically, the 5G-NR standard introduces

a set of system design improvements [3] mainly highlighted by:

- **Adaptive radio frame structure and numerology:** the 5G-NR [4] supports an agile frame design with a variable Transmission Time Interval (TTI) duration and Sub-Carrier Spacing (SCS), respectively. Hence, a 5G-NR radio frame is 10 ms, and consists of 10 sub-frames, each of 1 ms duration. Sub-frames are flexibly divided into  $2^n$ ,  $n = 0, 1, 2, \dots$ , slots. Accordingly, a radio slot is of 14 Orthogonal Frequency Division Multiplexing (OFDM) symbols duration when a normal cyclic prefix design is used. Within each slot, there can be several transmission opportunities - so called as a mini-slot based transmissions, e.g., transmissions based on 4-OFDM symbol mini-slot. Accordingly, the latency critical URLLC services are dynamically served with a shorter TTI duration (based on the mini-slot duration) and larger SCS, respectively. Alongside the smarter pipeline PHY processing, and the improved processing capabilities available for the 5G era, this fundamentally reduces both the transmission and processing delays, respectively, at the expense of the increased control overhead, due to the shorter transmissions. However, the latency-tolerant eMBB applications are dynamically scheduled with a larger TTI duration to increase the achievable spectral efficiency.
- **Multi-Quality of Service (QoS) Dynamic User Scheduling and Radio Resource Management (RRM) [5–7]:** the state-of-the-art proposals introduce agile RRM and dynamic user scheduling techniques for the 5G-NR multi-QoS deployments. Those contributions typically adopt a multi-objective optimization techniques towards achieving the diverse, and sometimes conflicting, QoS requirements of active UEs. 3GPP release-15 specifications consider the multi-QoS preemptive scheduling as the baseline Medium Access Control (MAC) technique for achieving the stringent radio latency requirements of the latency-critical traffic. It always prioritizes such traffic over other active latency-tolerant traffic by means of immediate pre-emption; however, recovering the capacity of the latter traffic QoS by smarter re-transmission and coding techniques.
- **Advanced channel coding:** in order to deliver higher performance



and efficiency, NR needs a new channel coding using larger coding block sizes. NR specifies an advanced Low-Density Parity-Check (LDPC) for the data channel using a quasi-cyclic structure where a smaller base matrix is used for the parity check matrix. A smaller base matrix means reduced coding latency and complexity as the code rates increase, while also supporting lower rates than LTE turbo codes. NR LDPC therefore provides a full rate compatibility with incremental redundancy and block length flexibility. Physical control channels typically have small block lengths and here Polar codes are proposed in NR. Polar codes are the first to achieve maximum channel capacity, closing the gap to the Shannon limit, and improve performance compared to LTE.

- **Optimized Waveform** [8]: The OFDM waveform has been adopted as the transmission waveform in LTE and has been inherited in 5G NR, thanks to its robustness against fading phenomena and ease of implementation. However, it suffers from several drawbacks, including high Peak-to-Average Power Ratio (PAPR) and high Out-Of-Band (OOB) emissions in the frequency domain. This implies a significant loss in spectrum efficiency. To overcome the above-mentioned limitations, an optimized waveform should be considered for the future stage of 5G. On one hand, the waveform candidate should inherit all the advantages of OFDM. On the other hand, the new waveform should be able to support flexibly configured numerologies and multi-numerology coexistence, to enable tailored services for different applications in heterogeneous scenarios.
- **Innovative Multiple Access Techniques** [9, 10]: Multiple access techniques play a key role in handling data traffic in multi-user systems as it directly determines the throughput performance by efficiently accommodating multiple users with the available resources. Consequently, a standard multiple access scheme is introduced as a specific feature for each generation of wireless network such as Time Division Multiple Access (TDMA) in 2G, Code Division Multiple Access (CDMA) in 3G, and Orthogonal Frequency Division Multiple Access (OFDMA) in 4G. These conventional schemes employ Orthogonal Multiple Access (OMA) techniques in which or-

thogonal resources such as time, frequency and code are assigned to different users to avoid mutual interference between them. Recently, Non-Orthogonal Multiple Access (NOMA) technique has been envisioned as a promising multiple access candidate to address these high data rate requirements as well as to support the massive connectivity in 5G and beyond networks. In contrast to OMA, NOMA can simultaneously allocate an available radio resource to more than one user which significantly enhances the system throughput due to frequency reuse within a cell. The available NOMA techniques can broadly be divided into two categories, namely, power-domain and code-domain NOMA:

1. **Power-domain NOMA** [11]: In this approach, different users are served at different transmit power levels according to their channel conditions to obtain the maximum gain in system performance. Therefore, the user with lower channel gain is served with higher transmit power whilst less transmit power is allocated to the user with high channel gain. To carry out this power domain multiplexing, the Base Station (BS) transmits a linear superposition of the signals of the users and at the receiver sides, multiuser detection algorithms such as Successive Interference Cancellation (SIC) are utilized to detect the desired signals.
2. **Code-domain NOMA** [12]: Unlike power-domain NOMA, which attains multiplexing in power domain, code-domain NOMA achieves multiplexing in code domain. Similar to the basic CDMA systems, code-domain NOMA shares the entire available resources in time and frequency. In contrast, code-domain NOMA utilizes user-specific spreading sequences that are either sparse sequences or nonorthogonal cross-correlation sequences of low correlation coefficient. Examples of code-domain NOMA are Low-Density Spreading-based CDMA (LDS-CDMA), Low-Density Spreading-based OFDM (LDS-OFDM), and Sparse Code Multiple access (SCMA).

## 1.2 Research Questions

---

In order to reach the ambitious requirements of 5G, research in diverse fields is required [13]. Therefore, in this dissertation, we make the efforts to answer the following research questions.

- Question 1: How to develop a flexible RRM framework that deals with multi-objective optimization techniques to achieve conflicting QoS requirements of UEs' services?
- Question 2: How to near-optimally predict traffic arrival at a gNB and use such information to boost the number of satisfied MTC devices and to reduce RA congestion during a bursty traffic arrival?
- Question 3: How to speed up the adoption of SCMA in practical networks?
- Question 4: How to efficiently learn MAC protocols in an automatic fashion without any prior agreement among the UEs?

## 1.3 Research Contributions and Thesis outline

---

To address the above mentioned research questions, the general objective of these Ph.D. research activities is to optimize the coexistence in the same OFDM grid among different type of services, to enhance the performances of the resource allocation for the mMTC scenario, and to automatically generate robust and generalized MAC protocols in the vision of Beyond 5G (B5G).

For each research activity, we carried out several studies and a critical analysis of the state of art. The objectives to be achieved had been clearly reported and mathematically formalized. When the analytical solution of the formulated problem is computationally high, we provided heuristic proposals or approaches based on Artificial Intelligence (AI) for the given problem. The goodness of the proposed solutions has been verified through a large number of simulations in comparison with other works available in the literature.

The main contributions of this Ph.D. Dissertation are summarized as follows.

- The first part considers the coexistence among a wide variety of services inside the same OFDM grid with the aim of providing flexibility. This flexibility brings with it a further challenge in RRM, i.e., how to best allocate the available band spectrum among the different non-orthogonal numerologies. Therefore, we propose a two-levels RRM framework that allocates the band spectrum to the numerologies on the basis of QoS requested by the diversified type of service, i.e., Best Effort (BE) services or Guaranteed Bit Rate (GBR) services with different priorities, and the channel condition experienced by the UE requiring that service. The goal is to maximize the number of satisfied GBR services, according to the priority, and the throughput of the BE services. Moreover, we implement a new Simulation Environment constituted of the defined framework and a physical level simulator publicly available. The simulator, named Vienna 5G Link Level (LL) Simulator, supports a single numerology scenario where some physical features can be set up, including the waveform (e.g., Orthogonal Frequency Division Modulation Universal Filtered Multi-Carrier, and filtered-OFDM), the channel model (e.g., Additive White Gaussian Noise, Pedestrian-A), and the sub-carrier spacing. We adapt this simulator to support the proposed multi-numerology scenario and variable Guard Band (GB) sizes between adjacent different numerologies. We identify the best (waveform, GB size) pair which reduces the INI phenomenon, maximizes the spectral efficiency while taking into account QoS requirements. The results of these activities have been published in the conference works [1, 14], and in the Computer Communications journal [15].
- In the second part, we focus on the optimization of the RRM for the mMTC services. We propose a new framework, customized for mMTC services, which includes a joint control of the dynamic resource allocation between the Physical Random Access Channel (PRACH) and to the Physical Uplink Shared Channel (PUSCH), and a new random access procedure based on an adaptive Access Class Barring (ACB) scheme that appropriately spreads random access re-attempts in time. In addition, to further increase the transmission efficiency, we adopt the SCMA technique for PUSCH re-

### 1.3. Research Contributions and Thesis outline

---

sources, because SCMA results as the most promising NOMA technique to support massive MTC connectivity with small-size data. The results are published in IEEE Internet of Things Journal [16]. Then, instead of improving the succeeded access attempts in the PRACH, as widely addressed in literature, we present the innovative idea to exploit the unused PUSCH resources to serve also the MTC devices that have failed their access attempt, by assigning them a pool of resources, in a contention-based mode. This study led to the publication in IEEE Communications Letters [17]. Moreover, all the schemes for the mMTC are mainly based on a grant-based Random Access (RA) procedure and proper load-aware access controls, e.g., ACB schemes, dynamic uplink radio resource allocation, and so on. therefore, the development of an efficient approach to estimate the traffic load is extremely important for the proper functioning of these access schemes. Therefore, we propose a current access attempts estimation, based on Deep Neural Network (DNN), which accepts as input only the information really available at the next generation Node B (gNB). Finally, to improve the adoption of SCMA in practical networks (especially mMTC scenario) we design an end-to-end SCMA en/decoding structure based on the integration between a state-of-the-art autoencoder architecture and a novel Wasserstein Generative Adversarial Network (WGAN) that improves the robustness to the channel noise. These two last studies have been published in Computer Networks journal [18] and in IEEE Wireless Communications Letters [19], respectively.

- In the third part, we address the problem of automatically learning innovative Medium Access Control (MAC) protocols catering to extremely diverse services with much better generalization capabilities. To do so, UEs are cast as Reinforcement Learning (RL) agents learning based on local observations while the Base Stations (BSs) are cast as an expert. To provide robustness, we leverage the concept of Observation Abstraction (OA) rooted in extracting useful information from the environment. The results have been accepted for publication in IEEE GLOBECOM 2022 [20].



---

**A QoS-aware and Channel-aware  
framework for Multi-numerology  
OFDM/UFMC/filtered OFDM Systems**

---

**2.1 Overview**

---

As presented in the Chapter 1, 5G networks are expected to provide wireless connectivity for a wide range of new applications, with the ability to be connected to the Internet anytime, anywhere, allowing access seamlessly, with any device [21]. Since the LTE resource grid has limited flexibility, it cannot support these heterogeneous requirements. Hence, to overcome this issue, the physical layer of the 5G cellular network, called 5G NR, introduces new concepts and building blocks to provide more flexible radio access technology [22]. This flexibility is provided mainly by the coexistence of multi-numerologies [23], where each numerology corresponds to one sub-carrier spacing  $\Delta f$  in the frequency domain. Consequently, in the OFDMA based 5G radio interface the time duration of a time slot scales with the chosen subcarrier spacing. Although this approach is useful to support a wide variety of services, it introduces new challenges for effectively and efficiently develop a RRM scheme that ensures the QoS requirements and maximizes system throughput.

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

In fact, in addition to how to correctly schedule the Physical Resource Blocks (PRBs) to the UEs, it is necessary to consider how to best divide the spectral resources among the different numerologies.

Generally speaking, the RRM scheme designed for multi-numerology scenarios aims to support all services belonging to different numerologies, guaranteeing QoS requirements, maximizing the system throughput and providing fairness among UEs that request the same type of service. However, these objectives are often in contrast to each other and, thus, optimal solutions are difficult to achieve and/or not feasible because they can require excessively long computation time. Then, to find sub-optimal but feasible solutions, the RRM problem could be approached by splitting the complex procedure into two simpler sub-procedures, each one managed by a distinct level of control. The first one, hereinafter referred to as *1st level allocation*, consists in a proper subdivision of the spectrum among the various numerologies. In the second level, *2nd level scheduling*, the PRBs allotted to each numerology are appropriately scheduled to the relative UEs. More specifically, we focus on the 1st level of the control framework and we propose a QoS-aware and channel-aware heuristic algorithm that does not require a long computation time. The aim of our control level was to determine "how many" and "which" spectrum portions have to be allocated to each numerology with the targets of satisfying the QoS requirements, in terms of priority and minimum Guaranteed Bit Rate (GBR), and to maximize the system throughput. To achieve this goals, our 1st level algorithm, called Channel-Aware Resource Allocation for Multi-numerology (CARAM) [1], dynamically splits the available bandwidth among the numerologies, taking into account the channel conditions. In addition, in this control level we integrated a simple dropping strategy that was activated in critical situations of overload. As regards the 2nd level control, for each numerology on the basis of the allocated bandwidth portions, it concludes the RRM process by distributing the PRBs among the UEs, by means of the well-known Proportional Fair (PF) algorithm [24], available in literature. In order to more efficiently address the traffic overload conditions and the time-varying channel quality, and to increase the flexibility in the adoption of different 2nd level algorithms, the core of the proposed framework (CARAM), has been completely revised (CARAM-new). Specifically, the CARAM-new



framework more efficiently addresses the traffic overload conditions and the time-varying channel quality. In addition, we increase its flexibility by including the adoption of different 2nd level algorithms, such as PF, Best Channel Quality Indicator (BCQI) and Highest Deviation (BCQI-HD) [25], that are all described in Section 2.7.

However, the main weakness of the proposed framework is that it neglects the Inter-Numerology Interference (INI), created by those adjacent radio resources that are no longer orthogonal as they are assigned to different numerologies. Moreover, it does not introduce any GB, only supports the conventional OFDM waveform, and its performance has only been analyzed with a simplistic Additive White Gaussian Noise (AWGN) channel. To reduce the INI effects some proposals are available in the literature, such as the introduction of a proper Guard Band (GB) between different numerologies, and the adoption of innovative Orthogonal Frequency Division Multiplexing (OFDM)-based waveforms, e.g., filtered OFDM (f-OFDM) [26], Universal Filtered Multi-Carrier (UFMC) [27] and windowed OFDM [28]. The optimal solution to the overall RRM problem is difficult to analytically derive or not viable as it requires excessively long calculation time. At this regard, powerful numeric simulators have been developed as tools for investigating and analyzing future wireless technologies. One of the most promising simulators is the Vienna 5G Link Level (LL) [29], a MATLAB-based link-level simulation tool that facilitates research and development in mobile communications. Therefore, taking into account the new and sophisticated RRM schemes and the introduction of advanced simulators, in this Chapter we additionally design a new Simulation Environment that emulates the complex multi-numerology scenario, including the INI phenomenon. In detail, as regards the RRM procedures, we adopt the CARAM-new framework and implement another second level scheduler [25]. Then, we appropriately integrate these control schemes to the Vienna 5G simulator which allows to modify several physical features, including the adopted waveform, the channel model (e.g., AWGN, Pedestrian-A), and the sub-carrier spacing. Finally, starting from the designed Simulation Environment, we carried out a case study that allowed us to emulate a multi-numerology scenario with specific and diversified requirements and evaluate what is the best (waveform, GB size) pair that reduces the

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

INI phenomenon and maximizes the spectral efficiency, taking into account the QoS requirements.

### 2.2 5G Downlink Frame Structure and Assumptions

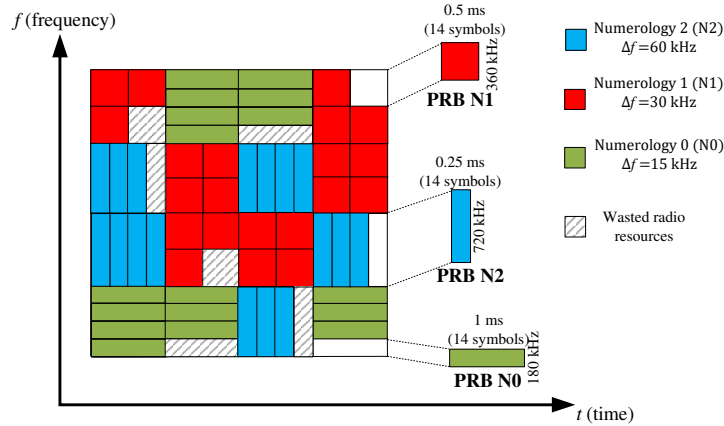
---

At the physical layer, 5G NR allows flexible bandwidth with a maximum channel bandwidth per carrier equal to 400MHz [4]. As in the LTE, in the 5G NR the OFDMA is used for the downlink transmission, where users are dynamically multiplexed on a time-frequency grid. Moreover, to meet the different requirements of the 5G usage scenarios, the support of different OFDM numerologies is introduced. The different OFDM numerologies are characterized by sub-carrier spacing  $\Delta f$ , ranging from 15 KHz to 480 KHz, calculated as follows:  $\Delta f = 15 \cdot 2^x, x \in \mathbf{X}_T = \{0, 1, \dots, 5\}$ . The correct numerology that should be used depends on the usage scenario and on UEs' requirements [30]. The radio resources are grouped into Physical Resource Blocks (PRBs), where one PRB spans 12 sub-carriers for the time of one time slot, corresponding to 14 OFDM symbols (normal cyclic prefix length).

In Fig. 2.1 we show an example reporting the first three numerologies of set  $\mathbf{X}_T$ . As depicted, there are sub-carrier spacing  $\Delta f$  of 15, 30 and 60 KHz, leading to PRB bandwidths of 180, 360, and 720 KHz, and time slots duration of 1, 0.5, and 0.25 ms, respectively. In the time domain, every TTI, the gNB shares-out the available band to the numerologies and schedules the PRBs to the related UEs. Differently from LTE, where a TTI is fixed to 1 ms, the 5G NR supports the adoption of a scalable TTI which can be shorter than 1 ms.

Since adjacent radio resources assigned to different numerologies are not orthogonal to each other, they cause an interference phenomenon known as INI that reduces the quality of communication experienced by users, and thus lowering the spectral efficiency. To reduce this phenomenon several proposals are available in literature, such as the introduction of a guard band between adjacent resources belonging to different numerologies [31], the adoption of an alternate waveform instead of OFDM (e.g., filtered OFDM [26] and windowed OFDM [28]), and/or a power control offset among different numerologies [32]. In this paper, the focus is not to choose the best solution that guarantees a tolerable

## 2.3. System Model and Problem Formulation



**Figure 2.1:** Multi-Numerology Resource Grid.

INI, but, in order to minimize the correlated effects, our strategy aims to reduce the number of areas affected by the INI trying to assign contiguous resources to the same numerology.

As regards the reporting of CQI values by the UEs, the procedures envisaged by the standard are considered [33]. In particular, we adopt the *higher layer configured sub-band reporting* method, where the gNB divides the downlink band into equally sized sub-bands, and reports a wideband CQI value for the whole system bandwidth, together with a CQI value for each sub-band.

The gNB, on basis of the CQI values received by a UE, begins the downlink transmission procedure by sending to it a scheduling command containing, inter alia, the Resource Block Allocation Type, the Resource Block Assignment, and the Modulation and Code Scheme (MCS) Index. The gNB adopts a proper mapping between the received CQI and the adopted MCS index, and an example is reported in [34]. Then, the transmission rate achievable by each UE depends on the MCS adopted for it in each PRB allocated to it.

## 2.3 System Model and Problem Formulation

We consider a downlink 5G NR network where the adopted resource grid of bandwidth  $BW$  allows both the dynamic adaptation of the numerology and the coexistence of multiple numerologies in a frame. There is a single gNB equipped with one transmit antenna, and  $N_{UE}$  UEs equipped with one receive antenna. Each numerology  $x$  belongs to  $\mathbf{X} \subseteq \mathbf{X}_T$ , and

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

we denote with  $N_x$  the cardinality of  $\mathbf{X}$ , and with  $\mathbf{X}[i]$  the  $i$ th element of  $\mathbf{X}$ . It is worth to underline that the maximum value of  $N_x$  is 3, because at most three numerology can be multiplexed over a band according to the 3GPP standardization [35]. Each numerology  $x$  is used by the gNB to manage different type of services, either a BE service or a GBR service with a proper priority. The GBR services belonging to the same numerology have the same priority, but may have a different GBR requirement. We define  $\mathbf{P} = [p_1, p_2, \dots, p_{N_x}]$  as a vector whose element  $p_i$  represents the priority of the services related to the numerology  $\mathbf{X}[i]$ .

We assume that each UE  $i$  requires only one BE or GBR service<sup>1</sup> starting from TTI  $l_{start,i} \in \{1, \dots, N_{TTI} - 1\}$ , where  $N_{TTI}$  the simulation time in terms of number of TTIs. Without loss of generality, in the following we use UE and service indifferently. Let  $\mathbf{U} = [1, \dots, N_{UE}]^T$  be the vector containing all UEs requesting a service in the entire simulation time arranged according to the order of arrival of requests, i.e.,  $l_{start,\mathbf{U}[i]} \leq l_{start,\mathbf{U}[i+1]}, \forall i = 1, \dots, N_{UE} - 1$ . We define also  $\mathbf{U}_x$  as a sub-vector of  $\mathbf{U}$  containing all UEs requiring a service belonging to numerology  $x, \forall x \in \mathbf{X}$ , arranged in the same order of  $\mathbf{U}$ . Clearly,

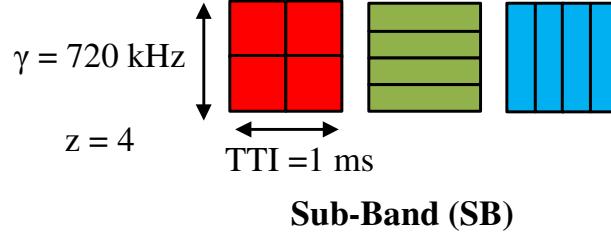
$$\mathbf{U} = \bigcup_{x \in \mathbf{X}} \mathbf{U}_x. \quad (2.1)$$

Moreover,  $\mathbf{R}_x = [R_{1_x}, R_{2_x}, \dots, R_{|\mathbf{U}_x|_x}]$  is a vector containing the GBR requirements of the UEs related to numerology  $x$ ,  $\mathbf{X}_{\text{GBR}}$  as a sub-vector of  $\mathbf{X}$  containing all numerologies  $x$  so that  $R_{i_x} > 0 \forall i \in \{1, 2, \dots, N_{UE}\}$  sorted in descending order of priority, and  $\mathbf{X}_{\text{BE}}$  as the relative complement vector of  $\mathbf{X}_{\text{GBR}}$  with respect to  $\mathbf{X}$ . We assume that  $|\mathbf{X}_{\text{BE}}|=1$ . For simplicity, we term the numerologies in  $\mathbf{X}_{\text{GBR}}$  as *GBR numerologies* and the one in  $\mathbf{X}_{\text{BE}}$  as *Best Effort numerology*.

For each TTI, in the 1st control level, the gNB should appropriately assign the radio resources to each numerology. To efficiently allocate the time-frequency resource grid without waste, it is first necessary to determine the granularity of the minimum amount of resources that can be assigned to a numerology. In the frequency domain, one Sub-Band (SB) of size  $\gamma$  corresponds to the minimum common multiple (mcm) among

<sup>1</sup>We note that the system model continues to be valid even in the case of one or more users requesting to transmit  $N$  different services. In this case, this user could be modeled as  $N$  users experiencing the same CQIs for the entire simulation length, each one requesting to transmit one out of the  $N$  services.

### 2.3. System Model and Problem Formulation



**Figure 2.2:** 1 Sub-Band (SB), i.e., granularity of the minimum amount of resources that can be assigned to each numerology.

the PRB bandwidths of the numerologies in  $\mathbf{X}$ . In the time domain, the minimum TTI is the mcm among the different time slot lengths. Therefore, the total number of available sub-bands is  $N_{SB} = \left\lfloor \frac{BW}{\gamma} \right\rfloor$ , and we define  $\mathbf{K}$  as the following vector  $\mathbf{K} = [1, \dots, N_{SB}]$ . In the example shown in Fig. 2.1, it follows  $\gamma = 720$  KHz, the minimum TTI lasts 1 ms, and consequently, one SB corresponds to  $z = 4$  PRBs for each numerology, as depicted in Fig. 2.2. Then, we define  $\mathbf{L}$  as the vector containing all the available PRBs, i.e.,  $|\mathbf{L}| = z \cdot |\mathbf{K}|$ .

For each TTI  $l$ , we denote with  $\mathbf{U}^{(l)}$  the sub-vector of  $\mathbf{U}$  containing the UEs requiring a generic service starting from TTI  $l_{start,i} \leq l$ . We assume that, the UEs in  $\mathbf{U}^{(l)}$  report to the gNB both the wideband CQI and the CQI values for each sub-band in  $\mathbf{K}$ . Let  $\mathbf{c}_i^{(l)} = [CQI_{i,1}^{(l)}, \dots, CQI_{i,N_{SB}}^{(l)}]$  be the vector containing the sub-band CQI values reported from UE  $i$  in TTI  $l$ , with  $i$  in  $\mathbf{U}^{(l)}$ .  $\mathbf{C}^{(l)} = [\mathbf{c}_1^{(l)}, \dots, \mathbf{c}_{|\mathbf{U}^{(l)}}^{(l)}]^T$  is the matrix containing all CQI values received by the gNB during TTI  $l$ , and  $\mathbf{w}^{(l)} = [wCQI_1^{(l)}, \dots, wCQI_{|\mathbf{U}^{(l)}}^{(l)}]^T$  the column vector containing the wideband CQI values reported from each UE in  $\mathbf{U}^{(l)}$  in the TTI  $l$ . Let  $\mathbf{A}^{(l)}$  be the 1st level Allocation Matrix of size  $N_x \times N_{SB}$  at TTI  $l$  whose binary elements  $a_{x,k}^{(l)}$  are indicators of the SBs assigned to the numerologies in TTI  $l$ :  $a_{x,k}^{(l)} = 1$  if SB  $k$  is assigned to the numerology  $x$ , otherwise  $a_{x,k}^{(l)} = 0$ . Therefore, the purpose of the 1st level control, for each TTI  $l$ , is to derive the proper values of  $a_{x,k}^{(l)}$  that maximizes the system throughput, while ensuring the QoS requirements.

Moreover, for each  $x \in \mathbf{X}$ , let  $\mathbf{U}_x^{(l)}$  be the sub-vector of  $\mathbf{U}^{(l)}$  containing the UEs requiring a service belonging to  $x$ ,  $\mathbf{K}_x^{(l)}$  be the sub-vector of  $\mathbf{K}$  containing all the SBs  $k$  so that  $a_{x,k}^{(l)} = 1$ ,  $\mathbf{L}_x^{(l)}$  a sub-vector of  $\mathbf{L}$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

containing all the PRBs related to the SB  $\in \mathbf{K}_x^{(l)}$ . More in detail, each element  $\mathbf{K}_x^{(l)}[i]$  corresponds to  $z$  elements of  $\mathbf{L}_x^{(l)}$ , i.e.,  $\mathbf{L}_x^{(l)}[(i-1) \cdot z + 1] = (\mathbf{K}_x^{(l)}[i] - 1) \cdot z + 1$ ,  $\mathbf{L}_x^{(l)}[(i-1) \cdot z + 2] = (\mathbf{K}_x^{(l)}[i] - 1) \cdot z + 2$ ,  $\dots$ , and  $\mathbf{L}_x^{(l)}[i \cdot z] = \mathbf{K}_x^{(l)}[i] \cdot z$ . We assume that for the UE  $i$  the CQI value reported on the sub-band  $k$ ,  $\mathbf{C}^{(l)}[i, k]$ , is constant for each PRB in the sub-band. Therefore, we define  $\mathbf{C}_x^{(l)}$  as a matrix of size  $|\mathbf{U}_x^{(l)}| \times |\mathbf{L}_x^{(l)}|$  containing all the CQI values  $\mathbf{C}_x^{(l)}[i, j]$  of the UEs of numerology  $x$ , related to each of the PRB  $j \in \mathbf{L}_x^{(l)}$ . It is clear that  $|\mathbf{L}_x^{(l)}| = |\mathbf{K}_x^{(l)}| \cdot z$ , for each TTI  $l$  and that the elements of  $\mathbf{C}_x^{(l)}$  are derived from those of  $\mathbf{C}^{(l)}$ , considering only the sub-bands  $k$  where  $a_{x,k}^{(l)} = 1$ . Then, the 2nd level control schedules the PRBs  $\in \mathbf{L}_x^{(l)}$  to UEs in  $\mathbf{U}_x^{(l)}$ . So, we define  $\mathbf{A}_x^{(l)}$  as the 2nd level Scheduling Matrix of size  $|\mathbf{U}_x^{(l)}| \times |\mathbf{L}_x^{(l)}|$ , whose binary elements:  $a_{i_x, j_x}^{(l)} = 1$  if PRB  $j_x$  is assigned to UE  $i_x$ , otherwise  $a_{i_x, j_x}^{(l)} = 0$ .

Given  $\mathbf{A}_x^{(l)}$ ,  $\mathbf{U}_x^{(l)}$ , and  $\mathbf{C}^{(l)}$ , the gNB knows the CQI values of the PRBs allocated to each user  $i_x \in \mathbf{U}_x^{(l)}$ , then determines which MCS it should use in the Physical Downlink Shared Channel (PDSCH). At this aim, we adopt the CQI-MCS Index mapping Table A.4-3 in [34]. The modulation order and code rate related to each MCS Index is derived by means of Table 5.1.3.1-1 in [33], which support QPSK, 16-QAM and 64-QAM. Then, the Transport Block Size ( $TBS_{i_x}^{(l)}$ ) related to UE  $i_x$  in TTI  $l$ , that is the number of information bits transmitted during TTI  $l$ , is estimated on basis of the number of allocated PRBs and on the MCS adopted for each PRBs, as reported in [33]. In order to verify the GBR requirement, we consider the standardized Default Averaging Window (DAW) [36] with  $\mathbf{W}$  corresponding to the time interval of the last  $N_{DAW}$  TTIs. For each TTI  $l$  and UE  $i_x$ , the related DAW is  $\mathbf{W}_{i_x}^{(l)} = [\max\{l_{start, i_x}, l - N_{DAW} + l_{start, i_x}\}, \dots, l]$ . Thus, the verification of the GBR requirement for each UE  $i_x$  exploits the following transmission rate:

$$T_{i_x, \mathbf{W}}^{(l)} = \sum_{t \in \mathbf{W}_{i_x}^{(l)}} \frac{TBS_{i_x}^{(l)}}{TTI \cdot N_{DAW}}, \text{ with } l \geq l_{start, i_x} \quad (2.2)$$

Clearly, the  $N_{DAW}$  value may be different according to the service requirement. In our work it was chosen in accordance with [36]. For

### 2.3. System Model and Problem Formulation

each numerology  $x \in \mathbf{X}_{\text{GBR}}$ , let  $\mathbf{S}_x$  be the user satisfaction vector of dimension  $|\mathbf{U}_x|$  with binary elements  $s_{i_x}$ . Since the bit rate should be guaranteed by the network over a sliding DAW, the Boolean value  $s_{i_x}$  is set as follows:

$$s_{i_x} = \begin{cases} 1, & \text{if } T_{i_x}^{(l)} \geq R_{i_x} \forall l \in \{N_{\text{DAW}} + l_{\text{start}, i_x}, \\ & \dots, N_{\text{TTI}}\} \\ 0 & \text{otherwise.} \end{cases} \quad (2.3)$$

The first goal is to maximize the number of satisfied UEs taking into account the priority of the required services. Hence, it is necessary to maximize the number of satisfied UEs belonging to  $|\mathbf{U}_x|$ , with  $x$  equal to  $\mathbf{X}_{\text{GBR}}[1]$ , i.e., the users requiring the service with the maximum priority. If all these UEs are satisfied, the target is to maximize the number of satisfied UEs in  $|\mathbf{U}_x|$ , with  $x$  equal  $\mathbf{X}_{\text{GBR}}[2]$ , and so on. This goal can be mathematically expressed as follows:

$$\max \left\{ \sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[1]}|} s_{i_{\mathbf{X}_{\text{GBR}}[1]}} + \sum_{x=2}^{|\mathbf{X}_{\text{GBR}}|} \sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[x]}|} s_{i_{\mathbf{X}_{\text{GBR}}[x]}} \cdot \prod_{y=1}^{x-1} \left[ \frac{\sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[y]}|} s_{i_{\mathbf{X}_{\text{GBR}}[y]}}}{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[y]}|} \right] \right\}. \quad (2.4)$$

The second goal of maximizing the Best Effort Numerology throughput can be expressed as follows:

$$\max \left\{ \sum_{i_x \in \mathbf{U}_{\text{XBE}}} T_{i_x} \right\}, \quad (2.5)$$

where  $T_i$  is the throughput experienced by the  $i$ th UE in the entire simulation length. It is calculated as

$$T_{i_x} = \frac{1}{N_{\text{TTI}}} \sum_{l=1}^{N_{\text{TTI}}} T_{i_x}^{(l)}, \quad (2.6)$$

where

$$T_{i_x}^{(l)} = \frac{TBS_{i_x}^{(l)}}{TTI}. \quad (2.7)$$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

Furthermore, the third goal is to maximize the fairness in throughput among BE UEs. At this aim, we consider the Jain index [37]:

$$\max \left\{ \frac{\left( \sum_{i_x \in \mathbf{U}_{\mathbf{X}_{\text{BE}}}} T_{i_x} \right)^2}{|\mathbf{U}_{\mathbf{X}_{\text{BE}}}| \sum_{i_x \in \mathbf{U}_{\mathbf{X}_{\text{BE}}}} T_{i_x}^2} \right\}. \quad (2.8)$$

It is evident that the objective functions (2.4), (2.5), and (2.8) are referred to the results obtained after applying the 2nd level scheduling algorithms.

With regards to the constraints, since we consider an Orthogonal Multiple Allocation (OMA) among different numerologies, each SB  $k \in \mathbf{K}$  must be allocated to one numerology at most. The first constraint can be formulated as:

$$\sum_{x=1}^{N_x} a_{x,k}^{(l)} \leq 1, \forall k \in \mathbf{K}, \forall l \in \{1, N_{TTI}\} \quad (2.9)$$

Second, for each numerology  $x \in \mathbf{X}$ , the PRB  $j_x$  can be assigned, by the 2nd level scheduler, to one UE at most. This can be expressed as:

$$\sum_{i=1}^{|\mathbf{U}_x^{(l)}|} a_{i_x, j_x}^{(l)} \leq 1, \forall j_x \in \mathbf{L}_x^{(l)}, \forall l \in \{1, N_{TTI}\} \quad (2.10)$$

In summary, the problem formulation for the optimal overall RRM, which provides the sub-band allotted for the different numerologies and the resource scheduling for the UEs, can be written as follows:

"Obtain the matrix  $\mathbf{A}$  and the matrices  $\mathbf{A}_x, \forall x \in \mathbf{X}$  with the best trade off among the objective functions (2.4), (2.5) and (2.8), subject to the constraints (2.9) and (2.10)."

This is a NP-hard problem and an analytical optimization solution will take significantly long computation time, which is unacceptable for allocation and scheduling procedures which should be made every TTI. Therefore, our approach consists in deriving a heuristic algorithm that does not require a long computation time and obtain sub-optimal solutions.

## 2.4 Benchmark Schemes

---

In this Section, we present the benchmark schemes considered in this paper. Other related works are reported in Section 2.7.



### 2.4.1 1st level allocation algorithms

In this subsection, we focus on the "1st level algorithms". The most straightforward solution, called Constant Band (CB) [30], is a channel-unaware algorithm that divides the radio spectrum among the numerologies, so that each of them obtains a number of sub-bands proportional to the number of UEs, regardless of the service priority and the QoS requirement. According to our system model, for each numerology  $x \in \mathbf{X}$ , the number of sub-bands allocated to  $x$  is  $N_{SB_x} = \left\lceil \frac{N_{SB}}{|\mathbf{U}|} |\mathbf{U}_x| \right\rceil$ . This is a very simple procedure which guarantees the lowest decision-making times. However, this solution offers very poor performance in terms of both system throughput and number of satisfied users. In fact, the same amount of resources is assigned to numerologies that could manage users with different CQI and GBR values.

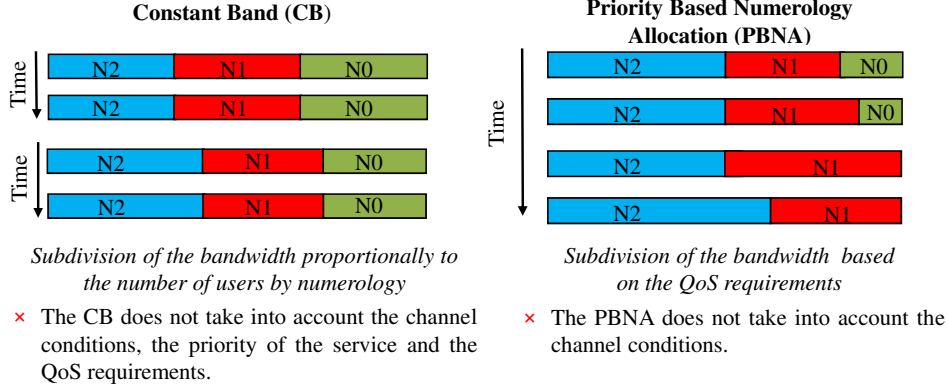
The Priority Based Numerology Allocation (PBNA) [30], unlike CB, takes into account the numerology priorities  $\mathbf{P}$  and the GBR requirements  $\mathbf{R}_x, \forall x \in \mathbf{X}$ , but the channel conditions are still not considered. As a first step, the algorithm assigns the resources to the highest-priority numerology in order to satisfy all services belonging to the first numerology. Then, if there are still resources available, the algorithm allocated the sub-bands needed to satisfy all services belonging to the second-best priority numerology. The procedure is iterated in the priority order until all the  $N_{SB}$  sub-bands are assigned or there are no more requests. Clearly, PBNA performs better than the CB in terms of satisfied UEs. However, as the CB, it does not take into account the UEs' CQI values, so it is an inefficient solution when the channel conditions are different among UEs.

For the sake of clarity, we summarize the main characteristics of the above 1st level schedulers in Fig. 2.3.

### 2.4.2 2nd level scheduling algorithms

Among the channel-aware 2nd level algorithms, some of the conventional well-known ones are BCQI, and PF [24]. The BCQI scheduler, every TTI  $l$ , assigns each resource to the user who has the maximum CQI value on it. According to our system model, it follows:

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.3:** 1st level allocation algorithms at a glance.

$$a_{i_x, j_x}^{(l)} = \begin{cases} 1, & \text{if } i = \arg \max_{i \in \mathbf{U}_x^{(l)}, j \in \mathbf{L}_x^{(l)}} \{ \mathbf{C}_x^{(l)}[i, j] \} \\ 0 & \text{otherwise.} \end{cases} \quad (2.11)$$

This is a very simple procedure, which maximizes the numerology throughput. However, it is highly unfair as users far away from the gNB with poor channel conditions do not get any resources.

The PF scheduler [24], like BCQI, does not take into account the GBR requirements but only the UEs channel conditions. It aims to find a good tradeoff between the system throughput and fairness, assigning every TTI  $l$  all  $N_{SB_x}$  SBs to the user  $i'_x$  calculated as follows:

$$i'_x = \arg \max \frac{T_{i_x}^{(l)}}{\bar{T}_{i_x}^{(l)}}, \quad (2.12)$$

where  $\bar{T}_{i_x}^{(l)}$  is recursively calculated as

$$\bar{T}_{i_x}^{(l)} = \beta \bar{T}_{i_x}^{(l-1)} + (1 - \beta) T_{i_x}^{(l-1)}, \quad (2.13)$$

and  $0 \leq \beta \leq 1$ . By using this strategy, users with good channel conditions obtain a large number of PRBs, while those far away from the gNB, generally with lower CQI, still get some resources. However, it is not suitable for real-time services and does not take into account GBR requirements.

In addition to the conventional 2nd level schedulers, we consider the BCQI Highest Deviation (BCQI\_HD) proposed in [25]. It is a QoS and

channel-aware scheduler that aims to achieve a good trade-off between maximizing the system throughput and improving the fairness among UEs, while ensuring the minimum data rate for GBR services. The entire procedure can be summarized as follows.

1. The scheduler initially assigns one PRB to each UE. The choice of PRB to allocate is carried out by taking into account the complete view of the channel conditions perceived by all UEs. In fact, before allocating one PRB to a UE, the scheduler verifies whether this choice does severely penalize any other UE. For doing so, the scheduler calculates, for each UE, a deviation value between the maximum and the second maximum CQI value, and assigns one PRB to the UE with the greatest value of deviation, that is, the UE which presents the highest deviation from its maximum reachable throughput.
2. Afterwards, it makes an estimate of the expected transmission rate in the considered TTI, i.e.,  $T_{i_x}^{(l)}$  for each user  $\forall i_x \in \mathbf{U}_x$ . If the  $T_{i_x}^{(l)} < R_{i_x}$  for at least one user  $i_x$ , the scheduler selects all users that are not satisfied and repeats the procedure 1) until  $T_{i_x}^{(l)} \geq R_{i_x}$ ,  $\forall i_x \in \mathbf{U}_x$ .
3. Finally, if there are still resources available and the UEs have still data to transmit, the scheduler assigns the PRBs cyclically to UEs, regardless of the channel conditions, until there are no more available PRBs.

This strategy guarantees better performance than the previous algorithms [25], but still has some drawbacks. Specifically, the BCQI\_HD does not take into account the average throughput achieved in the DAW by UEs, but aims only to guaranteeing the GBR requirement inside the single TTI  $t$ , without considering the past scheduling allocations, i.e.,  $T_{i_x,l}$ . For this reason, the integration between our 1st level algorithm presented in the next Section and the BCQI\_HD scheduler is not straightforward. In the Section 2.5.1, we present the strategy adopted to solve this issue.

## 2.5 The Proposed Control Framework

---

In this section, we present the proposed radio resource framework composed of two control levels, as shown in Fig. 2.4. As stated previously, the new proposal focuses on the 1st level allocation algorithm, termed CARAM-new. It is a QoS-aware and channel-aware algorithm that aims to maximize the BE throughput and the number of satisfied GBR users, according to the priority.

For each TTI  $l$ , the inputs are the CQI matrix  $\mathbf{C}^{(l)}$ , the vectors  $\mathbf{X}_{\text{GBR}}$ ,  $\mathbf{X}_{\text{BE}}$ ,  $\mathbf{U}_{\mathbf{x}}^{(l)}$ , and  $\mathbf{R}$ . The output is the appropriate allocation matrix  $\mathbf{A}^{(l)}$ . As depicted in Fig. 2.4, on the basis of  $\mathbf{A}^{(l)}$ , for each numerology  $x$ , a 2nd level scheduling algorithm assigns proper PRBs to each UE in  $\mathbf{U}_{\mathbf{x}}$ , giving as output the  $\mathbf{A}_{\mathbf{x}}^{(l)}$  matrix. Moreover, we assume a feedback reporting between the 1st and the 2nd control level with a periodicity of  $\Delta T$ .

The basic principle of CARAM-new is divided into two phase:

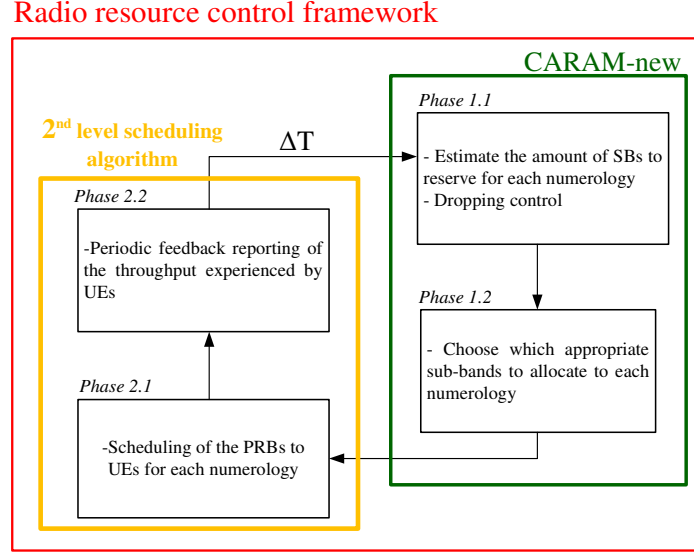
- Phase 1.1 evaluates *how many* minimum resources, in terms of SBs number, should be allocated to each numerology. In addition, it provides a periodic dropping control.
- Phase 1.2 determines *which* SBs should be assigned to each numerology with the aim of maximizing system throughput and reducing the INI.

In the following subsections, we describe step-by-step how each phase works.

### 2.5.1 Phase 1.1

To guarantee the users' requirements, the algorithm calculates the  $\mathbf{X}_{\text{GBR}}$  vector, which contains the GBR numerologies sorted by priority level. Periodically, the algorithm determines the number of SBs ( $N_{SB_x}$ ) necessary to satisfy the UEs in  $\mathbf{U}_{\mathbf{x}}^{(l)}$ , with  $x \in \mathbf{X}_{\text{GBR}}$ , starting from the first element of  $\mathbf{X}_{\text{GBR}}$ . Once all the  $N_{SB_x}$  values have been computed  $\forall x \in \mathbf{X}_{\text{GBR}}$ , if there are still resources available, they will be reserved to the numerology  $x \in \mathbf{X}_{\text{BE}}$ . The overall procedure of the proposed framework is reported in Fig. 2.5. Moreover, in order to compare the improved

## 2.5. The Proposed Control Framework



**Figure 2.4:** Two control levels framework for radio resource allocation and scheduling.

CARAM with the previous one [1], we report in Fig. 2.6 the related flow diagram. Let us start describing CARAM-new.

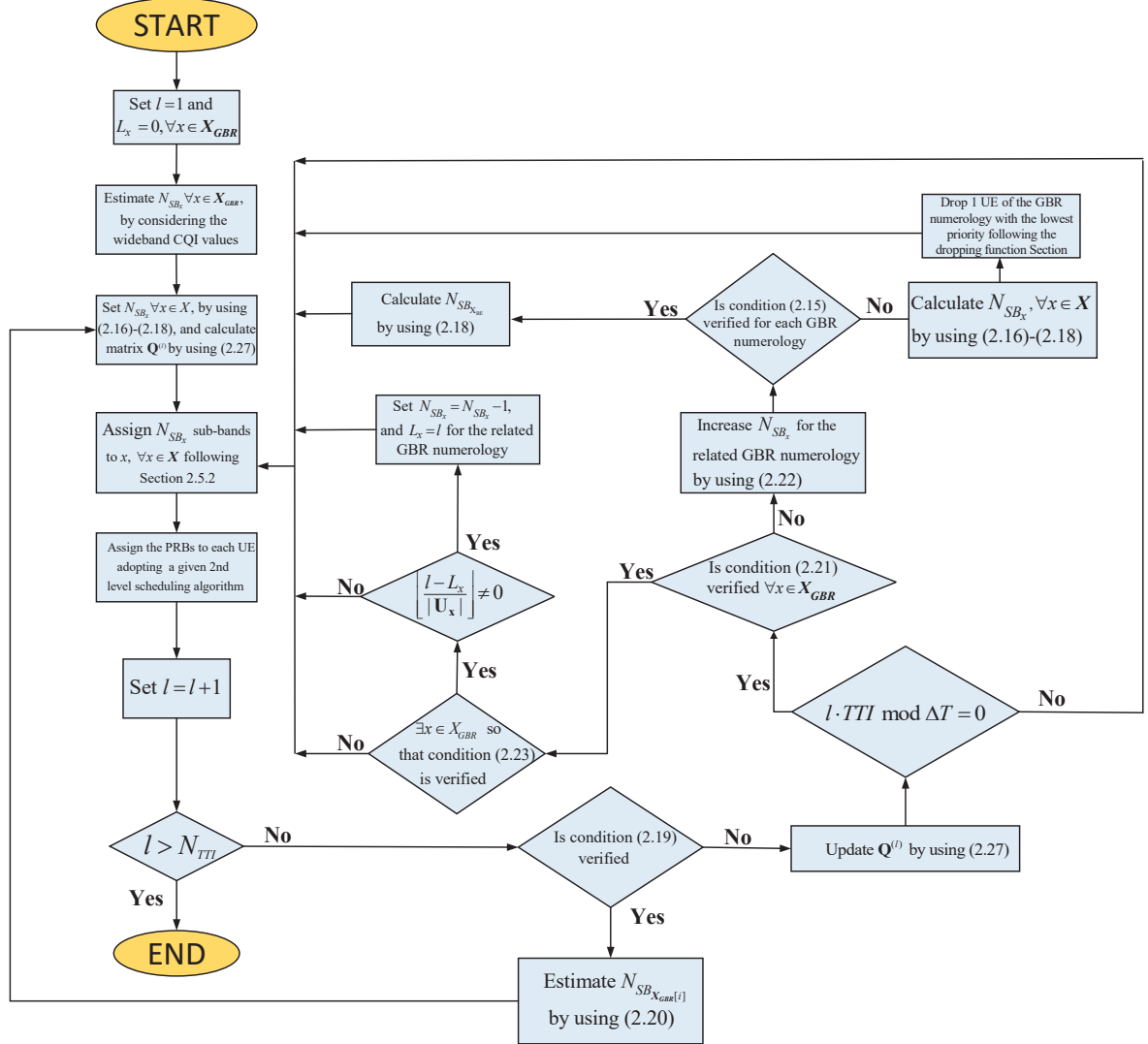
At the first TTI, unlike the procedure followed in our previous work [1], where the  $N_{SB_x}$  is initialized considering the best channel conditions, here we exploit the wideband CQI values in  $\mathbf{w}^{(1)}$  for initializing the  $N_{SB_x}$  values. Specifically, starting from the numerology with the highest priority, i.e.,  $\mathbf{X}_{\text{GBR}}[1]$ , for each UE  $i \in \mathbf{U}_{\text{GBR}}^{(1)}$ , we calculated the expected throughput per PRB by considering the wideband CQI value  $wCQI_i^{(1)}$ . Then, we calculated the number of PRBs necessary to UE  $i$  for guaranteeing the required  $R_{i_{\text{GBR}}[1]}$ . Finally, we define as  $N_{PRB_{\text{GBR}}^{(1)}}$  the sum of the number of PRBs required by each UE belonging to  $\mathbf{U}_{\text{GBR}}^{(1)}$  and determine

$$N_{SB_{\text{GBR}}[1]} = \left\lceil \frac{N_{PRB_{\text{GBR}}^{(1)}}}{z} \right\rceil. \quad (2.14)$$

Then, for the second numerology  $\mathbf{X}_{\text{GBR}}[2]$ , the algorithm operates in the same way. In general, for the numerology  $\mathbf{X}_{\text{GBR}}[i]$ , the following constraint should be guaranteed:

$$N_{SB_{\text{GBR}}[i]} \leq \begin{cases} N_{SB}, & \text{if } i = 1 \\ N_{SB} - \sum_{j=1}^{i-1} N_{SB_{\text{GBR}}[j]} & \text{otherwise.} \end{cases} \quad (2.15)$$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



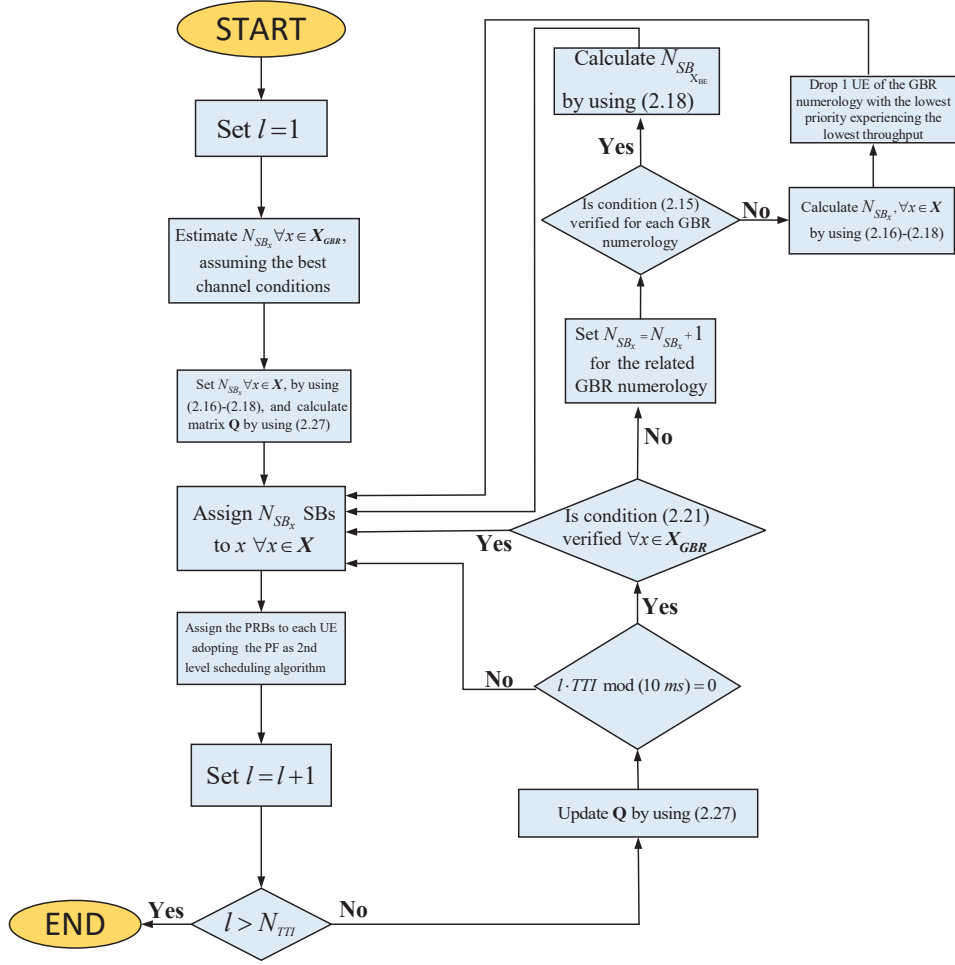
**Figure 2.5:** Radio Resource Management Flow diagram, adopting CARAM-new as 1st level allocation algorithm.

If condition (2.15) is not verified for one or more GBR numerologies, it means that we are operating under overload condition. In this case, the  $N_{SB_{X_{GBR}}}$  values are set as follows. The number of SBs allocated to Numerology  $X_{GBR}[1]$  is

$$N_{SB_{X_{GBR}[1]}} = \min \left\{ N_{SB}, N_{SB_{X_{GBR}[1]}} \right\}. \quad (2.16)$$

Then, the other number of SBs allocated to the GBR numerologies are

## 2.5. The Proposed Control Framework



**Figure 2.6:** Flow diagram of the preliminary control framework proposed in [1], adopting the previous version of CARAM as 1st level allocation algorithm.

calculated recursively as:

$$N_{SB_{\mathbf{X}_{GBR}^{[i]}}} = \min \left\{ N_{SB_{\mathbf{X}_{GBR}^{[i]}}, N_{SB} - \sum_{j=1}^{i-1} N_{SB_{\mathbf{X}_{GBR}^{[j]}}} \right\}. \quad (2.17)$$

Clearly, under this condition of overload some GBR services can not be satisfied. So, it will be necessary to drop some GBR services. Our approach is different from that of [1], and it is explained in Section 2.5.1. If there are still available SBs, they are allocated to the BE Numerology, i.e.,

$$N_{SB_{\mathbf{X}_{BE}}} = N_{SB} - \sum_{j=1}^{|\mathbf{X}_{GBR}|} N_{SB_{\mathbf{X}_{GBR}^{[j]}}}. \quad (2.18)$$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

It follows Phase 1.2, where the proper SBs are assigned to the related numerologies, and then the second level scheduler allocates the PRBs to the UEs.

Clearly, the amount of SBs to be allocated to each numerology should be updated during the simulation time. One of the main reasons is that, in general, the number of user requests may change during the simulation time. Therefore, for each TTI  $l > 1$ , the algorithm verifies whether there is a new GBR user in the system, i.e., if

$$\exists i \in \{1, \dots, |\mathbf{X}_{\text{GBR}}|\} \text{ so that } \left| \mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l)} \right| > \left| \mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l-1)} \right| \quad (2.19)$$

If so, for each numerology  $i$ th GBR numerology that satisfies this condition, we define  $\mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l)/(l-1)}$  as the complement vector of  $\mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l-1)}$  with respect of  $\mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l)}$ . Then, the algorithm calculates the sum of the number of PRBs necessary for the UEs in  $\mathbf{U}_{\mathbf{X}_{\text{GBR}}[i]}^{(l)/(l-1)}$  for guaranteeing the required  $R_{i_{\mathbf{X}_{\text{GBR}}[i]}}$  by considering the wideband CQI values in  $wCQI_i^{(l)}$ . The number of PRBs is termed as  $N_{PRB_{\mathbf{X}_{\text{GBR}}[i]}^{(l)/(l-1)}}$ . So, the algorithm sets

$$N_{SB_{\mathbf{X}_{\text{GBR}}[i]}} = N_{SB_{\mathbf{X}_{\text{GBR}}[i]}} + \left\lceil \frac{N_{PRB_{\mathbf{X}_{\text{GBR}}[i]}^{(l)/(l-1)}}}{z} \right\rceil. \quad (2.20)$$

Finally, all the  $N_{SB_x}$  values should be updated following (2.16)-(2.18).

We underline that, since the system needs to ensure the QoS requirements for the entire simulation length, the arrival of a new GBR service request is not the only check to be verified to update the number of SBs allocated to the GBR numerologies. Therefore, with a periodicity of  $\Delta T$ , we exploit the feedback between the two levels to checks for each GBR numerology if all users are satisfied, i.e.,  $\forall x \in \mathbf{X}_{\text{GBR}}$ , the algorithm checks whether

$$\sum_{i \in \mathbf{U}_{\mathbf{x}}^{(l)}} s_{i_x} = \left| \mathbf{U}_{\mathbf{x}}^{(l)} \right|. \quad (2.21)$$

If this condition does not occur for at least one numerology, the algorithm selects among them the GBR numerology with the highest priority, and increases the related  $N_{SB_x}$  value by as many units as necessary to assign



at least one PRB to each unsatisfied user, i.e.,

$$N_{SB_x} = N_{SB_x} + \left\lceil \frac{|\mathbf{U}_x^{(l)}| - \sum_{i \in \mathbf{U}_x^{(l)}} s_{i_x}}{z} \right\rceil. \quad (2.22)$$

Therefore, all the other  $N_{SB_x}$  values are updated accordingly. However, if condition (2.15) is no longer guaranteed for one GBR numerology, one of the associated services belonging to that numerology should be dropped. We empathize that this new approach regarding the increase of the  $N_{SB_x}$  values is an enhancement compared to our preliminary work [1], where  $N_{SB_x}$  is simply increased by one unit.

Otherwise, if condition (2.21) is verified, the algorithm verifies whether the amount of resources allocated in the previous TTIs to the GBR numerologies are excessively greater compared to those sufficient to guarantee the requirements the related UEs, the  $N_{SB_x}$  values should be properly decreased. This can occur when either the initial estimate of  $N_{SB_x}$  is oversized, or the channel conditions improve over time. We underline that the CARAM scheme [1] does not include any decrement strategy. We note that deciding to decrease the value of  $N_{SB_x}$  is not straightforward, as you risk not satisfying some users in subsequent TTIs. In the following subsection, we explain in detail how the decrement strategy works.

#### The $N_{SB_x}$ decrement strategy

Let define  $\mathbf{T}_{x, \mathbf{W}}^{1, (l)} = [T_{1_x, \mathbf{W}}^{1, (l)}, T_{2_x, \mathbf{W}}^{1, (l)}, \dots, T_{|\mathbf{U}_x|, \mathbf{W}}^{1, (l)}]$  as the vector containing the theoretical average throughput  $T_{i_x, \mathbf{W}}^{1, (l)}$  obtainable by each user of numerology  $x \in \mathbf{X}_{\text{GBR}}$  over window  $\mathbf{W}_{i_x}^{(l)}$  assuming that it received only one PRB until TTI  $l$ . We denote the UE  $i_x$  as "over-satisfied" if the condition  $T_{i_x, \mathbf{W}}^{(l)} - T_{i_x, \mathbf{W}}^{1, (l)} \geq R_{i_x}$  is true. This means that the UE has obtained from the scheduling algorithm one or more resources in addition to those needed to satisfy the GBR requirement.

Therefore, if exist one GBR numerology  $x$  so that

$$T_{i_x, \mathbf{W}}^{(l)} - T_{i_x, \mathbf{W}}^{1, (l)} \geq R_{i_x}, \forall i_x \in \mathbf{U}_x^{(l)}, \quad (2.23)$$

and if there have been no decrements in the last  $|\mathbf{U}_x^{(l)}|$  TTIs, then the algorithm properly decreases the related  $N_{SB_x}$  value of one unit, as shown

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

in Fig. 2.5. It is evident that our decrement strategy is very conservative. In fact, the decrement is carried out not only whether all UEs  $\in \mathbf{U}_x^{(l)}$  would be over satisfied, but also at most once every  $|\mathbf{U}_x^{(l)}|$  TTIs. In addition,  $N_{SB_x}$  is decremented only by one unit. Generally, it occurs that  $|\mathbf{U}_x^{(l)}| > z$ , where  $z$  represents the number of PRBs per SB (see Fig. 2.2), and therefore the  $N_{SB_x}$  decrease affects a small number of users.

We empathize that the proposed 1st level algorithm aims to adopt 2nd level scheduling algorithms available in literature without modifying them. So, it is necessary to introduce this conservative strategy due to the problems encountered in the scheduling algorithms. In fact, if our level 1 were less conservative, e.g., it would perform the decrement strategy only when half of the GBR UEs are over satisfied, then the 2nd level algorithm may erroneously reduce resources even to users who are non satisfied. This issue is expected by adopting PF and BCQI schedulers, since they are QoS-unaware. However, it can occur also by adopting the BCQI-HD although it is a QoS-aware scheduler, due to the fact that it aims to guarantee the GBR requirement inside a single TTI, without considering the past scheduling allocations, i.e.,  $T_{i_x, \mathbf{w}}^{(l)}$ .

It is worth to underline that the introduced decrement strategy allows in the first TTI to perform a more realistic estimate of  $N_{SB_x}$ , compared to the one of our previous work [1]. In fact, therein the algorithm assumes the maximum CQI values for each UE, corresponding to the minimum  $N_{SB_x}$  values. Instead, the improved strategy adopts the wideband CQI values and, consequently, initializes all the  $N_{SB_x}$  to values close to the optimal ones, but that may be either underestimated or overestimated. In the latter case, unlike CARAM, the CARAM-new algorithm is capable of properly decreasing the  $N_{SB_x}$  values.

### The dropping function

In this subsection, we describe the proposed criterion to select which service should be dropped. The aim is to drop the service belonging to the GBR Numerology with the lowest priority that is the "hardest" to satisfy. For doing so, first the algorithm selects the GBR Numerology  $x^*$  having the lowest priority among the ones with  $N_{SB_x} > 0$ . Then, for Numerology  $x^*$ , the algorithm determines the vector  $\mathbf{D}_{x^*}^{(l)}$  whose elements

$D_{i_{x^*}}^{(l)}$  are calculated as:

$$D_{i_{x^*}}^{(l)} = \begin{cases} \frac{R_{i_{x^*}} - T_{i_{x^*}, \mathbf{w}}^{(l)}}{T_{i_{x^*}, \mathbf{w}}^{1, (l)}}, & \text{if } R_{i_{x^*}} > T_{i_{x^*}, \mathbf{w}}^{(l)} \\ 0 & \text{otherwise.} \end{cases} \quad (2.24)$$

Finally, the UE

$$i' = \arg \max_{i \in \mathbf{U}_{\mathbf{x}^*}} \left\{ \mathbf{D}_{\mathbf{x}^*}^{(l)} \right\} \quad (2.25)$$

is dropped.

This strategy represents an improvement compared to the one used in our previous work [1], where the user with the lowest achieved throughput was dropped. In fact, that UE is not necessary the most difficult to satisfy, because this depends on the CQI values, which clearly affect  $T_{i_x}^{1, (l)}$  which is used in (2.24).

### Adaptation of the proposed CARAM-new to different 2nd level schedulers

In this subsection, we describe how we adapt the proposed CARAM-new to work well with different 2nd level schedulers, i.e., PF, BCQI, and BCQI-HD. This adaptation is necessary to allow the framework to obtain the best achievable performances without modifying the above schedulers.

We start describing the adaptation when the PF is adopted, that is, the same scheduler used in the preliminary work [1]. As reported in Section 2.7, given a numerology  $x$ , the PF algorithm allocates all the  $N_{SB_x}$  sub-bands only to one UE per TTI. So, the check on conditions (2.21) and (2.23) should be carried only after that the radio resources have been allocated to each UE at least once. As a result, this check should be performed periodically, with a periodicity of  $\left| \mathbf{U}_{\mathbf{x}}^{(l)} \right|$  TTIs, i.e., we set  $\Delta T = \left| \mathbf{U}_{\mathbf{x}}^{(l)} \right| \cdot TTI$ . This is an improvement compared to the previous work, where  $\Delta T$  was fixed to 10 ms regardless of the number of UEs.

As regards the BCQI-HD, unlike the PF, if there are enough PRBs, then it schedules all the UEs in the given TTI. Otherwise, only  $|\mathbf{L}_{\mathbf{x}}| = z \cdot N_{SB_x}$  UEs are scheduled in the Numerology  $x$ . As a result, the conditions (2.21) and (2.23) can be checked every TTI if there are enough

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

SBs allocated to  $x$ , otherwise the gNB should wait that all UEs have been scheduled. Specifically, we set:

$$\Delta T = \left\lceil \frac{|\mathbf{U}_{\mathbf{x}}^{(l)}|}{z \cdot N_{SB_x}} \right\rceil \cdot TTI. \quad (2.26)$$

We note that this choice improves significantly the responsiveness of the proposed framework compared to the previous work, especially when  $|\mathbf{U}_{\mathbf{x}}^{(l)}|$  is lowen that  $z \cdot N_{SB_x}$ .

Finally, as regards the BCQI, we remind that it schedules each PRB to the UE experiencing the best channel condition. This means that some UEs may be never scheduled, and thus conditions (2.21) and (2.23) may never been verified due to the intrinsic behavior of the scheduler. As a result, the CARAM-new without a proper adaptation produces that all the SBs will be allotted only to the GBR Numerology with the highest priority for trying to guarantee, in vain, the services of this Numerology . To overcome this issue, we adopt the following strategy. For the first TTI the CARAM-new is run without any changes. Then, we define  $\mathbf{U}_{\mathbf{x}_s}^{(l)} \subseteq \mathbf{U}_{\mathbf{x}}^{(l)}$  as the vector containing all UEs in  $\mathbf{U}_{\mathbf{x}}^{(l)}$  so that  $T_{i_x}^{(l)} > 0$ . So, starting from the second TTI, the scheduler updates  $\mathbf{U}_{\mathbf{x}_s}^{(l)}$ , on basis of the feedback received from the 2nd level scheduler, and adopts it instead of  $\mathbf{U}_{\mathbf{x}}^{(l)}$  for all the operations carried out in the first level. This adaptation permits the GBR Numerologies to obtain the amount of SBs necessary to satisfy the services with high CQI values, and the remaining resources are allotted to the BE Numerology, without waste.

### 2.5.2 Phase 1.2

In this phase, given the  $N_{SB_x}$  values for each numerology  $x \in \mathbf{X}$ , the CARAM-new algorithm determines *which* SBs to assign to each numerology with the aim of maximizing the system throughput and the spectral efficiency. To achieve this, we follow a channel-aware approach, in particular our algorithm is inspired by the procedure proposed in [30], where the services are typically 4G (FTP, streaming, VoIP) and therefore it is sufficient to allocate only one PRB to each user.

The algorithm makes use of the matrix  $\mathbf{Q}^{(l)}$ , which contains the value  $\mathbf{Q}^{(l)}[j, k]$  related to the  $j$ th numerology in  $\mathbf{X}$  and the  $k$ th sub-band in  $\mathbf{K}$ .

## 2.5. The Proposed Control Framework

For convenience, we define  $\mathbf{J} = \{1, \dots, N_x\}$ . Then,

$$\mathbf{Q}^{(l)}[j, k] = \frac{1}{|\mathbf{U}_{\mathbf{X}[j]}|} \sum_{i \in \mathbf{U}_{\mathbf{X}[j]}} m_{i,k}^{(l)}, \forall j \in \mathbf{J}, \forall k \in \mathbf{K}, \quad (2.27)$$

where  $m_{i,k}^{(l)}$  is a metric related to UE  $i$ , calculated on the SB  $k$ . More specifically, we assume as user metric the one used by the PF algorithm [38], which aims to reach a good trade-off between fairness and spectral efficiency. Therefore, we set  $m_{i,k}^{(l)}$  as:

$$m_{i,k}^{(l)} = \frac{d_{i,k}^{(l)}}{\bar{T}_i^{(l)}}, \quad (2.28)$$

where  $d_{i,k}^{(l)}$  is the estimated instantaneous data rate for UE  $i$  at the  $l$ th TTI over SB  $k$ , and  $\bar{T}_i^{(l)}$  is the average data rate achieved until time  $l$ , recursively calculated by using (2.13).

For the sake of clarity, we report the overall procedure in pseudo-code 1 and explain it by means of an example, consisting of several steps. We consider  $\mathbf{X} = \{0, 1, 2\}$ ,  $p_2 > p_1 > p_0$  and  $N_{SB} = 8$ . Furthermore, we assume  $\mathbf{X}_{\text{GBR}} = \{2, 1\}$ , and  $\mathbf{X}_{\text{BE}} = \{0\}$ . We also consider that on the basis of  $R_1$ , and  $R_2$ , the Phase 1.1 outputs are:  $N_{SB_1} = N_{SB_2} = 3$ . So, by using (2.18), it follows  $N_{SB_0} = 2$ . Obviously,  $\mathbf{J} = \{1, 2, 3\}$  and  $\mathbf{K} = \{1, \dots, 8\}$ .

### Step 1

In the first step, the algorithm calculates the matrix  $\mathbf{Q}^{(l)}$  which is based on the CQI values of matrix  $\mathbf{C}$  received from the gNB, and the users  $\in \mathbf{U}_{\mathbf{X}}^{(l)}, \forall x \in \mathbf{X}$ . We assume the following matrix  $\mathbf{Q}^{(l)}$ :

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & 0.87 & 0.93 & 0.93 & 0.91 & 0.89 & 0.91 & 0.95 \\ 0.89 & 0.97 & 1.00 & 0.88 & 0.87 & 0.91 & 0.85 & 0.91 \\ 0.85 & 0.86 & 0.97 & 0.91 & 0.97 & 0.87 & 0.89 & 0.91 \end{bmatrix}. \quad (2.29)$$

The first level allocation matrix  $\mathbf{A}^{(l)}$  is a zero matrix of size  $N_x \times N_{SB}$ .

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

### Pseudo-code 1 Phase 1.2 of CARAM-new

---

#### Definitions:

- $\mathbf{X}$ : set of numerologies, with cardinality  $N_x$ ;
- $\mathbf{J} = \{1, \dots, N_x\}$ ;
- $\mathbf{K}$ : set of sub-bands with cardinality  $N_{SB}$ ;
- $N_{SB_x}$ : number of sub-bands to be assigned to numerology  $x$
- $\mathbf{A}^{(l)}$ : First level allocation matrix.

#### Initialization:

- 1: calculate  $\mathbf{Q}^{(l)}$  by using (2.27)
- 2: set  $\mathbf{A}^{(l)}[j, k] = 0, \forall j \in \mathbf{J}, \forall k \in \mathbf{K}$

#### Start:

- 3: **while**  $\mathbf{J} \neq \emptyset$  **do**
  - 4:   calculate  $(j_M, k_M)$  by using (2.30)
  - 5:   set  $\mathbf{A}^{(l)}[j_M, k_M] = 1$
  - 6:   set  $\mathbf{K} = \mathbf{K} - \{k_M\}$  and  $\mathbf{Q}^{(l)}[j, k'_M] = -1, \forall j \in \mathbf{J}$
  - 7:   **while**  $\sum_{k=1}^{N_{SB}} \mathbf{A}^{(l)}[j_M, k] \leq N_{SB_x[j_M]}$  **do**
  - 8:     **if**  $k_M + 1 \notin \mathbf{K} \wedge k_M - 1 \notin \mathbf{K}$  **then**
  - 9:       **break**
  - 10:    **else**
  - 11:     set  $k'_M$  equal to one of the available adjacent sub-bands
  - 12:     **if**  $\mathbf{Q}^{(l)}[j_M, k'_M] \geq \max_{j \in \mathbf{J} - \{j_M\}} \mathbf{Q}^{(l)}[j, k'_M]$  **then**
  - 13:       set  $\mathbf{A}^{(l)}[j_M, k'_M] = 1$
  - 14:       set  $\mathbf{K} = \mathbf{K} - \{k_M\}$  and  $\mathbf{Q}^{(l)}[j, k'_M] = -1, \forall j \in \mathbf{J}$
  - 15:       set  $k_M = k'_M$
  - 16:     **end if**
  - 17:    **end if**
  - 18:    **end while**
  - 19:    set  $\mathbf{J} = \mathbf{J} - \{j_M\}$  and  $\mathbf{Q}^{(l)}[j_M, k] = -1, \forall k \in \mathbf{K}$
  - 20: **end while**
  - 21:  $\mathbf{J} = \{1, \dots, N_x\}$
  - 22: **while**  $\mathbf{K} \neq \emptyset$  **do**
  - 23:   **for**  $k \in \mathbf{K}$  **do**
  - 24:     **if**  $k \neq |\mathbf{K}| \wedge \max_{j \in \mathbf{J}} \mathbf{A}^{(l)}[j, k + 1] == 1$  **then**
  - 25:        $j_M = \arg \max_{j \in \mathbf{J}} \mathbf{A}^{(l)}[j, k + 1]$
  - 26:     **else if**  $k == |\mathbf{K}| \wedge \max_{j \in \mathbf{J}} \mathbf{A}^{(l)}[j, k - 1] == 1$  **then**
  - 27:        $j_M = \arg \max_{j \in \mathbf{J}} \mathbf{A}^{(l)}[j, k - 1]$
  - 28:     **if**  $\sum_{k^*=1}^{N_{SB}} \mathbf{A}^{(l)}[j_M, k^*] < N_{SB_x[j_M]}$  **then**
  - 29:       set  $\mathbf{A}^{(l)}[j_M, k] = 1$
  - 30:     **else**
  - 31:       Assign SB  $k$  to the non-satisfied numerology with the lowest priority
  - 32:     **end if**
  - 33:     set  $\mathbf{K} = \mathbf{K} - \{k\}$
  - 34:    **end if**
  - 35:    **end for**
  - 36: **end while**
-

### Step 2

The algorithm searches for the indices related to the highest metric value in matrix  $\mathbf{Q}^{(l)}$ , as follows:

$$(j_M, k_M) = \arg \max_{(j,k) \in \mathbf{J} \times \mathbf{K}} (\mathbf{Q}^{(l)}[j, k]), \quad (2.30)$$

where  $\times$  denotes the Cartesian product. Then, the SB  $k_M$  is assigned to numerology  $\mathbf{X}[j_M]$ , as reported in Line 5 of pseudo-code 1. In example (2.29), it follows  $(j_M, k_M) = (2, 3)$ , and consequently, the SB 3 is allocated to Numerology 1, i.e.,  $\mathbf{A}^{(l)}(2, 3) = 1$ . Finally, the algorithm sets in matrix  $\mathbf{Q}$  the column  $k_M$  equal to -1 and also deletes  $k_M$  from vector  $\mathbf{K}$  (see Lines 6 of pseudo-code 1). In the example, it follows:

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & 0.87 & -1 & 0.93 & 0.91 & 0.89 & 0.91 & 0.95 \\ 0.89 & 0.97 & -1 & 0.88 & 0.87 & 0.91 & 0.85 & 0.91 \\ 0.85 & 0.86 & -1 & 0.91 & 0.97 & 0.87 & 0.89 & 0.91 \end{bmatrix}, \quad (2.31)$$

and  $\mathbf{K} = \{1, 2, 4, \dots, 8\}$ .

### Step 3

In the third Step, the algorithm tries to assign the adjacent subbands to the same numerology, in order to reduce the INI phenomenon. The details are reported in pseudo-code 1 Lines 7-20. Specifically, if the number of allocated SBs is less than  $N_{SB_{\mathbf{X}[j_M]}}$  (i.e., the Numerology is not satisfied) and the adjacent SB has not been already assigned, then the algorithm checks whether the numerology  $\mathbf{X}[j_M]$  exhibits the highest metric value also in the adjacent SB. If so, then the algorithm assigns this SB to it. Therefore, the algorithm sets in matrix  $\mathbf{Q}^{(l)}$  the column  $k_M$  as -1, and also deletes it from  $\mathbf{K}$ . The procedure is executed iteratively until all the above conditions are valid. Looking at (2.31), it follows that SB 2 is assigned to Numerology 1, i.e.,  $\mathbf{A}^{(l)}(2, 2) = 1$ , the matrix  $\mathbf{Q}^{(l)}$  is update as follows:

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & -1 & -1 & 0.93 & 0.91 & 0.89 & 0.91 & 0.95 \\ 0.89 & -1 & -1 & 0.88 & 0.87 & 0.91 & 0.85 & 0.91 \\ 0.85 & -1 & -1 & 0.91 & 0.97 & 0.87 & 0.89 & 0.91 \end{bmatrix}, \quad (2.32)$$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

and  $\mathbf{K} = \{1, 4, \dots, 8\}$ . We note that SB 1 or SB 4 cannot be assigned to Numerology 1, since Numerology 1 does not exhibit the highest metric.

Then, the algorithm sets in matrix  $\mathbf{Q}^{(l)}$  the row  $j_M$  as equal to -1 and delete from  $\mathbf{J}$  the element  $j_M$ . It follows that  $\mathbf{J} = \{1, 3\}$  and

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & -1 & -1 & 0.93 & 0.91 & 0.89 & 0.91 & 0.95 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ 0.85 & -1 & -1 & 0.91 & 0.97 & 0.87 & 0.89 & 0.91 \end{bmatrix}. \quad (2.33)$$

The allocation carried out until this point is reported in Fig. 2.7a.

### Step 4

In Step 4, the algorithm simply repeats Steps 2-3 until all Numerologies get at least one SB. In the example (2.33), it follows that SB 5 is assigned to numerology 2 (corresponding to row 3), i.e.,  $\mathbf{A}^{(l)}(3, 5) = 1$ ,  $\mathbf{K} = \{1, 4, 6, 7, 8\}$ , and the matrix  $\mathbf{Q}^{(l)}$  becomes:

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & -1 & -1 & 0.93 & -1 & 0.89 & 0.91 & 0.95 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ 0.85 & -1 & -1 & 0.91 & -1 & 0.87 & 0.89 & 0.91 \end{bmatrix}. \quad (2.34)$$

The expansion of the band is stopped since both SB 4 and SB 6 do not exhibit the highest metric for the Numerology 2. As a result,  $\mathbf{J} = \{1\}$  and row 3 is set in matrix  $\mathbf{Q}^{(l)}$  equal to -1:

$$\mathbf{Q}^{(l)} = \begin{bmatrix} 0.93 & -1 & -1 & 0.93 & -1 & 0.89 & 0.91 & 0.95 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix}. \quad (2.35)$$

The allocation becomes as reported in Fig. 2.7b.

At this point, the algorithm repeats Steps 2-4 for the Numerology 0, assigning first SB 8 and then SB 7. We note that it does not allocate SB6 to Numerology 0 since  $N_{SB_0} = 2$ . Finally we obtain  $\mathbf{J} = \emptyset$ ,  $\mathbf{K} = \{1, 4, 6\}$ ,  $\mathbf{Q}^{(l)} = \emptyset$ , and

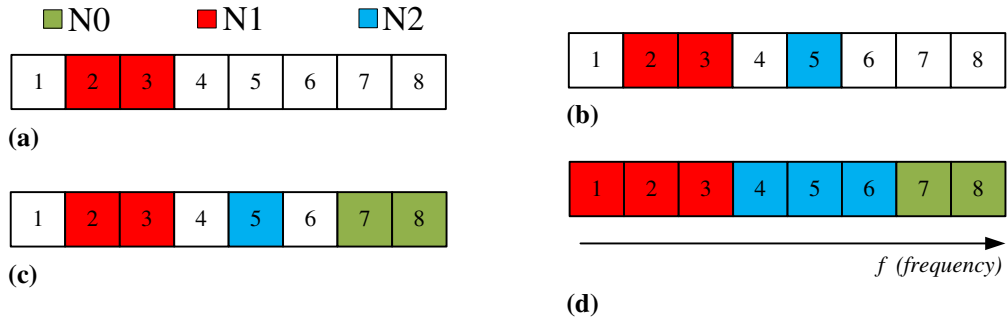


$$\mathbf{A}^{(l)} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}. \quad (2.36)$$

Fig. 2.7c shows the sub-bands assigned to all the numerologies at the end of Step 4. We note that the SBs 1, 4, and 6 are colored in white because are still available.

### Step 5

Step 5 is performed only if there are still resources available, i.e.,  $\mathbf{K} \neq \emptyset$ . Our approach is to reduce the alternation between different numerologies as much as possible, to minimize the phenomenon of INI. The details are reported in pseudo-code 1 Lines 21-36. First, the algorithm re-admits all the numerologies. Then, starting from the first SB not assigned, the algorithm tries to assign it to one of the numerologies not satisfied that have obtained one of the adjacent SBs (see Line 24-29 of pseudo-code 1). Otherwise, it assigns the SB to the non-satisfied numerology with the lowest priority. The algorithm proceeds the same way scrolling through all the SBs still available. In the example, the algorithm re-admits Numerology 1 and 2, and SB 1 is assigned to Numerology 1, while SBs 4 and 6 to Numerology 2. The final allocation result presented in Fig. 2.7d shows that our algorithm, when the numerology requirements allow it, performs a contiguous band allocation, in order to decrease the INI.



**Figure 2.7:** Steps of the SB allocation: (a) Step 3; (b) Step 4 first round ; (c) Step 4 second round; (d) final SB allocation.

### 2.5.3 Time complexity analysis

The proposed 1st level algorithm is designed to operate in near real-time, so it is important to evaluate its time complexity. Moreover, since the algorithm should work timely in any condition, we consider the worst case time complexity, i.e.,  $l_{start,i} = 1, \forall i \in \mathbf{U}$ . Consequently, we utilize the Big-O notation that permits to provide an upper bound of the runtime of the algorithm and guarantees that the algorithm will never take time longer than that. Since CARAM-new is composed of two phases, we first calculate the time complexity of each phase separately, and then, we consider the conjunction between them.

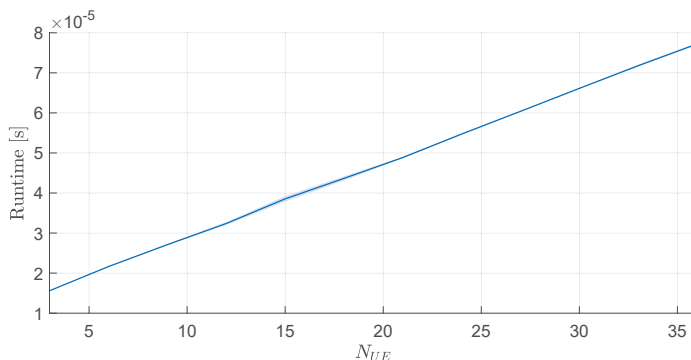
As regards Phase 1.1, the worst case corresponds to  $\mathbf{U}_{\mathbf{X}_{BE}} = \emptyset$ , i.e., all UEs require a GBR service. We remind that the algorithm calculates the following two quantities for each GBR UE. First, the expected throughput per PRB. Second, the number of PRBs necessary to guarantee  $R_i$ . Finally, the algorithm calculates  $N_{SB_x}$  for each GBR numerology  $x$ . It is straightforward to derive that the time complexity of this phase depends on the computation of the expected UE throughput values. Since each UE reports the wideband CQI value, the expected UE throughput is the same for each PRB, and so, it can be computed only once per UE. As a consequence, the Phase 1.1 belongs to  $O(N_{UE})$ .

As regards the Phase 1.2, we calculate its complexity by exploiting the pseudo-code 1. Therein, the algorithm assigns to which numerology each of the  $\mathbf{K}$  SBs is assigned. To do so, in the first part of this phase (lines 3-20) the algorithm assigns some of the SBs (at least  $|\mathbf{J}|$  SBs), while the other SBs (if still available) are assigned in the second part (lines 21-36). Since each iteration corresponds to assigning one SB, the whole second phase is directly proportional to the number of SBs, i.e.,  $|\mathbf{K}| = N_{SB}$ . Therefore, this phase belongs to  $O(N_{SB})$ .

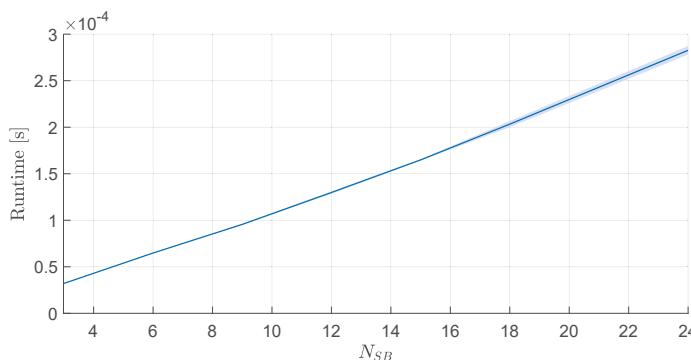
Finally, from the asymptotic analysis it is known that when adding functions, the order of the sum is just the sum itself. Therefore, the time complexity of the algorithm is  $O(N_{UE} + N_{SB})$ .

However, the Big O-notation provides only information on how the algorithm scale on the basis of the different inputs. So, since the algorithm should output its values within one TTI (i.e., 1 ms) we need to characterize the amount of time it takes to give the outputs. Therefore, we performed several simulation of the runtime for the two phases,

## 2.5. The Proposed Control Framework



**Figure 2.8:** Average runtime together with the standard deviation (shaded area) for the Phase 1.1 of the new-CARAM under different  $N_{UE}$  values.



**Figure 2.9:** Average runtime together with the standard deviation (shaded area) for the Phase 1.2 of the new-CARAM under different  $N_{SB}$  values.

by varying  $N_{UE}$  for the Phase 1.1 and  $N_{SB}$  for the Phase 1.2. The runtime simulations were run with a Intel(R) Core(TM) i9-10940X CPU @ 3.30GHz, with 14 cores, and 128GB RAM. Software environment includes Python 3.8.5 and benchit 0.0.6. Figs. 2.8 and 2.9 report the average runtime together with the standard deviation for Phase 1.1 and 1.2, respectively. We note that the runtime for the first phase is much less than the second phase, regardless on the number of UEs, and also the standard deviation of the first phase is very reduced compared to the second phase one. As expected, the trend in both Figures is linearly increasing with respect to  $N_{UE}$  and  $N_{SB}$ , respectively. The total runtime for the  $N_{UE}$  and  $N_{SB}$  parameters adopted in our simulations (see Table 2.1), is largely less than 1 TTI. In fact, the summed runtime, even considering the positive standard deviation, is less than 0.4 ms.

## 2.6 Performance Evaluation

---

In this Section, we evaluate the performance of the proposed 1st allocation algorithm by simulations in MATLAB environment. Results are averaged over 50 independent simulations, each of these consisting of  $N_{TTI}$  TTIs. We consider the downlink transmission scenario of a single gNB,  $\mathbf{X}_{GBR} = \{2, 1\}$ ,  $\mathbf{X}_{BE} = \{0\}$ ,  $p_2 > p_1 > p_0$ , and  $l_{start,i} = 1, \forall i \in \mathbf{U}$ . We consider also different UE ratios, where the UE ratio is defined as:

$$|\mathbf{U}_0| : |\mathbf{U}_1| : |\mathbf{U}_2|. \quad (2.37)$$

The overall parameters are reported in Table 2.1. We compare the performance of the proposed CARAM-new with CB, PBNA [30], and the previous CARAM [1], denoted as CARAM-old. However, to evaluate these 1st level control schemes in run-time, it is necessary to adopt a given 2nd level scheduling algorithm. At this regard, we choose the PF [24], BCQI, and BCQI-HD [25] algorithms.

The metrics used to evaluate our algorithm are described below.

- Throughput of the BE Numerology

The total throughput of Numerology 0 is

$$T_0 = \sum_{i \in \mathbf{U}_0} T_{i_0}, \quad (2.38)$$

where  $T_{i_0}$  is calculated by using (2.6).

- Number of satisfied UEs

It is the number of satisfied UEs, taking into account the priority of the Numerology. Specifically, it counts the number of satisfied UEs of the Numerology 2, and only if they are all satisfied it adds also the number of satisfied UEs of the Numerology 1, i.e.,

$$\sum_{i=1}^{|\mathbf{U}_2|} s_{i_2} + \sum_{i=1}^{|\mathbf{U}_1|} s_{i_1} \cdot \left[ \frac{\sum_{i=1}^{|\mathbf{U}_2|} s_{i_2}}{|\mathbf{U}_2|} \right]. \quad (2.39)$$

- Responsiveness

**Table 2.1:** *System Parameters*

Parameter	Value
Carrier Frequency	2 GHz
Pathloss Model	$128.1 + 3.76 \cdot 10 \log_{10}(d_{[km]})$
Minimum distance UE/gNB	200 m
Maximum distance UE/gNB	5 km
Channel Estimation	Ideal
gNB Antenna Gain	18 dBi
UE Antenna Gain	0 dBi
UE Noise Figure	7 dB
gNB Transmit Power	46 dBm
Number of UEs	$N_{UE} = 36$
UEs arrangement	uniform (1), non-uniform (2)
Number of SBs	$N_{SB} = 24$
UE ratio	1:1:1, 4:1:1, 1:1:4
Number of PRBs per SB	$z = 4$
SB bandwidth	$\gamma = 720$ kHz
$\beta$	0.9
TTI	1 ms
DAW duration	$N_{DAW} = 2000$ TTIs
Simulation time	$N_{TTI} = 15000$ TTIs
GBR numerologies	$\mathbf{X}_{GBR} = \{2, 1\}$
BE numerologies	$\mathbf{X}_{BE} = \{0\}$
Priority levels	$p_2 > p_1 > p_0$
GBR requirements for stationary homogeneous scenarios	$R_{i_x} = 750$ Kbps, $\forall i_x \in \mathbf{U}_{\mathbf{X}_{GBR}}$ $R_{i_x} = 2$ Mbps, $\forall i_x \in \mathbf{U}_{\mathbf{X}_{GBR}}$
GBR requirements for stationary heterogeneous scenario	$R_{i_x} \in \{2, 5\}$ Mbps, $\forall i_x \in \mathbf{U}_{\mathbf{X}_{GBR}}$
GBR requirements for dynamic homogeneous scenarios	$R_{i_x} = 1$ Mbps, $\forall i_x \in \mathbf{U}_{\mathbf{X}_{GBR}}$ $R_{i_x} = 3$ Mbps, $\forall i_x \in \mathbf{U}_{\mathbf{X}_{GBR}}$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

It indicates how long a GBR Numerology takes, on average, to reach the steady state. Given  $x \in \mathbf{X}_{\text{GBR}}$ , the responsiveness of the Numerology  $x$  is

$$\frac{1}{|\mathbf{U}_x|} \sum_{i=1}^{|\mathbf{U}_x|} \Delta l_{i,x} \cdot TTI, \quad (2.40)$$

where  $\Delta l_{i,x}$  is the time, in terms of number of TTIs, elapsed from when the  $i$ th UE requests a GBR service of Numerology  $x$  until the requirement is met.

### 2.6.1 Stationary analysis with homogeneous GBR requirements

We first analyze the performance of CARAM-new in a stationary scenario where all UEs are not in motion, and all the GBR services have the same GBR requirement. We assume two different traffic load conditions. For each UE  $i \in \mathbf{U}_x$ , with  $x \in \mathbf{X}_{\text{GBR}}$ , we set  $R_{i,x} = 750$  kbps for the medium load scenario, and  $R_{i,x} = 2$  Mbps for the high load scenario. Moreover, in order to test our system under different channel qualities, we consider two different user arrangements in the cell, as depicted in Fig. 2.10. As shown, in the arrangement 1 all UEs are uniformly distributed within the coverage radius of the gNB equal to 5 km, with a minimum distance of 0.2 km from the gNB. Conversely, in the arrangement 2, the users of Numerology 2 are uniformly distributed at a distance ranging in  $[2, 5]$  km from the gNB, while the other UEs in  $[0.2, 2]$  km. We note that the latter arrangement aims to assess the ability of control schemes to guarantee higher priority services, even when their UEs are experiencing the worst channel conditions. We point out that the performance of all benchmark schemes is tested under different UE ratios, namely 1:1:1, 4:1:1, and 1:1:4.

In Fig. 2.11, we assess all the combinations of the 1st and 2nd level algorithms in a medium traffic load scenario with UEs arrangement 1 and UE ratio 1:1:1. In detail, Fig 2.11a shows the number of satisfied UEs, Fig 2.11b the throughput of the BE Numerology, and Figs. 2.11c-d the responsiveness of the Numerology 2 and 1, respectively. In these and in the following figures, the heights of the bars represent the median values, while the vertical grey lines highlight the interval between the lower quartile (25th percentile) and the upper quartile (75th percentile).

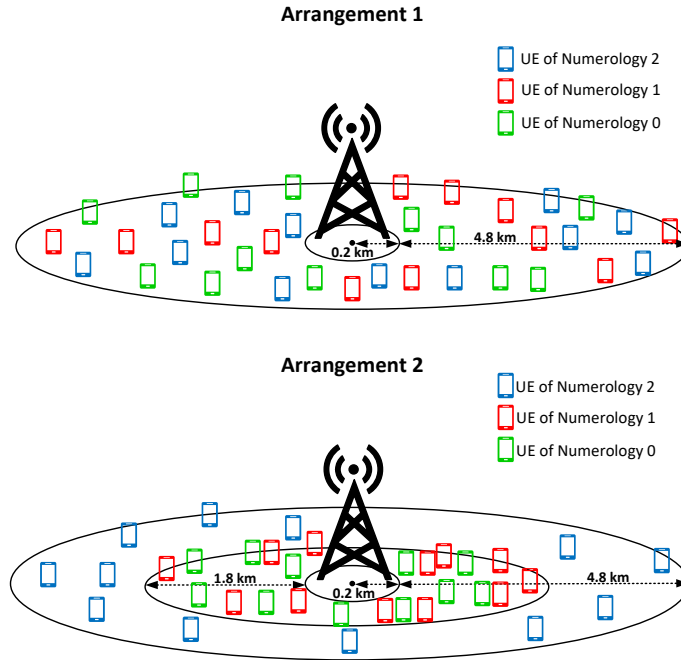


Figure 2.10: UEs Arrangements.

Let us start by comparing the performances achieved by the first level algorithms when the PF algorithm is adopted as the 2nd level scheduler. As shown in Fig. 2.11a, the CB exhibits the worst performance. This result was expected since the CB, on the basis of the number of UEs, splits the radio resources among the Numerologies without taking into account the QoS requirements. For this reason  $T_0 > 0$  (see Fig. 2.11b), but, on average, only 9 out of 12 UEs of Numerology 2 are satisfied and the same amount of UEs of Numerology 1, regardless of the priority requirements. Unlike CB, the other schemes take into account the priority and the GBR requirements, in fact,  $T_0 \neq 0$  only if all the GBR services are satisfied. Among them, the PBNA shows poor performance, since it is a channel-unaware algorithm. Conversely, thanks to the channel-awareness both the old and the CARAM-new achieve significantly better performance compared to the PBNA. However, the enhancements introduced in the CARAM-new ensure superior performance with respect to the CARAM-old, permitting not only to satisfy all the GBR services with a significantly lower responsiveness (see Figs. 2.11a, c-d), but also to obtain a non-zero  $T_0$  (see Fig. 2.11b).

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

Now, we analyze the performance achieved with the BCQI-HD. Thanks to the channel-awareness and the QoS-awareness of this 2nd level scheduler, the obtained results are significantly better for all the 1st level algorithms. In detail, even the CB is able to satisfy all GBR services, and also reaches a higher  $T_0$  value than PBNA and CARAM-old. However, even in this case, CARAM-new exhibits the best performance. We emphasize that the adaptation made in the proposed CARAM-new allows to obtain a significantly better responsiveness than the preliminary work. Finally, as regards the BCQI, our proposal achieves, in terms of number of satisfied UEs, the same poor results of PBNA and CARAM-old. This is an intrinsic limit of the adopted 2nd level scheduler that serves only the UEs perceiving the best channel quality. However, the adaptations introduced in the CARAM-new allow to obtain a  $T_0$  different from zero even if not all GBR UEs are satisfied.

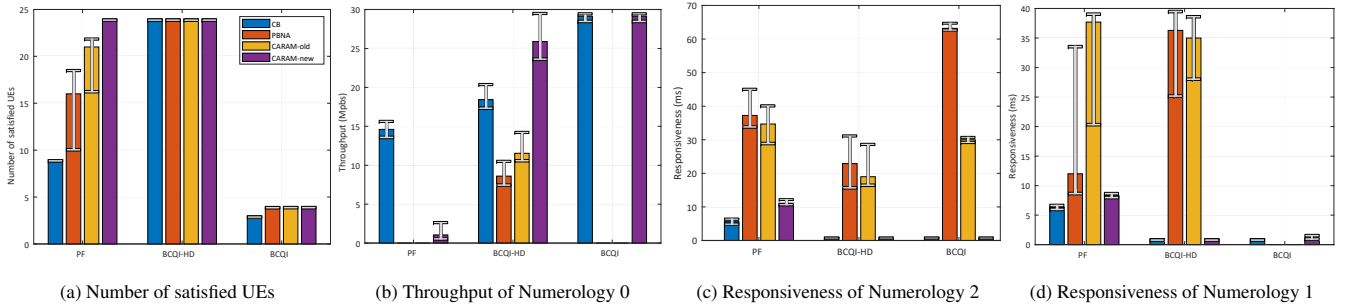
Now, we analyze the performance achieved with stringent characteristics, i.e., a high load scenario with Arrangement 2, UE ratio 1:1:1, and report it in Fig. 2.12. Due to the limiting channel conditions for the Numerology 2 and the high GBR requirements, the number of satisfied UEs is reduced for each configuration. As expected, the worst performance is achieved by the CB as it is not able to provide more radio resources to Numerology 2, which needs them. The only 2nd level scheduler that allows the CB to achieve the same number of satisfied UEs, with respect to the other 1st level schedulers, is the BCQI. This occurs because the BCQI serves only 4 UEs, on average, thus resulting in a non-critical scenario. As regards the CARAM-new, it reaches very important improvements compared to PBNA and CARAM-old, in terms of number of satisfied UEs when either the PF or the BCQI-HD schedulers are adopted, and an elevated  $T_0$  value when BCQI is adopted. The same improvements are valid for the responsiveness. We underline that this critical scenario does not allow even to CARAM-new to serve all GBR UEs, as in the previous case. Specifically, only adopting the BCQI-HD scheduler, our proposed scheme satisfies all services with the highest priority and some with lower priority. This means that, during the simulation time, the dropping function is performed. As a result, for the remaining GBR services, the condition (2.23) may be verified. So, the algorithm performs a decrement of the SBs allocated to the GBR Numerologies, thus permit-



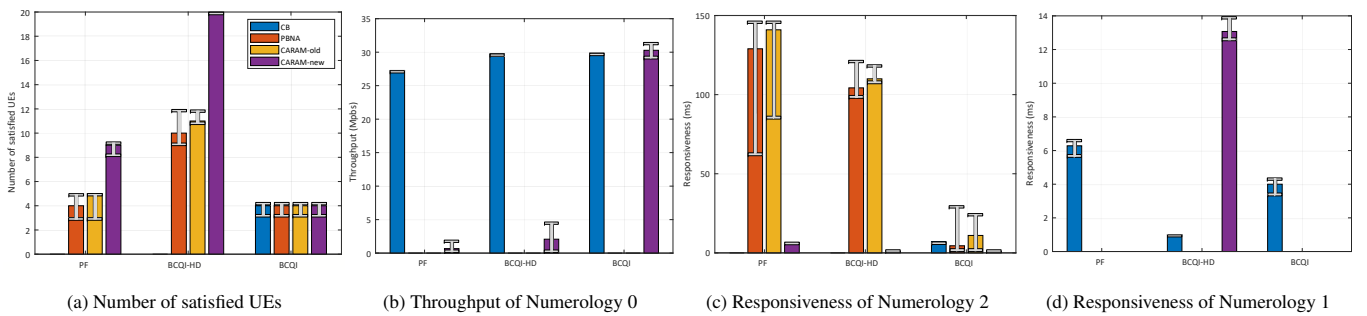
## 2.6. Performance Evaluation

ting the BE Numerology to obtain some radio resources. It is the reason why  $T_0$  is non-zero even though not all GBR services are guaranteed.

Finally, we report in Fig. 2.13 the performance achieved considering scenarios with the same characteristics as in Figs. 2.11 and 2.12, but with different UE ratios. In particular, Figs. 2.13a-b show the number of satisfied UEs and the throughput of the BE Numerology under a medium load scenario with Arrangement 1 and UE ratio 1:1:4 (i.e., the total number of UEs requesting a GBR service is 30, instead of 24 as in the case 1:1:1). As can be seen, the performances with this strong imbalance between numerologies are still in line with those reported in Fig. 2.11. Figs. 2.13c-d show the same KPIs under a high load scenario with Arrangement 2 and UE ratio 4:1:1 (i.e., the number of UEs with the highest priority is 6 instead of 12). Again, despite the imbalance between the number of UEs per numerology, the performances are like those reported in Fig. 2.12. The results reported in Fig. 2.13 show that the proposed CARAM-new is robust to changes in the UE ratios.

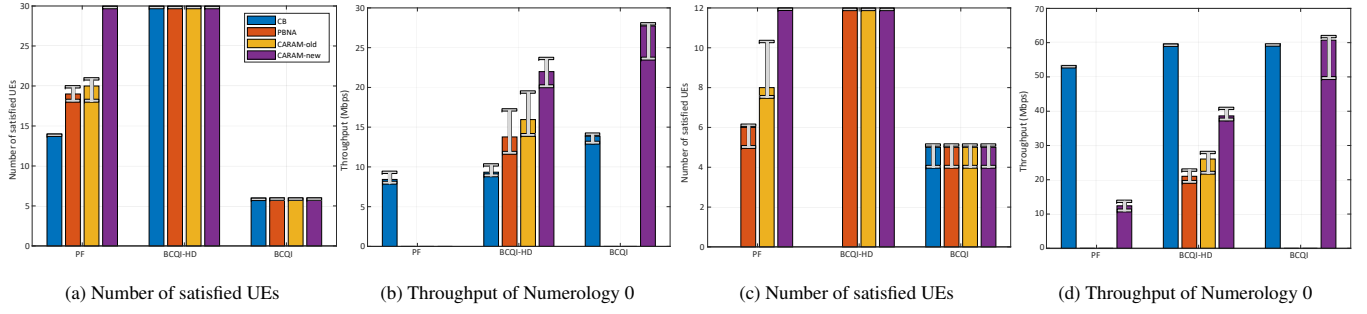


**Figure 2.11:** Stationary analysis with homogeneous GBR requirements and UE ratio 1:1:1. Medium load scenario with UEs arrangement 1.



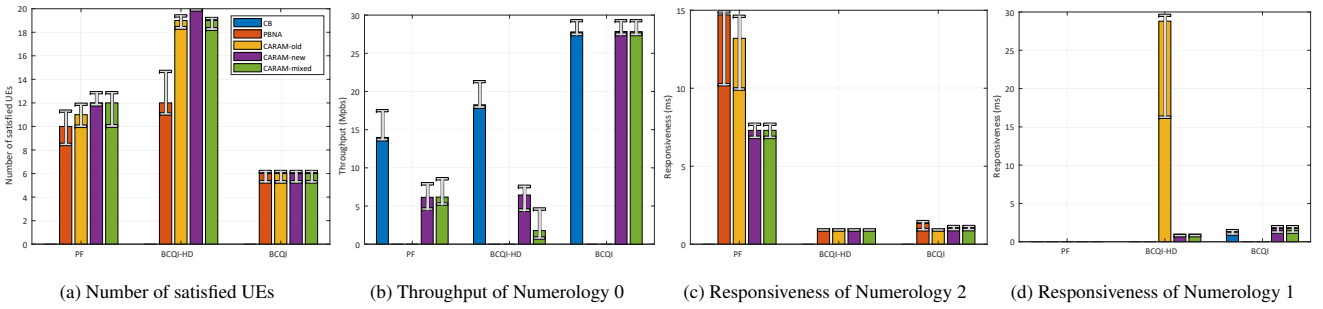
**Figure 2.12:** Stationary analysis with homogeneous GBR requirements and UE ratio 1:1:1. High load scenario with UEs arrangement 2.

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.13:** Stationary analysis with homogeneous GBR requirements. (a)-(b) UE ratio 1:1:4 and medium load scenario with UEs arrangement 1. (c)-(d) UE ratio 4:1:1 and high load scenario with UEs arrangement 2.

### 2.6.2 Stationary analysis with heterogeneous GBR requirements



**Figure 2.14:** Stationary analysis with heterogeneous GBR requirements.

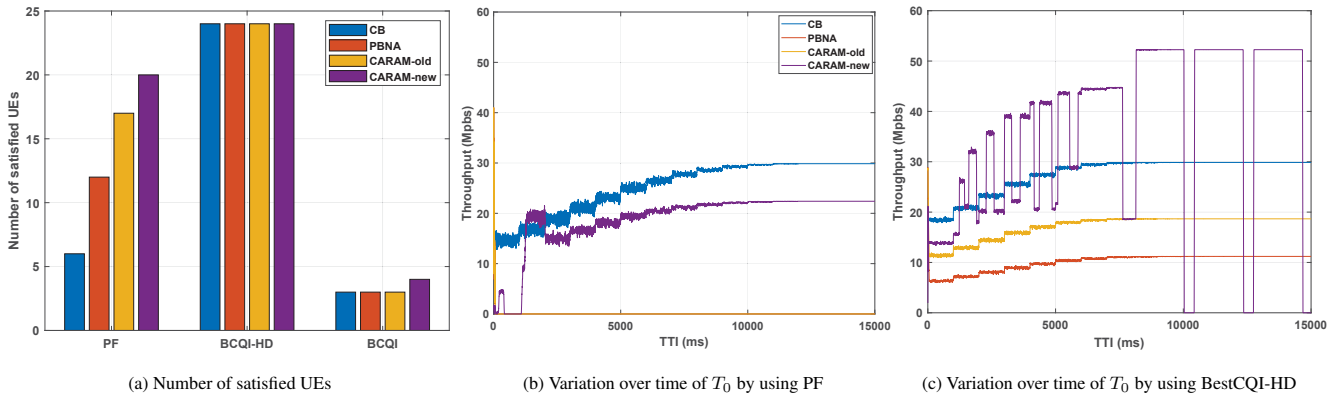
In this Section, we analyze the performance in a stationary scenario where all UEs are not in motion, and each GBR service requires a data rate  $R_{i,x}$  uniformly extracted in  $\{2, 5\}$  Mbps. The UE ratio is 1:1:1. The aim of this simulation setup is to highlight the advantages introduced by the new dropping function explained in Section 2.5.1. For this reason, we compare the performance of the proposed framework not only with CB, PBNA, and the CARAM-old, but also with a modified version of the CARAM-new, where only the dropping is the same of the CARAM-old version. We term this hybrid version as CARAM-mixed, and the performances are reported in Fig. 2.14. As shown, the introduction of this enhanced dropping function allows the CARAM-new with BCQI-HD to satisfy one more GBR UE and, at the same time, to almost triplicate the throughput of the BE Numerology. This occurs since the old dropping function aims to drop the GBR UE experiencing the lowest throughput, regardless of its GBR requirement. As a result, it may drop a UE re-

quiring 2 Mbps, experiencing a throughput equal to 1.8 Mbps, and that needs a single additional PRB to be satisfied, rather than a UE requiring 5 Mbps, experiencing 2.5 Mbps, and that needs 4 additional PRBs.

We note that the same improvement is not evident by adopting the PF, since the final result is to drop all UEs of Numerology 1, and also by adopting BCQI, due to the intrinsic limits of this 2nd level scheduler.

### 2.6.3 Dynamic analysis with homogeneous GBR requirements

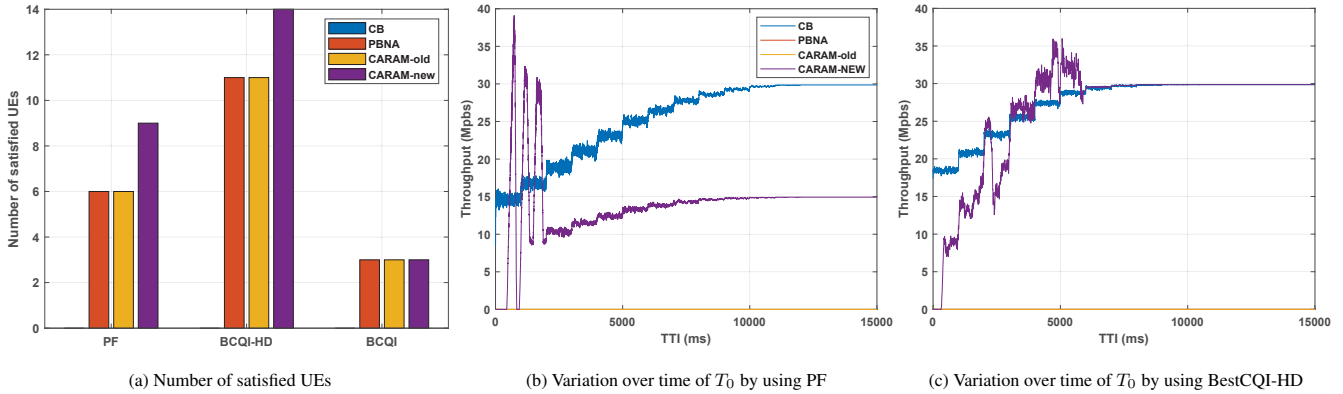
In order to assess the efficacy of new "decrement function", we analyze the performance of the proposed framework in a dynamic scenario where all UEs are in motion, and all the GBR services have the same GBR requirement. The UE ratio is 1:1:1. Since significant changes in the channel condition require a long simulation time, we force the simulation setup so that the CQI values reported to the gNB are increased by one unit every 1000 TTIs. Specifically, in the first TTI  $l = 1$ , all UEs are uniformly distributed in  $[2, 5]$  km. Starting from the second TTI, they are in motion at suitable constant speed towards the base station.



**Figure 2.15:** Dynamic analysis of a single simulation test with homogeneous GBR requirements. Medium-high load scenario.

Also in this analysis we introduce two different traffic load conditions, i.e., a medium-high load scenario with  $R_{i,x} = 1$  Mbps for each UE  $i \in \mathbf{U}_x$ , with  $x \in \mathbf{X}_{\text{GBR}}$ , and an overload scenario where  $R_{i,x} = 3$  Mbps. The results of a single simulation test related to the first traffic load condition are depicted in Fig. 2.15, where Fig. 2.15a reports the number of satisfied UEs, and Figs. 2.15b-c the variation over time at a TTI-time step of the Numerology 0 throughput by adopting PF, and BCQI-HD as

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.16:** *Dynamic analysis of a single simulation test with homogeneous GBR requirements. Overload scenario.*

second level scheduler, respectively. The number of satisfied UEs, under different configurations of 1st and 2nd level schedulers, is in line with the one obtained for the stationary scenario with medium load (see Fig. 2.11a). The only discrepancy is that, by adopting the PF, the number of satisfied UEs in Fig. 2.15a is slightly lower since the traffic load conditions are more demanding.

Now, we focus on the throughput of the BE Numerology when the PF is adopted (see Fig. 2.15b). As expected, the throughput achieved by adopting the CB is almost constant for each time interval of 1000 TTIs, and increases every 1000 TTIs due to the enhanced channel qualities. This behavior is valid until the channel conditions are the best achievable. Regarding PBNA and CARAM-old, the throughput is very high in the first TTI, since the number of SBs allocated to the GBR numerologies is the minimum needed to serve all UEs by assuming they are all experiencing the best possible channel condition. As simulation time increases, the number of SBs reserved to the BE Numerology decreases in order to serve all GBR services, and this trend occurs until the system reaches the steady state, giving all the SBs to the GBR Numerologies. As a result, after a few TTIs, the BE throughput becomes zero and does not change over the simulation time. As concerns the proposed CARAM-new in the first TTI the throughput of the BE Numerology is lower because the algorithm estimates, in a more accurate fashion, the number of SBs to allocate to the GBR numerologies. However, as the simulation time increases, even in this case, the throughput initially reaches zero. Since the

system is not able to satisfy all GBR services, some of them are dropped, and thus, before the system reaches a new steady state, some SBs are allocated to Numerology 0 (as shown in the interval [1, 1000] TTIs). After 1000 TTIs, since the channel quality increases, the CARAM-new reduces the number of SBs allotted to the GBR Numerologies, thanks to the decrement strategy provided in Section 2.5.1 and, consequently,  $T_0 \neq 0$  for all the remaining simulation time.

As regards the BE throughput by adopting the BCQI-HD, we report its variation over time in Fig. 2.15c. Unlike the previous scheduler, all the GBR services are satisfied in each configuration. So, the behavior is similar to the one described for the PF scheduler with the difference that in the steady state  $T_0 \neq 0$  for each 1st level algorithm. The only difference is the floating behavior of the proposed CARAM-new. It occurs since the granularity of the SBs allocation is low, and consequently, the steady state may be, as in this case, a non-static SBs dimensioning, where for several TTIs a large number of SBs are allotted to the GBR Numerologies while for other TTIs this value is low. This dynamic SBs allocation permits to obtain the optimal SBs dimensioning, on average, over time. We underline that only our scheme can provide this dynamic SBs allocation thanks to the joint increment and decrement strategy.

Hereinafter, in Fig. 2.16 we report the results of a single simulation test in the overload scenario. Also in this case the number of satisfied UEs, under the different configurations is similar to that obtained for the stationary scenario with high traffic load (see Fig. 2.12a). As a consequence, the throughput achieved by the BE Numerology is different from zero only when adopting CB and CARAM-new, for all the considered 2nd level schedulers (see Figs. 2.16b-c). In both the figures, our proposal shows an initial fluctuation in the throughput until it reaches the steady state with a constant SBs allocation.

In conclusion, all the simulation scenarios analyzed show that the best performances are achieved by adopting CARAM-new as 1st level allocation algorithm and BCQI-HD as 2nd level scheduling algorithm.

## **2.7 Related Works**

---

In this Section, we focus on other related works not reported in Section 2.4. Some of them focus on the first level, other ones on the second level, others still present a comprehensive approach which takes into account both levels.

As regards the latter category, in [39] the authors assume that the services requested by UEs are not mapped to the numerologies in a fixed manner. So, they consider a set of possible numerologies, where each numerology differs from the other ones in terms of sub-carrier spacing and cyclic prefix length. Then, they propose a heuristic algorithm to estimate which numerology should be allocated to each service, and to allocate the radio resources to UEs in order to maximize their average satisfaction degree. Despite the main idea being innovative, their results show that the improvements are significant only if the cardinality of the numerology set is high, at the cost of an increase in scheduling complexity and signaling overhead.

As regards the algorithms focusing on the first level, we present in this Section the Adaptive Numerology Resource Allocation (ANRA) [30] and the work proposed in [40]. Specifically, ANRA calculates a proper metric for each numerology  $x$  over each PRB, based on the average CQI values reported by UEs requiring a service belonging to numerology  $x$ . First, it considers the PRB where the maximum metric value is experienced (among all numerologies and PRBs) and assigns it to the related numerology. Second, it evaluates whether the highest metric value is also experimented by the same numerology on the adjacent PRBs. If so, it assigns them to this numerology. Otherwise, the algorithm temporarily excludes the previous numerology and repeats the above steps considering only the remaining numerologies. Finally, if there are still resources available, the scheduler re-admits all numerologies and performs all the steps again. The authors in [30] claim to propose a channel and QoS-aware allocation scheme for multi-numerology 5G networks aiming to achieve greater spectral efficiency, while maintaining fairness and guaranteeing users' QoS requirements. However, the solution therein proposed fails to achieve this goal in a realistic 5G scenario, characterized by stringent requirements and variable channel conditions. In fact, their

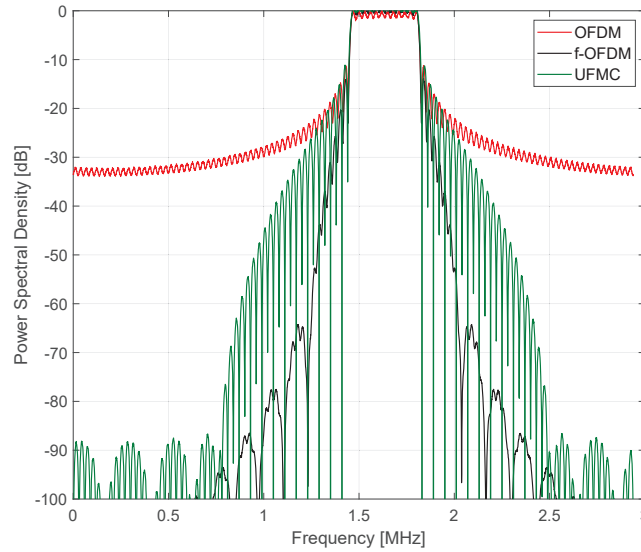
resource allocation procedure does not take into account any GBR, delay or priority requirements, nor does it adapt to the variation of channel and traffic conditions. In confirmation of this, in the performance evaluation section, the authors take into account just LTE type services for which only one PRB per UE is enough, static channel conditions, and the only performance parameter considered is the cell throughput. Finally, in [40], the authors propose a mixed-numerology scheduling scheme with a different goal. In fact, they aim to reduce the Peak-to-Average Power Ratio (PAPR) for both time-domain and frequency-domain numerology multiplexing scenarios. Their approach does not take into account the QoS requirement of services and the channel quality, but only how to best composite the signal of mixed-numerology so that the PAPR of the transmitted signal is minimized.

As regards the 2nd level scheduling algorithms, it should be noted that other scheduling algorithms with targets different from those considered in Section 2.4 are available in the literature. For example, the scheduling method proposed in [41] shows a new second level scheduling technique customized for multi-numerology scenarios. The aim is to improve the reliability of the system by protecting the users that may suffer more from the effect of the INI. In fact, the authors of [41] divide the bandwidth available for each numerology into an inner part (i.e., the central region), and an outer part (i.e., near the edge of the bandwidth, bordering on the bandwidth dedicated to a non-orthogonal numerology). Then, they allocate the PRBs belonging to the inner part to the users requiring a URLLC service and to the users at the cell edge, while the outer part to the other users. The aim of the proposed scheduler is to maximize the Signal-to-Interference Ratio (SIR) for the UEs at the cell edge and for URLLC services. However, it does not provide a new criterion for allocating the PRBs in the inner and the outer part of the available bandwidth, and does not take into account any GBR requirement.

### 2.7.1 New 5G NR waveforms

The OFDM waveform has been adopted as the transmission waveform in LTE and has been inherited in 5G NR, thanks to its robustness against fading phenomena and ease of implementation. However, it suffers from several drawbacks, including high peak-to-average power ratio (PAPR)

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.17:** *Out of band emission for the considered waveforms.*

and high Out-of-Band (OOB) emissions in the frequency domain [42], as shown in Fig. 2.17. In particular, the second drawback leads to large frequency guard bands needed at both the edges of the system bandwidth and also among different numerologies, for permitting the signal to reach enough attenuation and do not interfere too much. This implies a significant loss in spectrum efficiency. To overcome the above-mentioned limitations, a new waveform should be considered for the future stage of 5G. On one hand, the waveform candidate should inherit all the advantages of OFDM. On the other hand, the new waveform should be able to support flexibly configured numerologies and multi-numerology coexistence, to enable tailored services for different applications in heterogeneous scenarios [43]. Two important candidate waveforms among all those available are the f-OFDM and the UFMC that, as depicted in Fig. 2.17, show both a very reduced OOB emission in terms of Power Spectral Density.

### **filtered OFDM**

Filtered OFDM (f-OFDM) [26] is one of the most promising waveform that has been introduced to meet the new requirements introduced by 5G. In the f-OFDM, the assigned bandwidth is divided up into several portions, called subbands. In each subband, a conventional OFDM is inserted to suit the needs of certain type of service and the associated



## 2.8. THE VIENNA 5G Link Level (LL) SIMULATOR

---

channel characteristics. After that, a subband-based filtering is applied to suppress the inter-subband interference. In this way, f-OFDM is capable of overcoming the drawbacks of OFDM whilst retaining the advantages of it. First of all, with subband-based filtering, the requirement on global synchronization is relaxed and inter-subband asynchronous transmission can be supported. Secondly, with suitably designed filters to suppress the OOB, the GB consumption can be reduced to a minimum level. Third, within each subband, optimized numerology can be applied to suit the needs of certain type of services.

### Universal Filtered Multi-Carrier (UFMC)

The UFMC [27] proposes an alternative modulation scheme to OFDM, where a filtering operation is applied to a group of consecutive subcarriers according to different UEs' services requirements. The modulation technique processes these group of subcarriers individually. On each of those, an  $N$ -point Inverse Fast Fourier Transform (IFFT) is applied. Each subband passes through a filter, in order to reduce the OOB spectral emissions, and the responses from the different subbands are summed. Different filters can be applied, and one of the most used is the Dolph-Chebyshev. Moreover, unlike OFDM, the UFMC uses Zero-Postfix (ZP) instead of CP in order to avoid inter-symbol interference in a case of high delay spread channels. At the receiver side, the received signal is decomposed into two parts, named body and tail. First, the signal is transformed by copying the tail to the beginning of the signal, and then it is performed a  $2N$  points Fast Fourier Transform (FFT) on it. As result, the performance of UFMC is better than OFDM in terms of both OOB emissions and PAPR.

## 2.8 THE VIENNA 5G Link Level (LL) SIMULATOR

---

Various free simulation tools that allow detailed analysis of communication links and physical features are available. However, neither of these platforms support LL simulation of modern mobile communications networks, that is, 5G and beyond [44]. In particular, these new generation networks need the establishment of multiple communication links among BSs and UEs, and the implementation of new scheduling al-

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

gorithm utilizing both multi-numerology and new complex waveforms. The 5G Vienna LL simulator [29] offers a unifying platform for performance evaluation as well as co-existence investigation of candidate 5G physical layer schemes. It permits to provide great flexibility by supporting a broad range of simulation parameters. Therefore, many different combinations of physical layer settings are comparable. To enable investigation of 5G physical layer candidate methods, they support features such as new waveforms like filtered or windowed OFDM, UFMC, and different channel codes like Turbo coding, Low Density Parity Check (LDPC) coding or Polar coding. All of these schemes support any combination of channel coding rate and Quadrature Amplitude Modulation (QAM) alphabet size, which results in many different Modulation and Coding Scheme (MCS), and consequently Channel Quality Information (CQI) values. In addition, the employed physical layer schemes can be different for users of different cells such that their co-existence can be simulated. However, it does not intrinsically provide the coexistence of multiple numerologies for the same BS. In fact, it permits to use only one numerology, waveform and channel coding method per Base Station. Moreover, this is realized without implementing an underlying cellular geometry. There is neither a physical cell size nor a distance to the user being considered; the path-loss to a user is rather specified as an input parameter, leading to an average signal to noise ratio. As regards the system performances, they can be reported in terms of System throughput, throughput per UE, Bit Error Ratio (BER) or Frame Error Ratio (FER).

### 2.9 System Model and QoS-aware RRM framework

---

We consider a downlink 5G NR network where the adopted resource grid of bandwidth  $BW$  allows the coexistence of multiple numerologies in a frame and also allows the numerology to be adapted dynamically. There is a single cell equipped with one transmit antenna, and  $N_{UE}$  UEs, each equipped with one receive antenna. We assume that each UE requires only one service belonging to a numerology  $x$ , that characterizes by both the priority  $p_x$  and the GBR requirement  $R_x$ . We define the following vectors:

- $\mathbf{X} \subseteq \mathbf{X}_T$  with cardinality  $|\mathbf{X}| = N_x$ , as the vector containing each

numerology  $x$ ;

- $\mathbf{X}_{\text{GBR}} \subseteq \mathbf{X}$  as the vector containing the GBR numerologies sorted in descending order of priority.
- $\mathbf{X}_{\text{BE}}$  as the vector containing the Best Effort (BE) numerologies.
- $\mathbf{U} = [1, \dots, N_{UE}]^T$ , as the vector containing all UEs.
- $\mathbf{U}_x$  as the sub-vector of  $\mathbf{U}$  containing all UEs requesting a service belonging to numerology  $x$ .
- $\mathbf{V} = \{(p_1, R_1), (p_2, R_2), \dots, (p_{N_x}, R_{N_x})\}$  as the tuple containing all service priorities and the GBR requirements.
- $\mathbf{C} = [CQI_i, \dots, CQI_{N_{UE}}]$  as the vector containing the wideband CQI values reported by each UE.

For each TTI, first the gNB appropriately allots the radio resources to each numerology  $x$ . To do this efficiently, we determine the minimum radio resource unit that can be allocated to each numerology, and term it as Subband (SB). One SB occupies  $z$  in terms of bandwidth and lasts one TTI in the time domain. In the example shown in Fig. 2.2,  $z = 720$  kHz, one TTI lasts 1 ms, and, consequently, the number of PRBs per SB, denoted as  $N_{RB}$  is equal to 4, for each numerology. The total number of available SBs is  $N_{SB} = \lfloor \frac{BW}{z} \rfloor$ , and we define  $\mathbf{K}$  as the vector containing all subbands. Then, we consider  $g$  as the GB size that can be inserted between SBs that belong to different numerologies. For the sake of simplicity, we assume the same value of  $g$  independently from the subcarrier spacings which the numerologies belong to. We define  $\mathbf{A} \in \mathbb{B}^{N_x \times N_{SB}}$  as the assigned SBs matrix, where  $a_{x,k} = 1$  if SB  $k \in \mathbf{K}$  is assigned to the numerology  $x \in \mathbf{X}$ , otherwise  $a_{x,k} = 0$ . For each  $x \in \mathbf{X}$ , we define also  $\mathbf{K}_x \subseteq \mathbf{K}$  as the vector containing all SBs allotted to numerology  $x$ . Then, the PRBs of each numerology are assigned to the related UEs. Since each SB corresponds to  $N_{RB}$  PRBs, the number of PRBs that can be allocated to UEs is  $N_{RB} \cdot |\mathbf{K}_x|$ . So, we define  $\mathbf{A}_x \in \mathbb{B}^{|\mathbf{U}_x| \times (N_{RB} \cdot |\mathbf{K}_x|)}$  as the assigned PRBs matrix, where  $a_{i_x, r_x} = 1$  if and only if PRB  $r_x \in \{1, \dots, (N_{RB} \cdot |\mathbf{K}_x|)\}$  is assigned to UE  $i_x \in \{1, \dots, |\mathbf{U}_x|\}$ . Given the allocation matrices and the CQI vector

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems

---

C, the gNB knows the CQI values of the PRBs allocated to the UEs, and thus it can determine the Modulation and Coding Scheme (MCS) that will be used during the transmission. Specifically, we adopt the CQI-MCS Index mapping Table A.4-3 in [34]. Moreover, we adopt Table 5.1.3.1-1 in [33] to derive the modulation and the code rate related to each MCS Index. Then, the Transport Block Size ( $TBS_{i_x,t}$ ) related to UE  $i_x$  in the  $t$ th TTI, i.e., the number of information bits transmitted during TTI  $t$ , is estimated as reported in [33]. Finally, we can estimate the transmission rate of user  $i_x$  in the time of  $N_{TTI}$  TTIs as follows:

$$T_{i_x} = \sum_{t=1}^{N_{TTI}} \frac{TBS_{i_x,t}}{TTI \cdot N_{TTI}}. \quad (2.41)$$

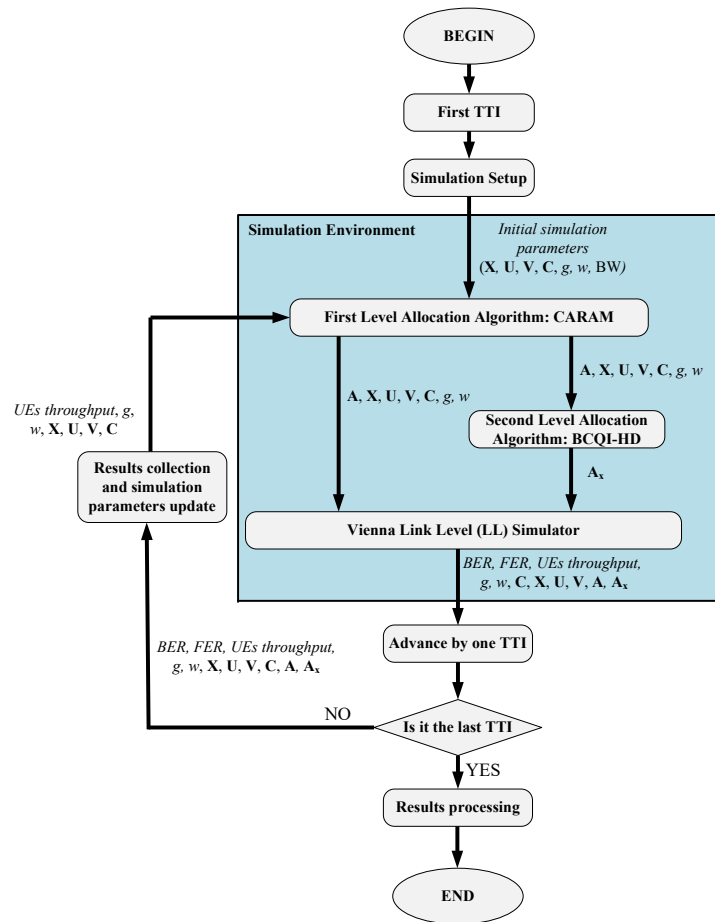
Given the transmission rate, we determine the satisfaction of UEs requesting a service belonging to  $\mathbf{X}_{\text{GBR}}$ . We set  $s_{i_x} = 1$  if  $T_{i_x} \geq R_x$ ,  $s_{i_x} = 0$  otherwise. Moreover, for the computation of  $T_{i_x}$  we adopted both the standard OFDM waveform, and the new candidate ones, i.e., f-OFDM and UFMC. So, for shortness we identifies the metric  $\psi$  measured by utilizing the different waveforms as  $\psi^{[O]}$ ,  $\psi^{[f]}$  and  $\psi^{[U]}$ , respectively. Finally, we define the set of considered waveforms as  $\mathbf{W} = \{O, f, U\}$ , with  $O = \text{OFDM}$ ,  $f = \text{f-OFDM}$ , and  $U = \text{UFMC}$ .

### 2.10 The Proposed Simulation Environment for the Multi-Numerology scenario

---

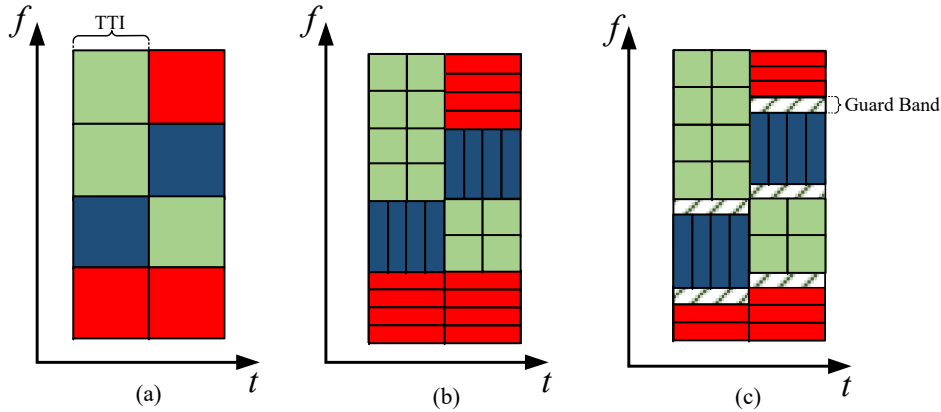
In this Section, we explain how the Simulation Environment is constituted, and the adaptation that has been made for permitting the correct functioning among the several parts of the Environment. The first part takes into account the problem of defining the most appropriate RRM scheme. Since we want to define a complete and in-depth Environment for a multi-numerology system, we need to consider an RRM scheme that can take into account the requirements of services related to each numerology and the channel conditions perceived by UEs requesting these services. To the best of our knowledge the only control framework that respects these requirements is the one proposed in [45]. By following this control framework, the problem of allocating the radio resources is divided into two levels. The first one is the CARAM, that permits to

## 2.10. The Proposed Simulation Environment for the Multi-Numerology scenario



**Figure 2.18:** Flow chart of the proposed Simulation Phase

calculate the Assigned SBs matrix  $\mathbf{A}$ . Conversely, as regards the second level, i.e., the computation of the Assigned PRBs matrix  $\mathbf{A}_x$ , the chosen algorithm is the one proposed in [25], named BestCQI\_Highest Deviation (BestCQI\_HD), that is, is a promising channel-aware and QoS-aware scheduling algorithm. Finally, the outputs coming from both CARAM and BestCQI\_HD are given to the 5G Vienna LL simulator, that permits to emulate in depth the behavior of a radio communication link. In particular, the 5G Vienna LL Simulator permits to implement different waveforms and to emulate the desired channel model. The detailed description of the Simulation Environment, comprehensive of the different inputs and outputs values for each part of the Environment is reported in Fig. 2.18. The simulation can be performed for several TTIs, depending on the time of study that want to be conducted. For each simulated TTI, it is possible to save the output and to updated the simulation parameters.



**Figure 2.19:** Radio resource allotment for the different phases of the Simulation Environment

The radio resource allocation resulting from the different parts of the Simulation Environment are reported in Fig. 2.19. In the first one, the CARAM subdivides the overall bandwidth in SBs (see Fig. 2.19a), each one assigned to a numerology with a specific subcarrier spacing (See Fig. 2.19b). Then, by means of the BestCQI\_HD the PRBs of each numerology are assigned to the corresponding UEs. Finally, Fig. 2.19c reports how the subdivided bandwidth is allotted including the introduction of the GB, i.e.,  $g$ . Regarding the latter phase, the Vienna 5G LL simulator needs an ad hoc approach to properly insert the GB.

### 2.10.1 Adaptation of the 5G Vienna LL Simulator scheduling for the multi-numerology scenario

As stated in Section 2.8, the considered LL simulator does not support the multi-numerology scenario. So, for implementing our Simulation Environment and to properly add the GB among the different numerologies, we designed the trick depicted in Fig. 2.20. In particular, in the LL simulator we define as many base stations as the number of considered Numerologies, i.e.,  $N_x$ . Each base station has its own resource grid of bandwidth  $BW$ , and the overall allocated bandwidth is the same as the one allocated by the RRM scheme. However, in the LL simulator we are adding to the RRM scheduled band an additional portion of bandwidth as GB. Consequently, in the LL simulator it is not possible to allocate all the SBs contained in  $\mathbf{A}$ . So, having  $\mathbf{A}$  and known the information  $g$  and

## 2.10. The Proposed Simulation Environment for the Multi-Numerology scenario

**Table 2.2:** *Parameters of the case study*

Parameter	Value
Carrier Frequency	2 GHz
gNB Antenna Gain	18 dBi
UE Antenna Gain	0 dBi
UE Noise figure	7 dB
gNB Transmit Power	46 dBm
BW	10 MHz
Number of subbands	$N_{SB} = 12$
GBR numerologies	$\mathbf{X}_{GBR} = \{2, 1\}$
GBR requirements	$R_1 = R_2 = 450$ Kbps
BE numerology	$\mathbf{X}_{BE} = \{0\}$
UE ratio ( $ \mathbf{U}_0  :  \mathbf{U}_1  :  \mathbf{U}_2 $ )	1:1:1
TTI length	1 ms
Number of TTIs for which the simulation is run	$N_{TTI} = 1000$

BW, we calculate the effective number of allocable SBs by the Vienna 5G LL simulator, termed  $N_{ASB}$ . If we assume  $h$  as the number of times one numerology changes to another, it simply follows that

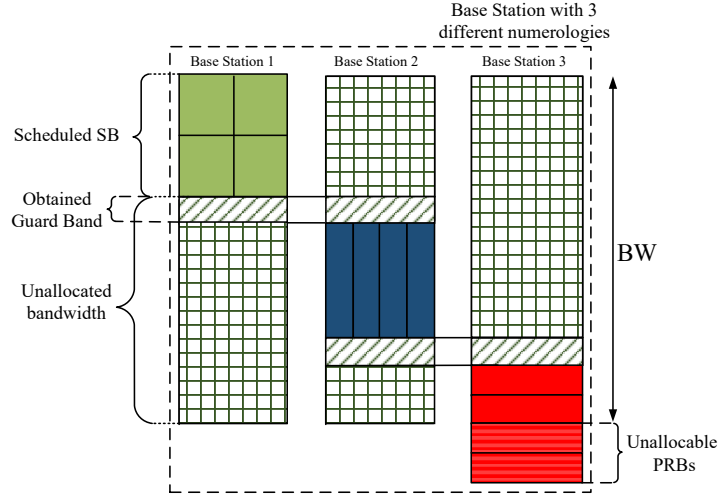
$$N_{ASB} = \left\lfloor \frac{BW - hg}{z} \right\rfloor. \quad (2.42)$$

However, since the LL simulator permits to schedule with the granularity of one PRB, it is important to calculate also how many PRBs cannot be allocated. Under the assumption that the portion of bandwidth dedicated to the GBs is subtracted from one numerology  $x^*$ , the number of unallocable PRBs for the numerology  $x^*$  is

$$N_{UPRB,x^*} = \left\lceil \left( \frac{BW}{b_{x^*}} - \frac{BW - hg}{b_{x^*}} \right) \frac{TTI}{T_{s,x^*}} \right\rceil. \quad (2.43)$$

Since we are considering a multi-numerology scenario where the services are categorized into GBR and BE, with the aim of preserving the quality of the GBR services, we choose as  $x^*$  the numerology related to the BE services. In detail, if more than one numerology is reserved for the BE services, then  $x^*$  is the BE numerology that minimizes (2.43). For example, in Fig. 2.20 the  $x^*$  numerology is the Numerology 0, that corresponds to  $\Delta f = 15$  KHz. Considering  $BW = 10$  MHz,  $g = 180$  KHz, and

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.20:** Adaptation procedure for the GB insertion in the proposed Simulation Environment

$h = 2$ , it follows  $N_{UPRB,0} = \left\lceil \left( \frac{10 \text{ MHz}}{0.18 \text{ MHz}} - \frac{10 \text{ MHz} - 2 \cdot 0.18 \text{ MHz}}{0.18 \text{ MHz}} \right) \right\rceil \frac{1 \text{ ms}}{1 \text{ ms}} = 2$ . Consequently, the number of PRBs that can be scheduled by the 5G Vienna LL simulator, termed  $N_{SPRB}$  is

$$N_{SPRB} = \sum_{x \in \mathbf{X}} \frac{TTI}{T_{s_x}} |\mathbf{K}_x| - \min_{x^* \in \mathbf{X}_{BE}} N_{UPRB, x^*}. \quad (2.44)$$

After having computed how many PRBs can be allocated, following the information reported in the  $\mathbf{A}_x$  matrix and considering the descending order of priority for the numerologies, we accordingly allocate the services that can be scheduled in the reserved bandwidth, i.e.,  $N_{SPRB}$ .

Finally, to virtually place the all UEs in such a way that they are scheduled in the resource grid of the same BS, the approach is to adjust the BSs' transmit power so that the power received by the UE from the proper Base Station in the scheduled PRBs (i.e., direct link) is the same as the interference links' power (i.e., the power received from all the other BSs in the same scheduled PRBs). By doing this, the interference channels and the direct channels become indistinguishable.

### 2.11 Case Study Analysis

Considering a multi-numerology scenario with specific and diversified requirements, as defined in Section 2.9, we analyze a Case Study which aims to determine the best  $(w, g)$  pair that allows on the one hand to



**Table 2.3:** Number of SBs per numerology versus (waveform, GB size) pair.

GB size	OFDM			UFMC			f-OFDM		
	N2	N1	N0	N2	N1	N0	N2	N1	N0
0 kHz	3	3	6	3	3	6	2	2	8
60 kHz	3	3	6	2	3	7	2	2	8
180 kHz	2	2	8	2	2	8	1	1	10
360 kHz	1	2	9	1	2	9	1	1	10
540 kHz	1	2	9	1	1	10	1	1	10
720 kHz	1	2	9	1	1	10	1	1	10

adequately reduce the INI phenomenon and, on the other hand, to pursue the following objectives.

The first one is to maximize the number of satisfied users, taking into account service priorities, i.e.,

$$\max_{g \in \mathbf{G}, w \in \mathbf{W}} \left\{ \sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[1]}|} s_{i_{\mathbf{X}_{\text{GBR}}[1]}}^{[w]} + \sum_{x=2}^{|\mathbf{X}_{\text{GBR}}|} \sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[x]}|} s_{i_{\mathbf{X}_{\text{GBR}}[x]}}^{[w]} \cdot \prod_{y=1}^{x-1} \left[ \frac{\sum_{i=1}^{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[y]}|} s_{i_{\mathbf{X}_{\text{GBR}}[y]}}^{[w]}}{|\mathbf{U}_{\mathbf{X}_{\text{GBR}}[y]}|} \right] \right\}. \quad (2.45)$$

The second objective is to maximize the throughput of the Best Effort numerology, i.e.,

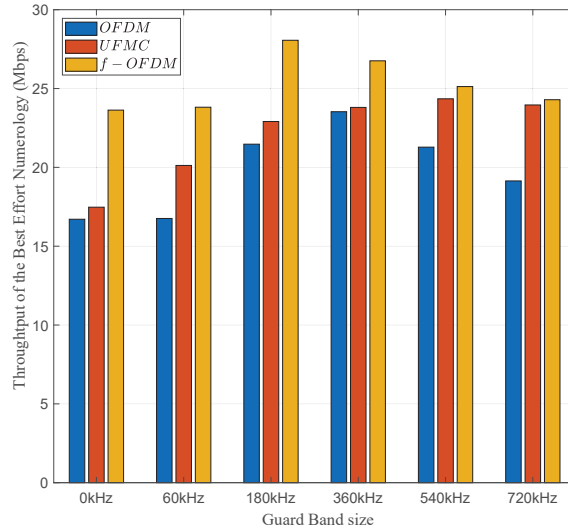
$$\max_{g \in \mathbf{G}, w \in \mathbf{W}} \sum_{x \in \mathbf{X}_{\text{BE}}} \sum_{i=1}^{|\mathbf{U}_x|} T_{i_x}^{[w]}. \quad (2.46)$$

In summary, starting from the initial parameters  $\mathbf{V}$ ,  $\mathbf{U}$ , and  $\mathbf{X}$ , the objective of this investigation is to find the optimal  $(g, w)$  pair, so that, by adopting the matrix  $\mathbf{A}$  and the matrices  $\mathbf{A}_x$ ,  $\forall x \in \mathbf{X}$  the objective functions (2.45) and (2.46) are maximized.

### 2.11.1 Simulation Configuration and Parameters

The simulation campaign involves the use of a cell with three different numerologies, i.e.,  $\mathbf{X}_{\text{T}} = \{0, 1, 2\}$ . The services with GBR requirement are mapped into  $\mathbf{X}_{\text{GBR}} = \{2, 1\}$ , and the priority order is  $p_2 > p_1$ . We consider  $N_{\text{UE}} = 12$  and the GBR requirement  $R_1 = R_2 =$

## Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems



**Figure 2.21:** *Throughput of the Best Effort Numerology vs different waveforms and guard band sizes.*

450 Kbps for all the services of the  $\mathbf{X}_{\text{GBR}}$  numerologies. The simulation has been carried out by adopting all possible (waveform, GB size) pairs, considering all waveforms in  $\mathbf{W}$ , and all guard band sizes in  $\mathbf{G} = \{0, 60, 180, 360, 540, 720\}$  kHz. The channel type is set to Pedestrian-A in which the speeds of individual users can vary from a minimum of 0 up to a maximum of 10 km/h. UEs are placed near the BS, so that the experienced Signal-to-Noise Ratio (SNR) is high, and consequently, the effects of INI dominates over the channel noise, that is the focus of our analysis. The other simulation parameters are reported in Table 2.2.

### 2.11.2 Performance Analysis

In the simulation setup considered, the RRM scheme allocates, for each GBR numerology, the minimum quantity of SBs necessary to satisfy all its UEs. The number of SBs depends almost exclusively on the INI perceived by UEs, which may require a greater quantity of PRBs to compensate for the low channel quality. Finally, the remaining SBs are allocated to the BE numerology.

In Table 2.3, for different (waveform, GB size) pairs, we report the amount of SBs allocated to each GBR numerology (denoted as  $N_2$  and  $N_1$ ) and to the BE numerology ( $N_0$ ). Therefore, we focus on the results that promise greater efficiency, i.e., those that require fewer SBs to

satisfy the GBR services. As shown, both the UFMC and the f-OFDM waveforms reach this optimal result. Specifically, the UFMC reaches it for  $g \in \{540, 720\}$  kHz, while the f-OFDM for  $g \in \{180, 360, 540, 720\}$  kHz.

Since the first goal is achieved for several  $(w, g)$  pairs, we evaluate for them the second goal (2.46). In this regard, in Fig. 2.21, we report the throughput of the BE numerology vs different waveforms and guard band sizes. As shown, the best result is achieved by adopting f-OFDM as waveform and 180 kHz as GB size.

## 2.12 Conclusion

---

In this Chapter, we addressed the problem encountered by a gNB in managing the radio resources in a downlink multi-numerology based 5G NR system. We assumed that UEs require either a GBR service with a given priority level or a BE service. Depending on the priority level, each service is mapped into a specific Numerology and may require a different GBR value. By following the higher layer configured sub-band CQI reporting method, we proposed a new channel-aware control framework consisting of two levels. The first one, termed CARAM-new, assigns SBs to the Numerologies with the goal of maximizing number of satisfied GBR services while respecting the service priority and the throughput of the BE numerology. Moreover, in order to reduce the INI phenomenon, CARAM-new tries to assign continuous SBs to the same Numerology. The second level assigns the PRBs allotted to each Numerology to UEs by exploiting the PF, BCQI, or BCQI-HD scheduling algorithms. The proposed CARAM-new is an extended version to our previous CARAM [1], that includes the following enhancements: an improvement of the dropping strategy, the introduction of a decrement strategy for the allotted SBs, an updated verification of the GBR requirement taking into account the DAW, and an adaptation to properly work with different 2nd level scheduling algorithm. We conducted a simulation campaign under different traffic loads and UEs arrangements, in both static and dynamic conditions. The CARAM-new performed considerably better than its previous version and other intuitive 1st level algorithms, in terms of both number of satisfied GBR services, according

## **Chapter 2. A QoS-aware and Channel-aware framework for Multi-numerology OFDM/UFMC/filtered OFDM Systems**

---

to the priority, and BE Numerology throughput. Moreover, the results highlighted that the best performances are achieved by adopting the pair CARAM-new and BCQI-HD as 1st and 2nd level algorithms, respectively. Driven from this good results, we additionally developed a new Simulation Environment for the analysis of multi-numerology scenarios with diversified QoS requirements. The proposed Simulation Environment exploits the interaction between two parties. On the one hand, the RRM control framework executes the procedures on the basis of both service requirements and the channel conditions experienced by the UEs. On the other hand, the LL simulator, named Vienna 5G LL simulator, permits to add and vary several physical features, such as the waveform and the guard band size. Moreover, starting from the defined Simulation Environment, we evaluated a Case Study, aiming to determine the best (waveform, GB size) pair that allows to reduce the INI phenomenon, in a given multi-numerology scenario. In the considered scenario the adoption of f-OFDM as waveform and 180 kHz as GB showed the best performance. We underline that this result comes from a single case study, but the powerful Simulation Environment here presented can be used for many other Case Studies adopting a multi-numerology scenario.

---

# CHAPTER 3

---

## Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

### 3.1 Overview

---

In this Chapter, we address the problems of traffic prediction, resource allocation, and en/decoding strategies for SCMA-based mMTC scenarios. As reported in the Introduction Chapter, this scenario is characterized by a large number of low-complexity and energy constrained MTC devices sending periodically very short packets with relaxed delay requirement, without or with very little human interventions. In general, in the mMTC scenario, because of the relaxed delay requirements, a contention-based Random Access (RA) procedure can be adopted. However, a huge number of MTC devices might simultaneously initiate their procedure to get access to the network causing a severe congestion problem in the system if conventional LTE RA procedure is adopted [46, 47]. Currently, in LTE / LTE-Advanced (LTE-A) uplink communication, the MTC device firstly performs a contention-based RA procedure by transmitting a preamble in the Physical Random Access Channel (PRACH). Then, if the RA attempt has been successful and there are available radio re-

### **Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios**

---

sources into the Physical Uplink Shared Channel (PUSCH), the MTC device transmits its data in it. Otherwise, the MTC device re-attempts the access procedure in successive RA cycles for a maximum number of times. When thousand or more devices simultaneously request to send data, massive RA attempts increase the preamble collision probability, thus decreasing the number of successful accesses. This problem could be alleviated by increasing the resources allocated to PRACH. However, due to limited uplink resources, the higher the amount of resources available in the PRACH, the lower the ones available in the PUSCH. Therefore, many MTC devices which have successfully completed their RA attempt, could not find enough transmission resources in the PUSCH.

As regards the transmission efficiency in the PUSCH resources one way is the adoption of a Non Orthogonal Multiple Access (NOMA) technique. In fact, by adopting a NOMA technique, more than one MTC device can share the same time-frequency resource by mainly exploiting the power-domain (e.g., PD-NOMA), or the code-domain (e.g., Sparse Code Multiple Access, SCMA). In the SCMA, two or more MTC devices are superposed to the same radio resource by adopting sparse multidimensional codewords. However, there are two major challenges for the adoption of the SCMA in practical networks. First, the design of optimal codebooks at the encoder side and, second, the implementation of efficient decoding algorithms at the receiver side. As regards the encoder side, in [48] the authors design a mother multi-dimensional constellation with a desired Euclidean distance profile to reduce the probability of errors at the receiver, and then apply specific unitary phase rotations to generate layer-specific codebooks. As regards the receiver point of view, the optimum multi-user detection problem can be solved by finding either the Maximum A Posteriori (MAP) [49] probability mass function (pmf) of all layers' transmitted symbols or the Message Passing Algorithm (MPA) [50]. The main weakness is that these optimum detectors suffer from huge complexity. To overcome this issue, in [51] the authors present a log domain implementation of the MPA, called Log-MPA, aiming to obtain a SCMA decoding algorithm with less computational complexity at the expense of a slight loss in the decoding performance. However, the complexity of this detector is still very high.

As regards the enhancement of the PRACH resources, i.e., to address

the problem of managing the access of a massive amount of RA attempts, in literature, several works (e.g., [52–58]) take into account the limit of PRACH resources and introduce further RA control strategies to decrease the preamble collision probability. These works are mainly based on Access Class Barring (ACB) schemes [52–54], dynamic Backoff Indicators (BIs) [55, 56], Contention Tree Algorithms (CTAs) [57], or on a mixture of them [58].

In the conventional ACB scheme [59], the Base Station (BS) periodically broadcasts an ACB factor  $p \in [0, 1]$  to the MTC devices. Every MTC device having data to be transmitted draws a uniform random number  $q \in [0, 1]$ , and transmits the preamble only if  $q \leq p$ . The main challenge of these ACB schemes is to dynamically adapt the ACB factor  $p$  according to the traffic load in the PRACH. In fact,  $p$  should be a small value in the case of bursty and heavy-loaded scenario in order to relieve congestion, while it should be a large value in the case of light traffic condition in order to efficiently use the uplink resources and do not delay inappropriately the access attempts. In [60] the authors propose a dynamic adaptation of the ACB factor  $p$  based on the estimation of the MTC devices that are in backoff, i.e., those which will re-attempt their RA procedure. In [61] the authors design a Q-learning algorithm to dynamically tune the ACB factor  $p$  such that it can rapidly react to the traffic changes using local information available at the BS.

Regarding the schemes adopting dynamic BIs, the aim is to dynamically adjust the backoff window  $B_W$  on the basis of the number of re-attempting devices and the available preambles.

As concerns CTAs [57], they attenuate the PRACH collisions by properly organizing the retransmissions according to a  $f$ -ary splitting-tree algorithm with depth equal to  $d$ . Specifically, the available preambles  $L$  are sequentially indexed and equally divided into  $f$  branches. Each MTC device randomly selects one branch and transmits a preamble sequence belonging to that branch. In case of collision, if the number of RA attempts already performed ( $N_A$ ) is lower or equal to  $d$ , it will retransmit in a specific RA cycle. The wait before retransmitting is proportional to  $f$  and  $N_A$ . Otherwise, if  $N_A > d$ , the MTC device initializes a new tree for further retransmissions. Obviously, the value of both  $f$  and  $d$  need to be properly set by the BS according to the PRACH traffic load and  $L$ .

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

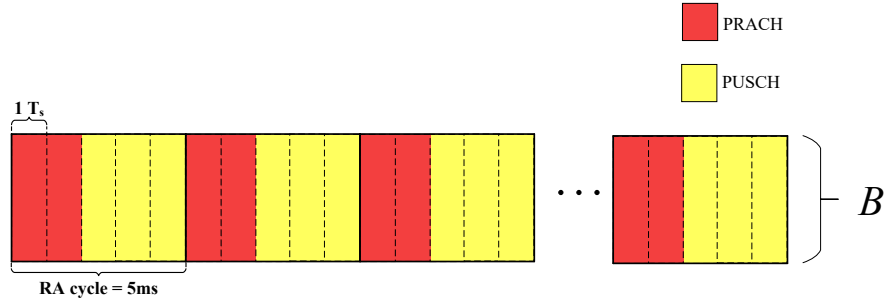
---

The hybrid scheme in [58] adopts jointly the principles of ACB, BI, and CTA. Specifically, the MTC device performs its RA attempt only after it has passed the ACB check. Then, if the RA attempt fails and  $N_A \leq d$ , it executes retransmissions according to the CTA. After that, if  $N_A > d$ , the MTC device will follow the random backoff according to the  $B_W$  value. However, in order to work well, all the works mentioned above and many other ones available in literature need to know the current PRACH traffic, denoted with  $M$ . It represents the total access requests per RA cycle, i.e., it is the sum of new access requests and the re-transmitting ones.

Despite these works aiming to maximize the number of successes in the PRACH, they do not consider the opportunity to dynamically dimensioning the uplink resources of the PRACH and the impact of the limited resources of the PUSCH. Other works (e.g., [62, 63]) focus exclusively on the optimal radio resource allocation between PRACH and PUSCH, based on the current traffic load. In [62] the authors, taking into account that the uplink resources are limited, propose a new control scheme which, before a RA cycle begins, allocates radio resources between PRACH and PUSCH, and broadcasts this configuration to all MTC devices. The main weakness of [62] is that the number of MTC devices attempting to access is considered well-known at the BS. In [63] we showed, by means of a long time evaluation (i.e., considering thousands of RA cycles), that the optimal PRACH/PUSCH resource allocation can be achieved only if a dynamic load-aware control is applied. So, in [63] we proposed a dynamic uplink resource dimensioning based on a predictive estimation of the total number of access attempts in the next RA cycle.

However, despite the proposal of a dynamic dimensioning, in both works no congestion control is applied in the PRACH, therefore the performances obtained can be significantly improved if PRACH/PUSCH resource allocation and access control are addressed together. To the best of our knowledge, only the authors of [64] did it. In that work, in each RA cycle, the optimal ACB factor, denoted as  $p^*$ , is derived together with a proper uplink resource allocation based on the hypothetical knowledge of the traffic load. For convenience, we term this control scheme  $ACB_{p^*}$ . Although  $ACB_{p^*}$  scheme has addressed the joint control of uplink re-





**Figure 3.1:** Uplink frame structure.

source dimensioning and ACB scheme in an mMTC scenario, it shows four significant weaknesses. First, the need to exactly know the amount of total access attempts to derive the proper resource dimensioning and the  $p^*$ . Second,  $ACB_{p^*}$  seems to achieve good performance only if a static evaluation in a single RA cycle is done, i.e., if the impact that the access re-attempts have on each MTC device is not taken into account. Third, since ACB factor  $p^*$  can be changed every RA cycle, the MTC device attempting to send a request has to remain continuously in the RX active state to listen the ACB factor, causing excessive energy consumption [65]. Fourth, the allocation problem is formulated in a generic SCMA-based network, without taking into account either the LTE or the 5G NR uplink frame structures.

Therefore, starting from [66], in this chapter we first investigate the access control problem in an mMTC usage scenario with a huge amount of MTC devices which sporadically send small-sized data, but in a highly synchronized manner [67]. We present a joint control of the RA in the PRACH and a new optimal dimensioning of uplink resources to be assigned to the PRACH and the PUSCH in a SCMA-based 5G NR. By exploiting the requirements of delay-tolerant services, our proposal is based on the clustering concept of successive RA cycles in order to properly spread the access re-attempts in time, in order to achieve two objectives. The first one is to reduce the preamble collision probability, thus increasing the successful access attempts. The second one is to increase the MTC device battery life, by activating, during the RA procedure, a discontinuous transmission that reduces the energy consumption.

Second, to further enhance the transmission performance in the advanced SCMA-based PUSCH resource allocation, we applied to a simple

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

dynamic resource allocation scheme [68] the feature of exploiting the unused PUSCH resources, if any, to serve also some MTC devices that have failed their access attempt. This system is termed as Enhanced Dynamic Uplink Resource Dimensioning (EDURD). More specifically, for a given RA cycle  $j$ , instead of limiting the number of attempting MTC devices by means of an ACB scheme, our idea was to serve a larger number of attempting MTC devices in  $j$ . This is achieved allowing the scheduler to allocate a pool of resources to a proper subset of collided preambles in a contention-based mode. This approach mitigates the congestion since it reduces the amount of re-attempting devices in the successive RA cycles.

Third, we address the problem of estimating the traffic load when the radio access follows a grant-based approach in an accurate and realistic way to improve the performances of all the works mentioned above and many other ones available in literature that need to know the current PRACH traffic load,  $M$ . We propose a new AI-based traffic load estimation method which exploits only the information related to the PRACH, available at the gNB. We chose a conventional Deep Neural Network (DNN) architecture used for regression tasks and we properly derived its topology for achieving the best accuracy with the lowest computational complexity.

Finally, to improve the benefits of the adoption of the SCMA in the PUSCH, we present a new SCMA encoder/decoder based on the Wasserstein Generative Adversarial Network (WGAN) [69]. It permits to select the optimal codebook and trains the decoder with the aim of obtaining good SER performance, under different  $E_b/N_0$  values.

## 3.2 Background

---

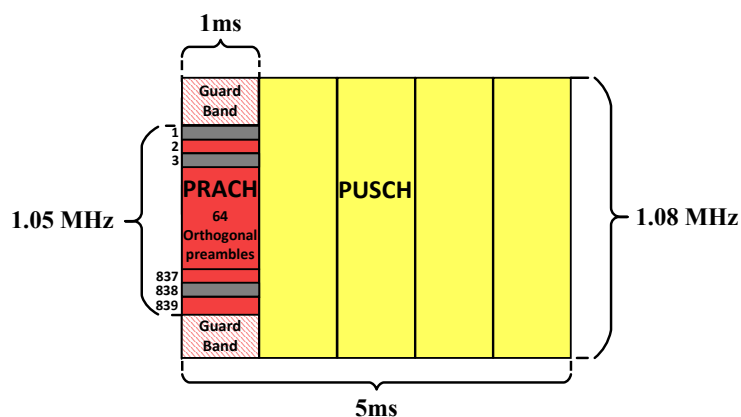
### 3.2.1 LTE Uplink

Fig. 3.2 shows a typical RA cycle of  $5 \text{ ms}$ <sup>1</sup>, where Uplink radio resources are divided into a PRACH subset and a PUSCH subset.

The PRACH is 1.08 MHz wide in frequency and has a time duration which depends on the RA Preamble Format (e.g., by using the com-

---

<sup>1</sup>In LTE systems different RA cycle lengths has been standardized based on the PRACH Configuration Index [3]. In particular, the typical PRACH Configuration Index ranges from 6 to 8, corresponding to RA cycle of 5ms [70, 71].



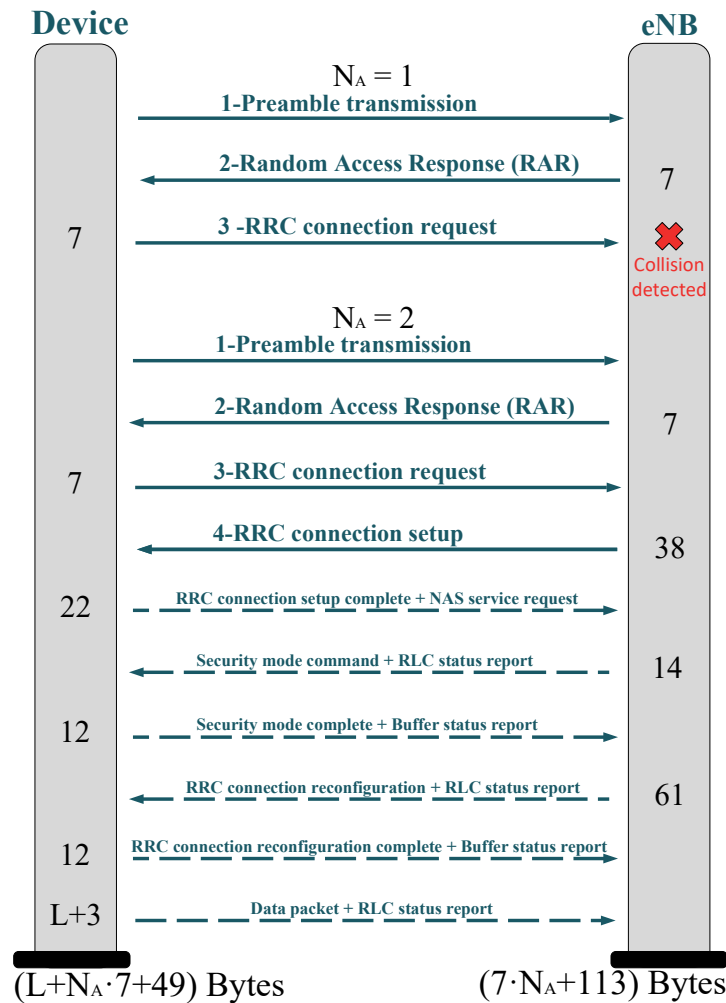
**Figure 3.2:** A typical RA cycle of 5ms.

mon Preamble Format 0, the PRACH is 1 ms in time). The PRACH access resources consist of 64 orthogonal preamble sequences which are mapped to 839 subcarriers of 1.25 kHz, and are generated from Zadoff-Chu sequences with zero-correlation zone [3]. The PUSCH consists of 72 sub-carriers of 15 kHz, is used to transmit user data and occupies the remaining available radio resources.

In the following, we describe the conventional contention-based RA procedure in the LTE/LTE-A system and show it in Fig. 3.3, in case the communication is successful on the second attempt. Before initiating their RA procedures, the devices receive periodically, from the Evolved Node B (eNB), the System Information Block (SIB) which contains, inter alia, broadcast information related to PRACH structure, the Backoff Window ( $B_W$ ) size for random backoff procedure, and the maximum number of access attempts ( $M_A$ ). The legacy RA procedure involves the following four-step message handshake between each device and the eNB.

*Step 1. Preamble transmission.* Each device randomly selects a preamble sequence out of the ones available for the contention-based procedure with equal probability and transmits it on the PRACH. Obviously, there is a non-zero preamble collision probability, since the same preamble can be selected by more than one device. Once the preamble transmission has been completed, the device increases its counter  $N_A$  by one.

*Step 2. RAR message.* After detecting the received preambles, the eNB can only detect whether a specific preamble has been transmitted or



**Figure 3.3:** Conventional 4-step RA procedure, when the communication is successful on the second attempt, and signaling messages for data packet transmission in LTE/LTE-A and eMTC technologies.

not, but it cannot recognize how many devices have transmitted it, i.e., if a collision has occurred. For each detected preamble, the eNB transmits the Random Access Response (RAR), which contains the Timing advance command, the temporary C-RNTI and the UL Grant. In particular, the latter field is of 27 bits, 18 of which are reserved for time and frequency resource allocation of the message transmission in Step 3 [72].

*Step 3. RRC connection request transmission.* After the RAR reception, each MTC device transmits on the scheduled radio channel the Radio Resource Control (RRC) connection request to set up the RRC connection with the eNB. However, if more than one UE has selected the same preamble in Step 1, they will receive the same RAR and send

their scheduled messages on the same radio channel, which makes the eNB hard to decode the received message correctly. In this case, the eNB recognizes a preamble collision.

*Step 4. RRC connection setup reception.* After correctly decoding the RRC connection requests, if there are available resources in the PUSCH for the device transmission, the eNB transmits the RRC connection setup message to the corresponding device.

We note that if the device does not receive either the RAR message or the RRC connection setup within the related predetermined time window, termed  $W_{RAR}$  and  $W_{CR}$ , respectively, then it reattempts the RA procedure inside the backoff window ( $B_W$ ) only if  $N_A \leq M_A$ .

If the 4-step RA procedure was successful, further signaling messages have to be exchanged between the UE and the eNB in order for the UE to initiate the transmission of the data packet. In detail, the device sends the RRC connection setup complete message and initializes the Non-Access Stratum (NAS) procedures to the Mobility Management Entity (MME), including security. Then, the UE will receive the RRC connection re-configuration message, which is used to establish the data radio bearer. Finally, it transmits the data packet.

### 3.2.2 eMTC and NB-IoT

Since the LTE is highly inefficient for supporting the MTC traffic characterized by sporadic infrequent transmission of small packets, the 3GPP has already introduced, in Release 13, a suite of two complementary technologies adapted for this type of traffic, denoted as enhanced Machine Type Communication (eMTC) and NarrowBand IoT (NB-IoT) [73]. Both technologies are optimized for granting lower complexity, and providing longer battery life, while seamlessly coexisting with other LTE services. As regards some technical specifications for the eMTC, it adopts 1.4 MHz receiver bandwidth, a reduced maximum transmission power (20 dBm), and extended discontinuous reception (eDRX) modes. However, the RA procedure and the data transmission are completely inherited from the conventional LTE/LTE-A shown in Fig. 3.3.

As regards the NB-IoT technology, it further pursues the goals of providing a cost-effective solution. In fact, some relevant NB-IoT features are: 200 kHz receiver bandwidth, half-duplex operation, new narrow-

### **Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios**

---

band (NB) physical channels for downlink and uplink, and improved power saving modes. In addition, the 4-step RA procedure has been inherited from the LTE/LTE-A, while the data packet transmission has been deeply modified to support small-sized data in a more efficient way. In particular, a Cellular IoT (CIoT) EPS optimization in the Control Plane (CP) has been designed to allow a piggyback data packet into the NAS service request, which is sent together with the RRC connection setup complete message. Thus, the NB-IoT CP CIoT EPS optimization allows the UE to transmit the data packet immediately after the 4-step RA procedure [74].

We emphasize that both the aforementioned RA procedures and data transmissions are connection-oriented. Therefore, although NB-IoT proposes a re-engineering of the conventional RRC procedures, it suffers from the RRC connection setup overhead.

#### **3.2.3 5G Uplink frame structure**

Current 5G radio interface proposals intend to use, in uplink as well as in downlink, Orthogonal Frequency Division Multiple Access (OFDMA) [4]. Introducing flexible numerologies is a key feature to satisfy the stringent requirements of 5G NR for reliability, latency, and data rate. These numerologies would differ in subcarrier spacing from 15 kHz up to 480 kHz, Cyclic Prefix (CP) length, TTI length, etc. As regards the radio frame composition, the time slot duration ( $T_s$ ) is typically equal to 14 OFDM symbols, each one having a duration dependent on the subcarrier spacing (e.g., for 15 kHz subcarrier spacing, one  $T_s$  is equal to 1ms). One sub-frame duration is fixed to 1 ms and 10 sub-frames make one frame, which is 10 ms long. The proper numerology could be selected on basis of the service type (eMBB, URLLC, and mMTC), and link type (uplink or downlink). In addition, at the aim of achieving better multiplexing performance, 3GPP gives the opportunity to apply NOMA techniques [75], which are based on the OFDMA grid, but use different domain (e.g., power domain or code domain).

### 3.3 Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

---

#### 3.3.1 System Model

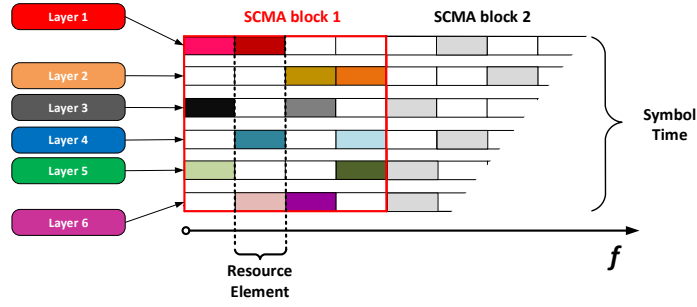
As reported in the 5G implementation guidelines provided by GSMA in [76], 5G networks can be deployed in different deployment options, where StandAlone (SA) options consist of only one generation of radio access technology (i.e., LTE or NR), and Non-StandAlone (NSA) options consist of NR radio cells combined with LTE radio cells using dual connectivity to provide radio access. In this paper, we consider a transmission scenario based on SA Option 2, where the radio access network consists of only one next Generation NodeB (gNB) connected to the 5G Core and  $N_{MTC}$  MTC devices. Each MTC device performs a contention-based access procedure to request resources for data transmission. We do not adopt a NSA option because the dual connectivity technology is not suitable for supporting low-complexity and energy-constrained MTC devices. Moreover, among the SA options, we choose option 2 because it is the only one that adopts the NR radio access, and therefore it is able to fully support the mMTC scenario.

As regards the radio interface, we adopt the smallest 5G NR numerology (i.e., subcarrier spacing of 15 kHz, which corresponds to  $T_s = 1$  ms), because small subcarrier spacing results in longer symbol duration and lower overhead. Therefore, delay-tolerant services, such as MTC services, can benefit from small subcarrier spacing to reduce bandwidth consumption. In addition, we assume that each RA cycle occupies the same amount of uplink resources of LTE, i.e., 72 subcarriers of 15 kHz and a fixed time length,  $T_{ra} = 5$  time slots (see Fig. 3.2). Also, we adopt Preamble Format 0 (i.e., a preamble lasts  $1 T_s$ ). In the frequency domain, PRACH and PUSCH occupy the entire considered bandwidth, while in the time domain PRACH lasts  $T_{pr} \in \{1, 2, \dots, T_{ra} - 1\}$  time-slot units and PUSCH  $T_{pu}$  slots. In summary, for each RA cycle we have:

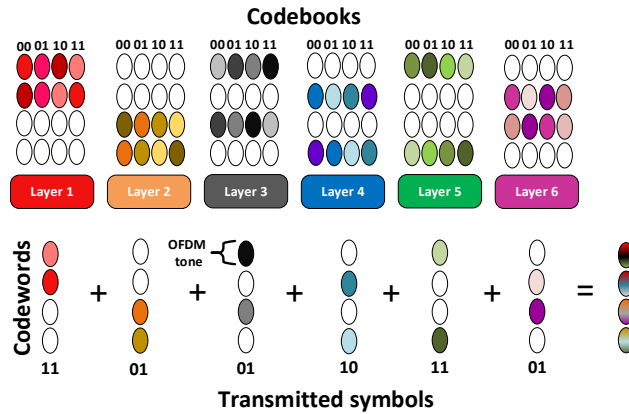
$$T_{ra} = T_{pr} + T_{pu}. \quad (3.1)$$

We consider only  $L_0$  out of 64 preambles available for the RA contention-based procedure in  $1 T_s$ . So, the total number of preamble sequences ( $L$ )

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



(a) Multiple Access in one SCMA block.



(b) Bit-to-codeword mapping.

**Figure 3.4:** Example of multiple access and bit-to-codeword mapping of an SCMA encoder with  $Q = 4$ ,  $S = 2$ ,  $K_{max} = 3$ ,  $L_{SB} = 6$ ,  $I = 4$ .

in a RA cycle is  $L = L_0 T_{pr}$ .

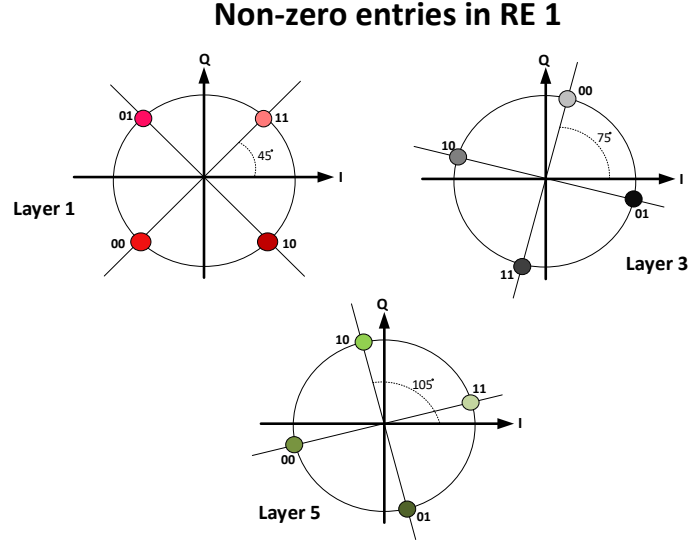
#### SCMA Overview

To further increase the transmission efficiency, we adopt the SCMA technique for PUSCH resources. The SCMA technique, defined in [77], can be regarded as a code division Non-Orthogonal Multiple Access (NOMA) scheme, i.e., it allows multiple transmissions on the same Resource Element (RE), as shown in Fig. 3.4a. Multiple SCMA layers are defined, and one or more layer can be assigned to a device or data stream. In Fig. 3.4a, on the first RE, the first part of each transmitted symbol belonging to Layers 1, 3, and 5 are superposed.

The superposition of different levels in a single RE is allowed by the use of different codebooks, as shown in Fig. 3.4b, which are built on multidimensional constellations, instead of the conventional linear spreading, typical of the traditional Code Division Multiple Access (CDMA).



### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



**Figure 3.5:** SCMA Codebooks. Example of the first dimension of the codebook in RE1.

In addition, SCMA uses sparse spreading to reduce the number of symbol collisions. For instance, in Fig. 3.4a the number of superposed layers in a RE is 3 instead of 6 (traditional case with non-sparse spreading). The superposition pattern per RE is typically statically configured and indicated by the factor graph  $\mathbf{F}$ , which is a binary matrix of size  $Q \times L_{SB}$ , whose element  $f_{i,j} = 1$  if and only if Layer  $j$  transmits inside RE  $i$ . For instance, the factor graph  $\mathbf{F}$  in Fig. 3.4a is:

$$\mathbf{F} = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 & 0 \end{bmatrix}. \quad (3.2)$$

According to SCMA encoder, one  $SCMA_{block}$  is the minimum resource quantity which can be shared by different layers. Inside a single  $SCMA_{block}$ , a layer can be associated to the transmission of one symbol. One  $SCMA_{block}$  occupies  $Q$  subcarriers in one symbol time, so  $Q$  is the number of REs in one  $SCMA_{block}$ . Also, we denote with  $S$  the number of REs which a layer occupies respect to  $Q$ , and  $K_{max}$  the maximum number of overlapped layers with different codebooks in one RE. So, the

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

number of different layers per  $SCMA_{block}$  is:

$$L_{SB} = \binom{Q}{S}, \quad (3.3)$$

and

$$K_{max} = \binom{Q-1}{S-1} = \frac{L_{SB}S}{Q}. \quad (3.4)$$

The gNB receives, over each Resource Element, the sum of the symbols sent by  $K_{max}$  MTC devices, each one by using the codebook assigned to it. To detect the signals transmitted by the MTC devices, the iterative Message Passing Algorithm (MPA) is performed [78]. The complexity order of this method is given by  $O(I^{K_{max}})$  per RE for each iteration, where  $I$  is the number of constellation points [79].

Let us analyze an example of an SCMA encoder in Fig. 3.4a, where  $Q = 4$ ,  $S = 2$ ,  $L_{SB} = 6$ , and  $K_{max} = 3$ . In order for the receiver to can decode the superposed data, in transmission at least  $K_{max}$  different constellations must be adopted. Each constellation should be  $S$ -dimensional and each dimension should be composed of  $I$  points. Several constellations are available in literature [48, 80] which typically are designed on basis on a multi-dimensional mother constellation with a good Euclidean distance profile. The mother constellation is then rotated to achieve a reasonable product distance. A simple example is shown in Fig. 3.5, where we show a constellation adopted for the first non-zero entry (i.e., the first dimension of the constellation) of Layers 1, 3, and 5, which are superposed in the same RE 1. The three constellations are generated from the mother constellation QPSK by applying the phase rotations  $\{0, \frac{\pi}{6}, \frac{\pi}{3}\}$ , as reported in [77]. Then, all symbols transmitted in the same RE do not match each other. Also, we note that, in this example, the same information bits are retransmitted in a second RE ( $S = 2$ ) to reduce the bit error probability at the receiver side.

Since the PUSCH is allocated over 72 subcarriers and each time slot contains 14 OFDM symbols per subcarrier, the number of  $SCMA_{block}$  per time slot is:

$$N_{SB} = \left\lfloor \frac{72}{Q} \cdot 14 \right\rfloor. \quad (3.5)$$

Finally, the number of layers in one RA cycle,  $L_{RAC}$  is:

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

---

$$L_{RAC} = N_{SB}L_{SB}T_{pu}. \quad (3.6)$$

#### The proposed basic two-step RA procedure

In parallel to both the eMTC and NB-IoT technologies analyzed in Section 3.2.2, that adopt a connection-oriented approach, connectionless transmission protocols for machine-type traffic have been proposed in literature at the aim of avoiding the RRC connection setup overhead. These protocols allow the MTC device to transmit small packets without the establishment of radio bearers [81–83]. So, the signaling message exchange between the network and the MTC device that is used to set up the RRC connection and to establish device and security contexts is deleted, implying a device battery power saving.

Among the available proposals, we consider the 2-step connectionless packet transmission procedure described in [82]. By adopting this procedure, the MTC device transmits, immediately after the reception of the RAR message (Step 2), its data packet together with an UL context containing all necessary information related to the device identity, PDN-ID, and security. Thus, once the gNB receives the data packet, it forwards the packet to a connectionless access gateway, which inspects the context header, verifies integrity, performs decryption, and, based on stored state information, forwards the packet to the expected network entity. The main drawback is that the collision detection occurs after the packet has been transmitted, i.e., regardless of whether the access attempt has been successful or not, the device sends the data packet. Taking into account this procedure, we propose a new improved connectionless 2-step RA procedure and data transmission shown in Fig. 3.6. Like in [82], the MTC device transmits, after Step 2, the data packet piggybacked with the UL context. Conversely, at the aim of overcoming the issue of sending the data packet regardless of whether the access attempt has been successful or not, we adopt the early Preamble Collision Detection (e-PACD) technique, proposed in [84], where the gNB can detect in Step 1 whether a preamble has been affected by collision or not. In detail, each device randomly selects one preamble among those available for contention-based procedure and transmits a tagged preamble, consisting of both the selected preamble and a tag sequence. By

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

exploiting the received tagged preambles, the gNB can detect, for each received preamble, whether a collision has been occurred by extrapolating the tags associated to it and verifying if more than one tag has been transmitted.

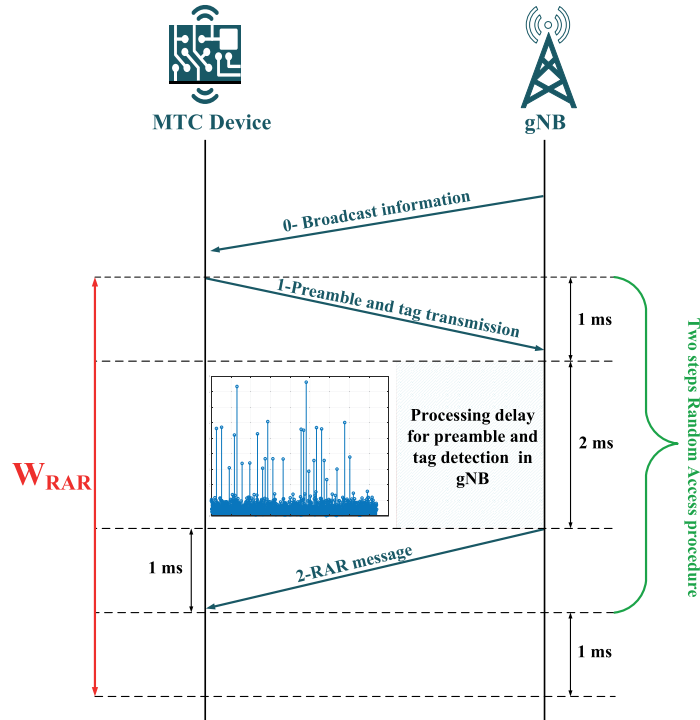
Moreover, because for the machine type communications a very small amount of data is expected [85], we assume for each transmission request the same upper bound value, sent in broadcast, denoted as  $\theta_{max}$  bits. So, in Step 2, the gNB assigns to each successful access attempt, the PUSCH resources needed to satisfy the maximum data transmission ( $\theta_{max}$ ), if available.

Consequently, the MTC device which has received the message response from the gNB, within the  $W_{RAR}$  window, transmits its data in the PUSCH of the next RA cycle. We note that the data can be transmitted in the next RA cycle since we set  $W_{RAR} = 5\text{ms}$ . In fact, thanks to the two-step procedure, this time is sufficient to guarantee the transmission of the tagged preambles (including cyclic prefix and propagation delay), the processing delay at the gNB side, the transmission of the RAR message, and 1ms of margin time. Conversely, if the message has not been received from the gNB, within the  $W_{RAR}$  window, then it will reattempt only if  $N_A \leq M_A$ .

The proposed 2-step RA procedure results not only in very few messages being exchanged but also in a very low signaling overhead, as depicted in Fig. 3.7 and described in the following. In order to reduce the signaling overhead, we assume that SCMA blocks are assigned to MTC device in multiples of SCMA Block Groups (SBGs), i.e., one SBG is the smallest unit of resources that can be allocated to an MTC device. An SBG corresponds to one SCMA block assigned for the time of 1  $TTI$ . Moreover, we assume that Matrix  $\mathbf{F}$ , a set of  $L_{SB}$  codebooks and the mapping between codebook and the assigned layer are static and known by both sides of the communication link (e.g., sent in broadcast in the SIB).

As regards the number of bits associated to the SCMA allocation information for the MTC device, it results that, when  $S = 2$ , the 18 bits reserved for time and frequency allocation in the conventional RAR (as reported in Section 3.2.1) are enough to deliver the information, regardless of the value of  $Q$ . Therefore, no additional bits need to be sent by

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

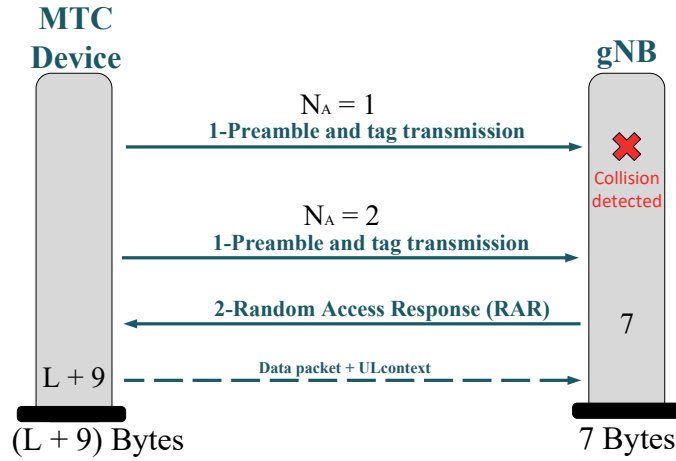


**Figure 3.6:** *The proposed basic two-step RA procedure.*

the gNB in DL. The proof is reported in Appendix 3.3.7. Furthermore, as in [81], we assume that the UL context header sent together with the data packet is of 9 bytes. Globally, the total amount of bytes necessary to transmit  $L$  bytes of data is equal to  $L + 9$  bytes in UL and 7 bytes in DL, regardless the number of access attempts, i.e.,  $N_A$ .

In order to underline the advantages of our two-step RA procedure, we compare the signaling overhead with the traditional four-step RA procedure reported in Section 3.2.1, and the one adopted by NB-IoT. By adopting the LTE procedure depicted in Fig. 3.3, the amount of bytes needed to transmit  $L$  bytes of data is equal to  $(L + N_A \cdot 7 + 49)$  in UL and  $(N_A \cdot 7 + 113)$  in DL. We note that this legacy RA procedure is the same adopted by the enhanced MTC (eMTC) technology. Conversely, the NB-IoT technology adopts a slightly simplified 4-step procedure consisting of  $(N_A \cdot 8 + L)$  bytes in UL and  $(N_A \cdot 7 + 8)$  in DL [81]. In conclusion, our proposal not only shows a lower amount of signaling with respect to both eMTC and NB-IoT (if  $N_A > 1$ ). In addition, by adopting both the eMTC and the NB-IoT RA procedure, the MTC device for each RA attempt needs to perform 2 transmissions, independently whether the procedure has been succeeded or not. Instead, by adopting our proposal,

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.7:** Two-step RA procedure and data packet transmission.

the number of transmission per RA attempt is reduced to 1. So, the connectionless concept helps to reduce energy consumption in the MTC device, mainly because in this case there is a lower number of transmissions performed by the radio module in comparison with the eMTC and NB-IoT [86].

#### Traffic Model

The Beta traffic model is related to a scenario where a large number of MTC devices access the network in a highly synchronized manner [67]. Therefore, we assume that each MTC device generates data at time  $t \in [0, T_{arrival}]$  following the Beta distribution [87] and that each activation time is independent of each other. The Beta distribution is characterized by the following Probability Density Function (PDF):

$$f(t) = \frac{t^{\alpha-1} (T_{arrival} - t)^{\beta-1}}{T_{arrival}^{\alpha+\beta-1} \text{Beta}(\alpha, \beta)}, t \in [0, T_{arrival}], \quad (3.7)$$

where

$$\text{Beta}(\alpha, \beta) = \int_0^1 t^{\alpha-1} (1 - t)^{\beta-1} dt. \quad (3.8)$$

Given  $N_{MTC}$  devices, we derive a discrete random process which counts the number of data arrivals per RA cycle inside the time period  $[0, T_{arrival}]$ . Let  $\{N(\omega, n)\}$  denote this discrete random process, where  $\omega \in \Omega$  is an experiment and  $n = 1, 2, \dots, \lfloor \frac{T_{arrival}}{T_{ra}} \rfloor$  is a point in time

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

representing the  $n$ th RA cycle. Then, the function  $N(\bar{\omega}, n)$  is a realization of the random process related to the outcome  $\bar{\omega}$ . For each  $n$  value,  $N(\bar{\omega}, n)$  is the sum of the arrivals in the interval time  $[(n-1)T_{ra}, nT_{ra}]$ .

We can define the temporal auto-correlation function  $r_{NN}(\bar{\omega}, l)$  at the time-lag  $l$  between two points in time.

$$r_{NN}(\bar{\omega}, l) = \sum_{n=1}^{\lfloor \frac{T_{arrival}}{T_{ra}} \rfloor - l} N(\bar{\omega}, n)N(\bar{\omega}, n+l), \quad (3.9)$$

for  $l = -\lfloor \frac{T_{arrival}}{T_{ra}} \rfloor + 1, \dots, \lfloor \frac{T_{arrival}}{T_{ra}} \rfloor - 1$ . By simulation, we estimated that  $r_{NN}(\bar{\omega}, l) = r_{NN}(l)$ , and the mean value  $\mu_n(\bar{\omega}) = \mu_n$ , for a very large number of experiments, thus we assume that the process is ergodic in the wide sense.

The auto-covariance can be also expressed as a function of the time-lag  $l$  as follows:

$$c_{NN}(l) = r_{NN}(l) - \mu_n^2. \quad (3.10)$$

The Pearson Auto-Correlation Coefficient (PAC) is defined as

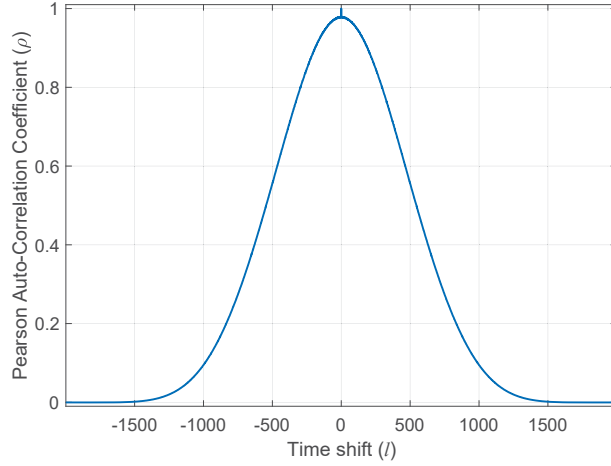
$$\rho(l) = \frac{c_{NN}(l)}{c_{NN}(0)}. \quad (3.11)$$

As it is known, the range of the PAC is  $-1 \leq \rho(l) \leq 1$ . In particular, when  $|\rho(l)| \geq 0.8$  the sequence is highly correlated [88]. Fig. 3.8 shows  $\rho(l)$  when Simulation Parameters in Table 3.6 are adopted. As shown,  $|\rho(l)| \geq 0.8$  when  $|l| \leq 296$ , i.e., the number of arrivals is highly correlated in 296 consecutive RA cycles.

#### 3.3.2 Problem Formulation

In this section, we formulate the problem of the joint control, in a single RA cycle, of the uplink radio resource dimensioning and the random access. In detail, given the number of MTC devices ( $M$ ) which are performing the contention-based RA procedure, we aim to derive the optimal trade-off between the radio resources reserved to the PRACH and to the PUSCH (i.e., the  $T_{pr}^*$  value) jointly with the optimal ACB factor,  $p^*$ . Specifically, in the first subsection we formulate the problem when all MTC devices perform the access attempt without any barring. Then,

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.8:** *Pearson Auto-Correlation Coefficient (PAC) of the random process  $\{N(\omega, n)\}$ .*

in the second subsection, the problem is re-formulated when the conventional ACB scheme is introduced for properly reducing the number of MTC devices attempting the RA procedure.

#### **Uplink resource dimensioning**

For deriving the optimal  $T_{pr}^*$ , we need to estimate the number of successful accesses in PRACH ( $N_{SA}$ ) and the number of available resources in the PUSCH based on SCMA technique ( $N_{ST}$ ) as function on  $T_{pr}$ . As regards  $N_{SA}$ , we note that it is a random variable whose mean value ( $\bar{N}_{SA}$ ) is derived in Appendix 3.3.7, on basis of  $T_{pr}$  and  $M$ , as follows:

$$\bar{N}_{SA} = M \left( 1 - \frac{1}{L_0 T_{pr}} \right)^{M-1}. \quad (3.12)$$

Clearly, given  $M$ , the higher  $T_{pr}$ , the higher  $\bar{N}_{SA}$ .

As regards  $N_{ST}$ , we derive the maximum number of available transmissions in the PUSCH, each one consisting of  $\theta_{max}$  bits, as follows:

$$N_{ST} = \left\lfloor \frac{L_{RAC}}{\left\lceil \frac{\theta_{max}}{r \log_2(I)} \right\rceil} \right\rfloor = \left\lfloor \frac{N_{SB} L_{SB}}{\left\lceil \frac{\theta_{max}}{r \log_2(I)} \right\rceil} \right\rfloor (T_{ra} - T_{pr}), \quad (3.13)$$

where  $\log_2(I)$  is the number of information bits sent for each symbol of the constellation and  $r$  is the code rate. It is evident that, the higher  $T_{pr}$ , the lower  $N_{ST}$ .



### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

The above results confirm the following major problems.

1. The PRACH resources are not sufficient when thousand or more devices simultaneously request to send data. In fact, as reported in (3.58), for very high  $M$  values, the smaller  $T_{pr}$ , the faster  $\bar{N}_{SA}$  approaches zero. This problem could be alleviated by increasing the access resources, i.e.,  $T_{pr}$ .
2. The number of resources available for the PUSCH decreases with increased  $T_{pr}$ . So, also if we have a high value of successful accesses, by adopting a high  $T_{pr}$  value, a part of these ones could not be satisfied due to lack of resources in the PUSCH.

Considering the previous issues there is the need to find a good trade-off between the amount of access resources, contained in  $T_{pr}$ , and the ones used for data transmission, contained in  $T_{pu}$ . For the sake of simplicity, in the problem formulation we assume that the MTC devices, which have successfully passed their access attempt in the PRACH, will transmit their data in the PUSCH of the same RA cycle (if available), instead of the next RA cycle, as reported in Section 3.3.1. Therefore, thanks to this assumption, the amount of successful communications ( $S_{MTC}$ ) inside an RA cycle can be derived as:

$$S_{MTC} = \min(N_{SA}, N_{ST}). \quad (3.14)$$

Since  $N_{SA}$  is a random variable, we approximate (3.14) as:

$$\tilde{S}_{MTC} = \min(\bar{N}_{SA}, N_{ST}). \quad (3.15)$$

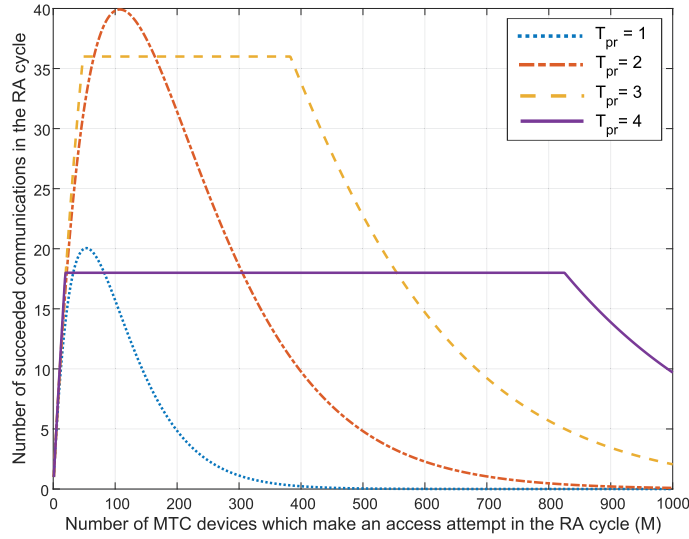
In Fig. 3.9 we plot (3.15), that is, the variation of  $\tilde{S}_{MTC}$  versus  $M$  for any  $T_{pr}$  value. We observe that the straight lines are caused by the limited PUSCH resources. So, the goal can be formulated as follows:

*"Given  $M$  value, find the optimal  $T_{pr}$  value so that*

$$T_{pr}^*(M) = \arg \max_{T_{pr} \in \{1, \dots, T_{ra}-1\}} \{ \min(N_{ST}, \bar{N}_{SA}) \}, \quad (3.16)$$

*where the dependence on  $M$  is contained in  $\bar{N}_{SA}$ ."*

Choosing  $T_{pr} = T_{pr}^*(M)$ , we have the optimal resource dimensioning in the RA cycle and consequently the maximization of  $\tilde{S}_{MTC}$  value, denoted as  $\tilde{S}_{MTC}^*(M)$ .



**Figure 3.9:** Number of succeeded MTC devices transmission vs number of MTC devices which carry out the RA procedure for different  $T_{pr}$  values when  $\theta_{max} = 160$  bits.

This optimal uplink resource dimensioning should be iterate for each RA cycle, on the basis of the traffic load, so that a dynamic load-aware PRACH and PUSCH resource dimensioning is applied. We underline that, although a static solution with respect to a dynamic one apparently does not show significant improvements, if a temporal analysis is carried out considering thousands of RA cycles, then the gain becomes important. In fact, the number of retransmissions may become so large that only the resource dimensioning  $T_{pr} = 4$  can prevent the system from collapsing. On the other hand, a static dimensioning  $T_{pr} = 4$  would be highly inefficient because it should be adopted only in presence of very high traffic load. In this regard, in our conference work [63], we showed the considerable advantage of a Dynamic Uplink Resources Dimensioning (DURD) over any static one, by analyzing the distribution of the new access attempts in the space of 10 s, i.e., 2000 RA cycles, and we faced the challenge of predicting the expected traffic load.

### Joint dimensioning of the uplink resources and ACB scheme

Despite the benefits of the optimal uplink resource dimensioning (3.16), it does not consider the possibility of reducing the preamble collision probability. For this reason, we considered the ACB technique [59] proposed by 3GPP to redistribute the access requests of MTC devices

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

through time. In fact, by means of this technique, the number of access requests per RA cycle is reduced, thus also the number of RA failed attempts is reduced. The goal of an optimal access control based on the conventional ACB scheme is to dynamically derive the optimal ACB factor ( $p^*$ ) which maximizes the amount of successful communications, i.e.,  $S_{MTC}^*(M)$  in our study. Then, this ACB factor value must be continuously provided by the gNB in broadcast. By integrating the conventional ACB scheme in our uplink resource dimensioning, the amount of MTC devices ( $M$ ) performing the RA procedure could be properly reduced.

Given  $\theta_{max}$ , from (3.15) a curve for each  $T_{pr}$  value is obtained, and we define  $T_{pr}^*$  as:

$$T_{pr}^* = \arg \max_{\forall M, T_{pr}} \{ \min(N_{ST}, \bar{N}_{SA}) \}. \quad (3.17)$$

Also, we denote the maximum  $\tilde{S}_{MTC}$  value as  $S_{MAX}$ , that is,

$$S_{MAX} = \max_{\forall M} \left[ \tilde{S}_{MTC}(M, T_{pr}^*) \right]. \quad (3.18)$$

The value of  $M$  that leads to  $S_{MAX}$  is denoted as  $M_{max}$  and is derived in Appendix 3.3.7 for a generic  $\theta_{max}$  value. We underline that, given  $\theta_{max}$ , it follows that  $T_{pr}^*$ ,  $M_{max}$ , and  $S_{MAX}$  are uniquely defined.

Finally, fixed  $\theta_{max}$ , the problem of the optimal uplink resource dimensioning in the presence of an ACB scheme can be re-formulated as follows:

*"Given the  $M$  value*

- *if  $M > M_{max}$ , find the optimal ACB factor  $p^*$  so that  $M$  is reduced to  $M_{max}$ ; the optimal uplink resource dimensioning is  $T_{pr}^*$ ;*
- *if  $M \leq M_{max}$  the ACB scheme is not needed, i.e.,  $p = 1$ , and the optimal uplink resource dimensioning is obtained by using (3.17)."*

For instance, in the case considered in Fig. 3.9 with  $\theta_{max} = 160$  bits, we obtain  $T_{pr}^* = 2$ ,  $S_{MAX} \cong 40$ , and  $M_{max} \cong 108$ .

However, the above approach presents the following significant drawbacks.

1. In order to derive the proper ACB factor  $p$ , there is the need to perfectly know  $M$  for each RA cycle.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

2. Despite this resource dimensioning maximizing the number of succeeded communications in a single RA cycle, the impact that the failed access attempts have on each MTC device is not considered. For instance, when a proper ACB factor  $p$ , so that  $M \cong M_{max} = 108$ , and  $T_{pr} = 2$ , has been communicated by the gNB, only 40 out of 108 communications have succeeded. This means that 68 MTC devices have failed their access attempt, thus their own access attempt counter is increased by 1 unit and due to this iterative process, the counter could reach  $M_A$  quickly.
3. Since ACB factor  $p$  can be changed every RA cycle, the MTC device attempting to send an access request has to remain continuously in the Reception (RX) active state to listen the ACB factor as long as the transmission has been successful or the maximum number of attempts has been reached. It leads to an excessive energy consumption and does not allow the MTC device to apply Discontinuous Transmission and Reception (DTX/DRX) for power saving [65].

To overcome these issues, a new control framework is proposed where the random access procedure and a dynamic uplink resource dimensioning are jointly considered. This control framework is described in the next Section.

#### 3.3.3 Overview of the proposed Control Framework

The objectives of our control framework are:

- to minimize the number of failed attempts, by blocking, for an adaptive time period, the MTC devices which should perform an access re-attempt;
- to maximize the number of succeeded communications over a time period much longer than a single RA cycle, e.g., on the order of a few seconds;
- to minimize the MTC device energy consumption.

The control framework is based on the following features:

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

---

1. a predictive estimation of the total number of MTC device access attempts for each RA cycle;
2. a variable-length superframe, called Cluster (see Fig. 3.11), which consists of a sequence of  $G$  RA cycles having the same uplink resource dimensioning, i.e.,  $T_{pr}$  is constant;  $G$  is the cluster size and represents the superframe length in RA cycle units;
3. an ACB scheme which regulates the number of RA re-attempts based on the Cluster size  $G$ .

We consider that, at the end of each Cluster, the RRM entity computes the suitable size  $G$  for the next Cluster, and the related  $T_{pr}$  value, by means of the algorithm described in the next Section. We assume that the gNB sends the broadcast information inside a SIB with a periodicity of  $5\text{ms}^2$  during the  $T_{ra}$ th time slot of each RA cycle. Particularly, these information are related to the next RA cycle, i.e., the size  $G$  of the Cluster to which the next RA cycle belongs, the  $T_{pr}$  value, and the number  $j$  of the RA cycle within the Cluster.

As regards the MTC device which needs to transmit data, the following phases are carried out.

1. After listening to the broadcast information ( $T_{pr}$ ,  $j$ , and  $G$ ), it starts the two-step RA procedure, by transmitting one of the  $L$  available preambles in the next RA cycle. The counter  $N_A$  is increased by one unit. This event is called "*access attempt*". Let us note that the  $T_{pr}$  is needed to know  $L$ , while size  $G$  and  $j$  are ignored in this phase. If the RA procedure is successful ("*succeeded access attempt*") and there are available resources in the PUSCH, the data will be transmitted in the PUSCH ("*succeeded communication*"). Otherwise, the event, called "*failed transmission*", occurs if either the RA procedure is successful but there are no available resources in the PUSCH ("*PUSCH resource lack*") or the RA procedure is not successful ("*failed access attempt*"), then it follows phase 2.
2. The MTC device, on basis of the previously received information ( $j$  and  $G$ ), waits for the last RA cycle of the current Cluster in order to

---

<sup>2</sup>This choice complies with the options provided by the standard [89]

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

listen the updated broadcast information related to the next Cluster. After that, it draws a uniform random number  $h \in \{1, 2, \dots, G\}$ . Then, it performs its "*random access re-attempt*" only in the  $h$ th RA cycle of the current Cluster.

3. If the RA procedure fails again, the MTC device follows phase 2 iteratively until either the communication will be satisfied ("*successful communication*") or the maximum number of attempts ( $M_A$ ) will be exceeded, i.e.,  $N_A > M_A$  ("*failed communication*").

The complete flow diagram of the MTC device behavior is reported in Fig. 3.10, where the color red represents the 2-step random access procedure, while the color green is the ACB scheme.

We underline that our control framework does not adopt the traditional backoff window, because it is implicit in the proposed ACB procedure. In fact, our ACB scheme allows each MTC device to perform only a single access attempt in the same Cluster, thus the overall amount of RA access attempts are spread in time on basis of a dynamic  $G$  value. It is clear that  $G$  and  $T_{pr}$  need to be properly computed, in order to obtain the best performance.

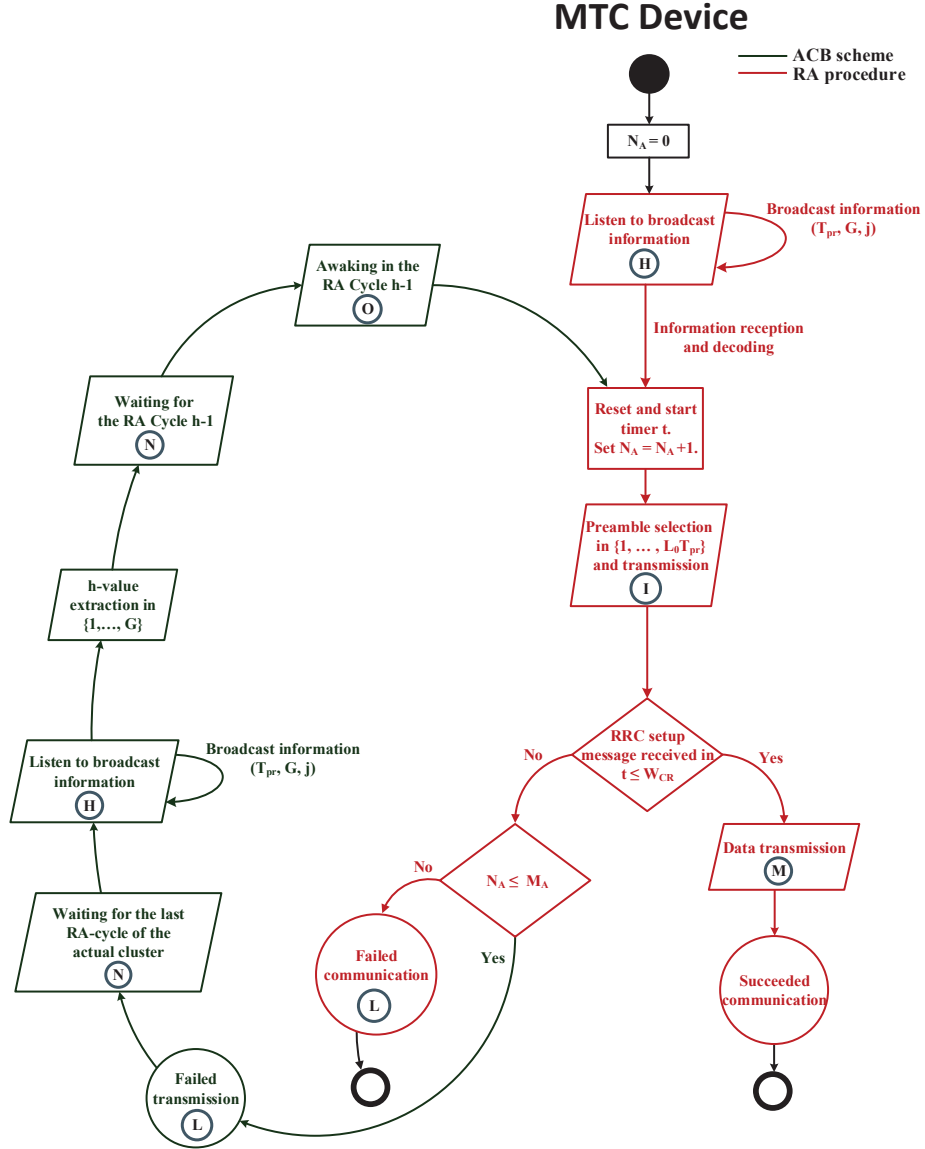
As regards the energy consumption, we observe that, inside a Cluster, each MTC device which needs to re-attempt the access should be active only at the end of the current Cluster (to receive the information related to the next Cluster) and in the selected RA cycle  $h$  of that Cluster. This strategy allows the MTC device to save a considerable amount of energy compared to the conventional ACB scheme, in which the MTC devices need to constantly listen to the broadcast information, since the ACB factor can change every RA cycle.

#### 3.3.4 Enhanced Dynamic Clustering Dimensioning Algorithm

In this section, we present the enhanced Dynamic Cluster Dimensioning Algorithm (eDCDA), that runs in the RRM entity in order to calculate the parameters of the  $i$ th Cluster, i.e.,  $T_{pr_i}$  and  $G_i$ . It is an improvement of the Dynamic Cluster Dimensioning Algorithm (DCDA) presented in [66].

The eDCDA is composed of two distinct phases, as shown in Fig. 3.12.

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



**Figure 3.10:** Flow diagram of the MTC device behavior which adopts the proposed joint control scheme.

#### Phase 1

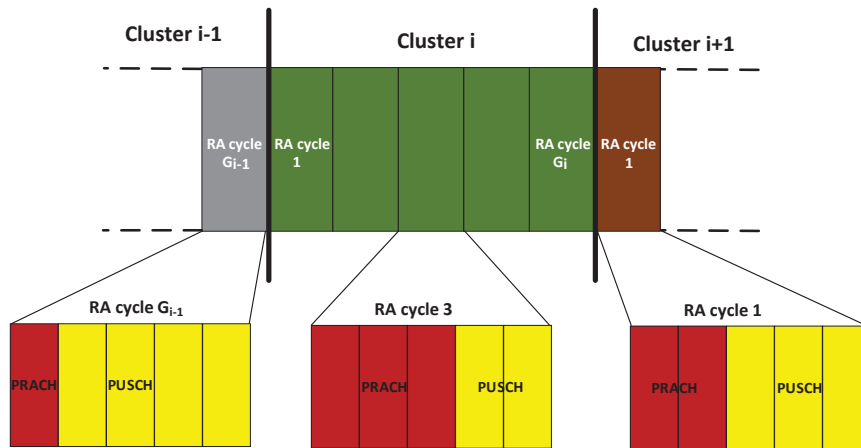
The first phase has the task of estimating the number of expected access attempts in Cluster  $i$ . This phase takes in input the following parameters:

- $G_{i-1}$  = size of the  $(i - 1)$ th Cluster;
- $G_{i-2}$  = size of the  $(i - 2)$ th Cluster;
- $\mathbf{S}_{i-1}$  = vector of size  $G_{i-1}$  containing the values  $s_{i-1}^j$  representing

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

**Table 3.1:** *Simulation Parameters*

Parameter	Symbol	Value
Preambles reserved for contention-based procedure	$L_0$	54
Number of total MTC devices	$N_{MTC}$	5000 : 50000
Number of data transmission requests per MTC device		1
Maximum number of RA procedure attempts	$M_A$	10
Data transmission size	$\theta_{max}$	160 bits
Arrival distribution	Beta	$\alpha = 3, \beta = 4$ $T_{Arrival} = 10s$
Time slot	$T_s$	1 ms
Time Transmission Interval	$TTI$	1 ms
RA cycle duration	$T_{ra}$	5 time slots
Contention resolution window size	$W_{RAR}$	5 ms
Simulation Length	$T_{sim}$	10s
SCMA parameters	$Q$	4
	$K_{max}$	3
	$S$	2
Backoff Window	$B_W$	{20, 40, 80, 160} ms
Constellation Points	$I$	4
Code rate	$r$	1
Multiplicative factor	$\delta$	1.1
Threshold value	$\tilde{\delta}_{rTHR}$	20



**Figure 3.11:** *Examples of Cluster, each one has its own  $T_{pr}$  dimension.*

the number of successful preambles, i.e., the succeeded access attempts, in the RA cycle  $j$  of Cluster  $i - 1$ ;



### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

- $\mathbf{C}_{i-1}$  = vector of size  $G_{i-1}$  containing the values  $c_{i-1}^j$  representing the number of collided preambles, i.e., the failed access attempts, in the RA cycle  $j$  of Cluster  $i - 1$ ;
- $\mathbf{R}_{l_{i-1}}$  = vector of size  $G_{i-1}$  containing the values  $r_{l_{i-1}}^j$  representing the number of failed transmissions due to lack of PUSCH resources in RA cycle  $j$  of Cluster  $i - 1$ ;
- $\tilde{\mathbf{U}}_{i-2}$  = vector of size  $G_{i-2}$  containing the estimated values  $\tilde{u}_{i-2}^j$  representing the sum of the number of failed transmissions due to the lack of PUSCH resources and the failed access attempts in the  $j$ th RA cycle of Cluster  $i - 2$ ;
- $T_{pr_{i-2}}$  and  $T_{pr_{i-1}}$  represent the PRACH dimensioning of Cluster  $i - 2$  and Cluster  $i - 1$ , respectively.

The output parameters of the first phase are:

- $\tilde{\mathbf{U}}_{i-1}$  = vector of size  $G_{i-1}$  containing the estimated value  $\tilde{u}_{i-1}^j$ ;
- $\tilde{n}_{r_i}^1$  = number of estimated new access requests in RA cycle 1 of Cluster  $i$ .

The first output parameter is derived as follows:

$$\tilde{\mathbf{U}}_{i-1} = \tilde{\mathbf{T}}_{r_{i-1}} - (\mathbf{S}_{i-1} - \mathbf{R}_{l_{i-1}}), \quad (3.19)$$

where  $\tilde{\mathbf{T}}_{r_{i-1}}$  is a vector containing the estimated total access requests, including the new access requests and the re-transmitting ones, for each RA cycle of Cluster  $i - 1$ . In order to estimate the elements of this vector, we consider that, with respect to RA cycle  $j$  of Cluster  $i - 1$ , the RRM knows the amount of collided preambles  $c_{i-1}^j$ , but it does not know how many access attempts has been occurred in the collided preambles. For these reasons, to estimate the total access requests, we set:

$$\tilde{t}_{r_{i-1}}^j = s_{i-1}^j + k_{i-1}^j c_{i-1}^j, \quad (3.20)$$

where  $k_{i-1}^j$  is called collision coefficient; obviously,  $k_{i-1}^j \geq 2$  because a preamble collision occurs when at least two MTC devices choose the

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

same preamble in the PRACH. From an empirical study, we derived the relationship between  $k^{j-1}$ ,  $c_{i-1}^j$  and  $T_{pr_{i-1}}$  as follows:

$$k^{j-1} = 2 \cdot 1.97^{\frac{c_{i-1}^j}{L_0 T_{pr_{i-1}}}}. \quad (3.21)$$

Equation (3.21) was found to be sufficiently accurate for our purpose, as proved in [63]. Considering (3.20), in general, vector  $\tilde{\mathbf{T}}_{r_{i-1}}$  can be easily calculated as follows:

$$\tilde{\mathbf{T}}_{r_{i-1}} = \mathbf{S}_{i-1} + \mathbf{k}_{i-1} \circ \mathbf{C}_{i-1}, \quad (3.22)$$

in which  $\circ$  is the Hadamard product operator,  $\mathbf{k}_{i-1}$  is a vector of size  $G_{i-1}$ , whose  $j$ th element is calculated by using (3.21).

The second output parameter,  $\tilde{n}_{r_i}^1$ , is calculated as an approximation of the value  $\tilde{n}_{r_{i-1}}^{G_{i-1}}$ . This choice has been done because the two terms are very highly correlated, in fact they belong to two consecutive RA cycles although they are in two different Clusters. The term  $\tilde{n}_{r_{i-1}}^{G_{i-1}}$  can be derived from the vector  $\tilde{\mathbf{T}}_{r_{i-1}}$ . In fact,  $\tilde{t}_{r_{i-1}}^{G_{i-1}}$  is the total amount of active MTC devices attempting to access to the PRACH resources of the RA cycle  $G_{i-1}$ , i.e., the sum of the re-attempting MTC devices  $\left(\tilde{o}_{r_{i-1}}^{G_{i-1}}\right)$  which failed in the Cluster  $i - 2$  and the new arrived ones  $\left(\tilde{n}_{r_{i-1}}^{G_{i-1}}\right)$ :

$$\tilde{t}_{r_{i-1}}^{G_{i-1}} = \tilde{o}_{r_{i-1}}^{G_{i-1}} + \tilde{n}_{r_{i-1}}^{G_{i-1}}. \quad (3.23)$$

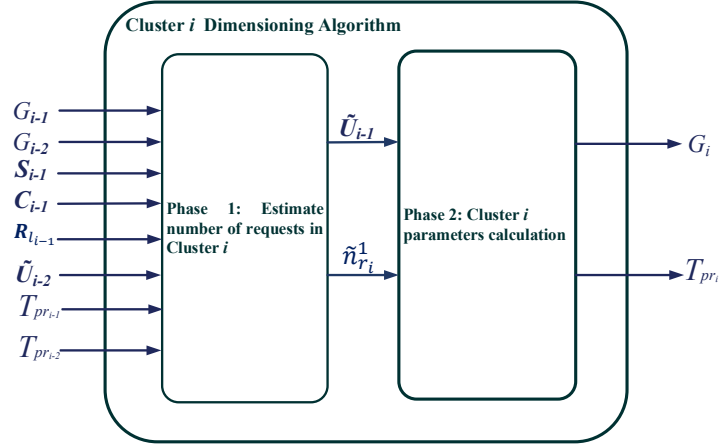
We underline that, by using the proposed ACB procedure, the total failed transmissions in the Cluster  $i - 2$  are distributed, on average, in equally manner among all the RA cycles of Cluster  $i - 1$ . So, the expected value of access re-attempts, in the RA cycle  $G_{i-1}$ , is derived as follows:

$$\tilde{o}_{r_{i-1}}^{G_{i-1}} = \frac{\sum_{k=1}^{G_{i-2}} \tilde{u}_{i-2}^k}{G_{i-1}}. \quad (3.24)$$

Finally, the second output value is:

$$\tilde{n}_{r_i}^1 \simeq \tilde{n}_{r_{i-1}}^{G_{i-1}} = \tilde{t}_{r_{i-1}}^{G_{i-1}} - \tilde{o}_{r_{i-1}}^{G_{i-1}}. \quad (3.25)$$

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



**Figure 3.12:** Enhanced Dynamic Cluster Dimensioning Algorithm (eDCDA) phases.

#### Phase 2

The second phase takes in input the two parameters coming out from the first phase. This phase has the task of calculating the final parameters of the algorithm, i.e.,  $G_i$  and  $T_{pr,i}$ . Compared to the DCDA proposed in [66], in this paper we propose an improved version of the second phase. In [66], on basis of the input values  $\tilde{U}_{i-1}$  and  $\tilde{n}_{r_i}^1$ , we proposed a heuristic criterion which, starting from a Cluster size  $G = 1$ , iteratively increases  $G$  until the following exit condition is satisfied:

$$\tilde{S}_{MTC}^* \geq \tilde{\sigma}_{r_i}^j = \frac{\sum_{j=1}^{G_{i-1}} \tilde{u}_{i-1}^j}{G}, \quad (3.26)$$

i.e., the number of succeeded transmission is at least equal to the amount of access re-attempts belonging to the RA cycles of Cluster  $i$ . The above criterion aims to reduce the probability that the number of reattempts for the single MTC device reaches the  $M_A$  value.

In overall, the improved version of this phase reduces the number of iterations, manages several exceptions that may occur (e.g.,  $\tilde{n}_{r_i}^1 = 0$ ), and improves the Cluster size aiming to satisfy, on average, not only the amount of re-attempting devices but also a proper percentage of new access attempts. The second phase algorithm is summarized in pseudo-code 2.

With the aim of reducing the number of iterations, instead of initializing  $G = 1$ , as in [66], the new algorithm starts by calculating an appro-

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

### Pseudo-code 2 Control Framework

---

#### Definitions:

- $\delta$ : multiplicative factor;
- $G_i$ : chosen Cluster size;
- $G_{MAX}$ : maximum allowed Cluster size value;
- $\tilde{\sigma}_{r_{THR}}$ : threshold value of  $\tilde{\sigma}_{r_i}^j$  when  $\tilde{n}_{r_i}^1 = 0$ ;
- $T_{pr_i}$ : number of PRACH time slots reserved for the chosen Cluster size;
- $T_{pr}$ : number of time slots available per PUSCH;
- $\tilde{\mathbf{U}}_{i-1}$ : vector containing  $\tilde{u}_{i-1}^k$  values which represent the estimated number of failed transmissions.

#### Initialization:

- 1:  $P'_N = 0$
- 2: calculate  $S_{MAX}$  by using (3.18)
- 3: initialize  $G$  by using (3.27)
- 4: **while true do**
- 5: calculate  $\tilde{\sigma}_{r_i}^j = \frac{\sum_{k=1}^{G_i-1} \tilde{u}_{i-1}^k}{G}$
- 6: calculate  $\tilde{t}_{r_i}^j$  by using (3.28)
- 7: set  $M = \tilde{t}_{r_i}^j$ , calculate  $T_{pr}^*$  by using (3.17), and the related  $\tilde{S}_{MTC}^*$  value
- 8: **if**  $G < G_{MAX}$  **then**
- 9:     **if**  $\tilde{n}_{r_i}^1 > 0$  **then**
- 10:         calculate  $P_N$  by using (3.29)
- 11:         **if**  $P_N > 0$  **then**
- 12:             **if**  $P_N > \delta P'_N$  **then**
- 13:                  $P'_N = P_N$  and  $T'_{pr} = T_{pr}^*$
- 14:                  $G = G + 1$
- 15:             **else**
- 16:                 choose  $G_i = G - 1$  and  $T_{pr_i} = T'_{pr}$
- 17:             **break**
- 18:         **end if**
- 19:         **else**
- 20:              $G = G + 1$
- 21:         **end if**
- 22:         **else**
- 23:             **if**  $\tilde{\sigma}_{r_i}^j \geq \tilde{\sigma}_{r_{THR}}$  **then**
- 24:                  $G = G + 1$
- 25:             **else**
- 26:                 choose  $G_i = G$  and  $T_{pr_i} = T_{pr}^*$
- 27:             **break**
- 28:             **end if**
- 29:         **end if**
- 30:         **else**
- 31:             choose  $G_i = G$  and  $T_{pr_i} = T_{pr}^*$
- 32:             **break**
- 33:         **end if**
- 34: **end while**

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

appropriate initial value of  $G$ , as follows.

$$G = \min \left( \left\lceil \frac{\sum_{k=1}^{G_{i-1}} \tilde{u}_{i-1}^k}{S_{MAX}} \right\rceil, G_{MAX} \right). \quad (3.27)$$

The first function of the minimum represents the smallest Cluster size so that the number of estimated access re-attempts per RA cycle ( $\tilde{\sigma}_{r_i}^j$ ) are at most equal to  $S_{MAX}$ . The second term is introduced to prevent the Cluster from reaching a size so large as to lose the high correlation of the new arrivals inside a Cluster. This value, called  $G_{MAX}$  has been calculated in Section 3.3.1 and it is equal to 296.

Starting from this  $G$  value (3.27), for each iteration (i.e., for each increase of  $G$ ) the gNB estimates the future total access requests in the  $G$  RA cycles of Cluster  $i$ , as reported in Lines 5-6. So, for the generic  $j$ th RA cycle ( $j \in \{1, \dots, G\}$ ), this value is equal to:

$$\tilde{t}_{r_i}^j = \tilde{\sigma}_{r_i}^j + \tilde{n}_{r_{i-1}}^1, \quad (3.28)$$

where  $\tilde{\sigma}_{r_i}^j$  is derived by the total estimated failed transmissions of Cluster  $i - 1$ , which will be re-transmitted in Cluster  $i$ , subdivided in equally manner in the  $G$  RA cycles (3.26). As regards the second addend of (3.28), by exploiting the high correlation of the arrival process for  $G \leq G_{MAX}$ , we assume that the amount of new access requests is the same for each  $j$ th RA cycle belonging to Cluster  $i$ . In summary, we assume that  $\tilde{t}_{r_i}^j$  is the same for each RA cycles of the Cluster. Now, by imposing  $M = \tilde{t}_{r_i}^j$  (see Line 7), the gNB calculates  $T_{pr}^*$  by using (3.17), and the related estimated number of succeeded communications,  $\tilde{S}_{MTC}^*$ .

Now, in order to improve the Cluster size, we introduce a new parameter,  $P_N$ , utilized for establishing a percentage of new access attempts to be satisfied. Let us note that  $P_N$  is not a previously set threshold but it will be derived dynamically. The key parameter  $P_N$  is defined as follows:

$$P_N = \frac{\tilde{S}_{MTC}^* - \tilde{\sigma}_{r_i}^j}{\tilde{n}_{r_i}^1}, \quad (3.29)$$

where  $P_N \in ] - \infty, 1]$ , because  $P_N > 1$  leads to the absurd that the number of succeeded transmissions exceeds the total number of access attempts.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

If  $P_N < 0$ , then not even the base exit condition (3.26) is satisfied. So,  $G$  is increased (see Line 24).

If  $P_N = 1$ , the configuration is the best one, because it satisfies the total amount of  $\tilde{tr}_i^j$  and a further increase of  $G$  do not lead to any improvement.

Finally, if  $0 < P_N < 1$ , we verify whether, by increasing the Cluster size, there is a significant improvement in the  $P_N$  value. With the aim of quantifying this improvement, we introduce a  $\delta$  value which represents a multiplicative gain factor ( $\delta > 1$ ). For the sake of simplicity, we denote as  $G' = G + 1$ , the new configuration and as  $P'_N$  the new arrivals' percentage of success obtained by adopting the dimensioning  $G'$ . If  $P'_N > \delta P_N$ , then we consider the obtained improvement as substantial. In this case,  $G$  is iteratively increased until  $P'_N \leq \delta P_N$ . In the latter case, the algorithm ends and the previous configuration is selected as output.

Finally, we focus on a possible case in which  $\tilde{n}_{r_i}^1 = 0$ . This exception, not considered in [66], leads to very inefficient outputs. More specifically, the Cluster size becomes inappropriately large since (3.26) is satisfied only for  $\tilde{\sigma}_{r_i}^j = 1$ , i.e., only one access request per RA cycle is admitted. In the following, we report the analytical proof.

Given  $\tilde{n}_{r_i}^1 = 0$ , therefore  $M = \tilde{\sigma}_{r_i}^j$  in (3.17), condition (3.26) can be written as

$$\min(\bar{N}_{SA}, N_{ST})|_{T_{pr}=T_{pr}^*(\tilde{\sigma}_{r_i}^j)} \geq \tilde{\sigma}_{r_i}^j, \quad (3.30)$$

that means

$$\begin{cases} \bar{N}_{SA}|_{T_{pr}=T_{pr}^*(\tilde{\sigma}_{r_i}^j)} \geq \tilde{\sigma}_{r_i}^j \\ N_{ST}|_{T_{pr}=T_{pr}^*(\tilde{\sigma}_{r_i}^j)} \geq \tilde{\sigma}_{r_i}^j. \end{cases} \quad (3.31)$$

We evaluate when the first condition of system (3.31) is satisfied, i.e., by using (3.77) the values of  $\tilde{\sigma}_{r_i}^j$  that satisfy the following inequality:

$$\tilde{\sigma}_{r_i}^j \left(1 - \frac{1}{L_0 T_{pr}}\right)^{\tilde{\sigma}_{r_i}^j - 1} \geq \tilde{\sigma}_{r_i}^j, \quad (3.32)$$

that simply involves

$$\left(1 - \frac{1}{L_0 T_{pr}}\right)^{\tilde{\sigma}_{r_i}^j - 1} \geq 1. \quad (3.33)$$

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

This condition is satisfied when either  $\tilde{\sigma}_{r_i}^j = 1$  or  $L_0 T_{pr} = \infty$ . Obviously, the only viable solution is  $\tilde{\sigma}_{r_i}^j = 1$ . We evaluate whether the latter solution is also valid for the second condition, therefore, it is a solution for the whole system (3.31). Clearly,  $N_{ST} \geq 1$  is satisfied accordingly because for each RA cycle at least one resource is available in the PUSCH. In summary, the use of exit condition (3.26) when  $\tilde{n}_{r_i}^1 = 0$  would lead to two evident inefficient conditions:

1. a huge waste of the overall resources in each RA cycle of Cluster  $i$  to serve only one MTC device that transmits at most  $\theta_{max}$  bits;
2. the useless creation of large sized Cluster, that delays the successive possible communications.

At the aim of overcoming these issues, we introduce a value  $V$  utilized for determining an exit condition less stringent than (3.26). In particular, instead of satisfying all the access re-attempts, we consider a percentage loss  $V$ . This new exit condition can be expressed as:

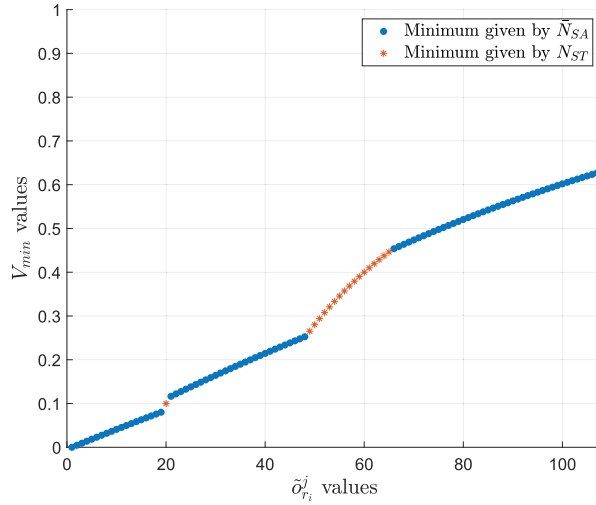
$$\tilde{S}_{MTC}^* \geq (1 - V) \tilde{\sigma}_{r_i}^j. \quad (3.34)$$

So, we need to define the  $V$  value. We consider that, since the ideal value is  $V = 0$ , we should find a proper minimum value of  $V$  ( $V_{min} > 0$ ). At this aim, by substituting (3.15) in (3.34), we obtain:

$$V_{min} = \begin{cases} 1 - \left(1 - \frac{1}{L_0 T_{pr}^*(\tilde{\sigma}_{r_i}^j)}\right)^{\tilde{\sigma}_{r_i}^j - 1} & \text{if } \bar{N}_{SA} > N_{ST} \\ \frac{\tilde{\sigma}_{r_i}^j - N_{ST}}{\tilde{\sigma}_{r_i}^j} & \text{otherwise.} \end{cases} \quad (3.35)$$

Let us note that  $V_{min}$  is only a function of  $\tilde{\sigma}_{r_i}^j$ . In fact,  $\bar{N}_{SA}$ ,  $N_{ST}$ , and  $T_{pr}^*$  are function of  $M$ , and  $M = \tilde{\sigma}_{r_i}^j$ . Equation (3.35) is shown in Fig. 3.13. Fixed a  $V_{min}$  value, i.e., the maximum loss in terms of satisfied access re-attempts, we get the maximum  $\tilde{\sigma}_{r_i}^j$  value which satisfies (3.34). For instance, if  $V_{min} = 0.1$ , we obtain that the related maximum  $\tilde{\sigma}_{r_i}^j$  value is equal to 20 (i.e.,  $G$  is increased until  $\tilde{\sigma}_{r_i}^j \leq 20$ ) and the related  $\tilde{S}_{MTC}^*$  value is at most equal to 18. As shown in Fig. 3.13,  $V_{min}$  is a monotonically increasing function of  $\tilde{\sigma}_{r_i}^j$ , then there is not a clear preferential  $V_{min}$  value.

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.13:** Representation of  $V_{min}$  vs  $\tilde{\sigma}_{r_i}^j$  when  $\theta_{max} = 160$  bits.

At this aim, by analyzing the curve, we note that there are some slope changes, in which the value  $V_{min}$  increases rapidly. In order to quantify these increases, we calculate from (3.35) the following incremental ratio:

$$\frac{\Delta V_{min}}{\Delta \tilde{\sigma}_{r_i}^j} = V_{min}(\tilde{\sigma}_{r_i}^j + 1) - V_{min}(\tilde{\sigma}_{r_i}^j), \quad (3.36)$$

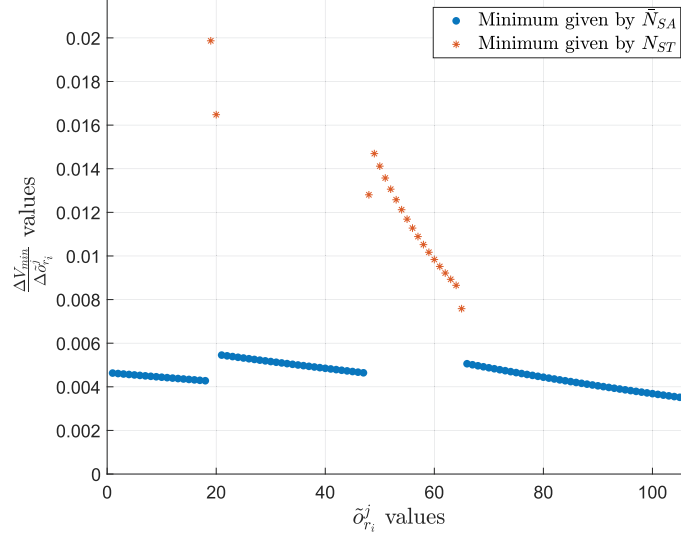
which is depicted in Fig. 3.13. It is clear that the high variations of  $V_{min}$  occur in the transition from the first to the second condition of the system (3.35). We denoted as  $\tilde{\sigma}_{r_{THR}}$  the point corresponding to the maximum variation. Therefore, we set that when  $\tilde{n}_{r_i}^1 = 0$ , the exit condition becomes  $\tilde{\sigma}_{r_i}^j < \tilde{\sigma}_{r_{THR}}$  and  $G$  is no longer incremented. Similar behavior is obtained for different  $\theta_{max}$  values (e.g., 80, 320, and 640 bits) which we do not report here for space reason and the suitable  $\tilde{\sigma}_{r_{THR}}$  values can be derived. In the specific case depicted in Fig. 3.14, i.e.,  $\theta_{max} = 160$  bits, it results  $\tilde{\sigma}_{r_{THR}} = 20$ . Additionally, it is verified that  $T_{pr}^* = 4$  and the related  $S_{MTC}^*$  value is 18, which corresponds to the maximum available transmissions in the PUSCH, i.e., there is no waste of resources in the PUSCH.

### 3.3.5 Performance Evaluation

In this section, we evaluate the performance of our control framework in comparison with DCDA proposed in [66], with the  $ACB_p^*$  scheme [64] and with static uplink resource allocations, under different backoff win-



### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



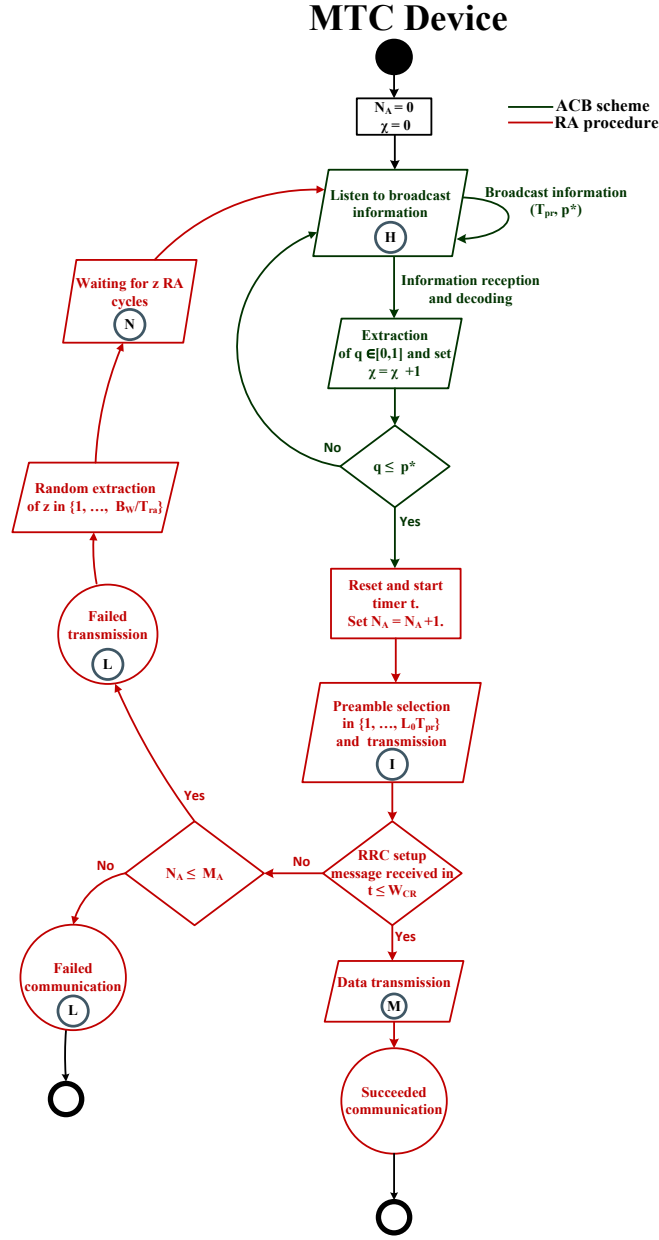
**Figure 3.14:** Representation of the incremental ratio of  $V_{min}$  with respect to  $\tilde{\delta}_{r_i}^j$  vs  $\tilde{\delta}_{r_i}^j$ .

dows and  $T_{pr}$  values. We denote by  $S_i$  the static dimensioning with  $T_{pr} = i$ . The set of parameters used in the simulations is provided in Table 3.6. For a fair comparison, in any control scheme, the PUSCH resources are assigned based on the SCMA technique with the same configuration and the same channel conditions<sup>3</sup>. Specifically, the parameters of the SCMA encoder have been chosen so as to reduce the decoder complexity at the gNB, which is the main bottleneck of the SCMA systems. In fact, the complexity order is  $O(64)$  per RE unit for each iteration. In this case, we obtain 6 transmission layers inside an SCMA block composed of 4 REs, i.e., an overloading factor  $\lambda$  of 150%. On the other hand, adopting a slightly more performing configuration (e.g.,  $Q = 5$ ,  $S = 2$ , that means  $K_{max} = 4$ ,  $L_{SB} = 10$  and  $\lambda = 200\%$ ) involves a much larger computational complexity (i.e.,  $O(256)$  per RE unit for each iteration).

We underline that, by using dimensioning  $S_1$ , the number of succeeded communications values achieved depend exclusively on the very reduced PRACH resources, therefore the multiplexing technique (i.e., traditional SC-FDMA or SCMA) adopted for the PUSCH is irrelevant. Thus, the performance of  $S_1$  is the same to the one of traditional LTE.

<sup>3</sup>In this section, performances are evaluated by considering ideal channel conditions. The considered energy-constrained MTC devices do not follow any channel quality reporting schemes, thus the radio resource allocation performed by these channel-unaware MAC protocols. So, the comparative performance analysis of the different MAC protocols does not depend on the characteristics of the radio channel.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.15:** Flow diagram of the MTC device behavior which adopts the adapted  $ACB_{p^*}$ .

The simulations have been made in MATLAB environment and results are averaged on 50 independent simulations, each one relating to a different realization of the discrete random process of the data arrivals inside the time period  $[0, T_{arrival}]$ .

This section is organized as follows. In Subsection 3.3.5 we describe how the  $ACB_{p^*}$  scheme has been adapted to our system model. Next, in

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

In Subsection 3.3.5 we analyze the Connectivity Capability, that is one of the most challenging target in the mMTC scenario [90], and the Service Delay. Then, in Subsection 3.3.5 we analyze the MTC device's energy consumption, that is another very important parameter for this type of devices [91]. Finally, in Subsection 3.3.5 we quantify the advantages of the eDCDA in comparison with the traditional DCDA when some exceptions occur.

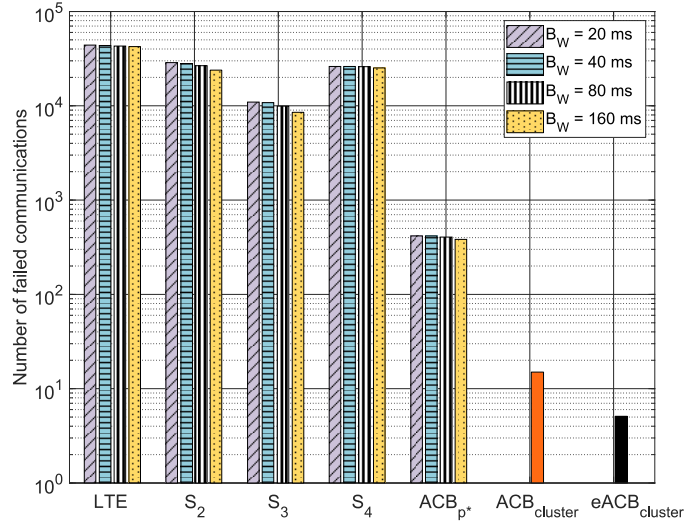
#### Adaptation of the $ACB_{p^*}$ scheme

We compare our solution with the  $ACB_{p^*}$  scheme [64], which is a joint ACB scheme and uplink resource dimensioning. Unlike our control framework, in that work, Xue *et al.* focuses on a single RA cycle in a generic SCMA network. In addition, the dimensioning is based on the amount of MTC devices attempting the contention-based random access  $M$  which is assumed perfectly known. So, in order to fairly assess the two control schemes, we adapt  $ACB_{p^*}$  to our dynamical analysis based on the standard 5G Uplink frame structure as follows:

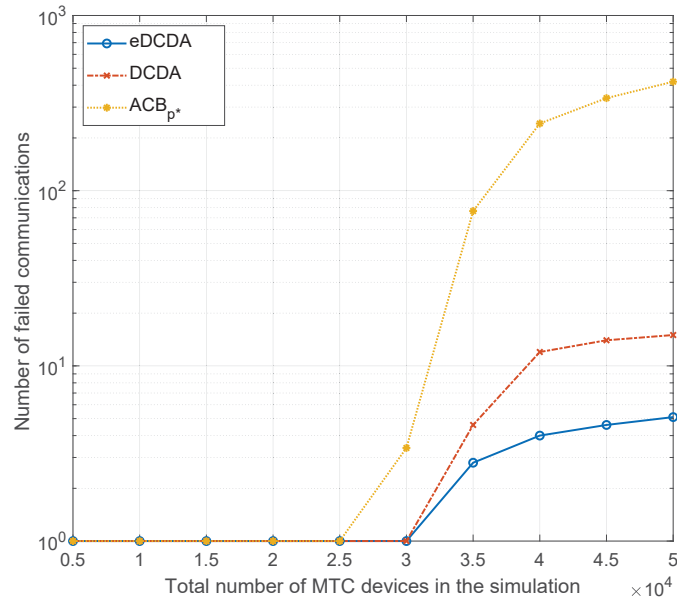
- with regard to the uplink resource dimensioning, the sizing of PRACH and PUSCH takes into account the constraints of the 5G frame structure;
- we inherit their ideal assumption that the  $M$  value is exactly known for the next random access cycle;
- we assume, for conducting a medium-term analysis, that the gNB calculates the ideal value of the ACB factor  $p$  for each RA cycle and transmits it;
- we adopt in  $ACB_{p^*}$  control scheme the proposed 2-step RA procedure and the same  $\theta_{max}$  value.

We note that these choices do not penalize the benchmark scheme but, on the contrary, they allow the control system the best functioning. The MTC device behavior of the adapted  $ACB_{p^*}$  is described in the flow diagram of Fig. 3.15. Compared to our control framework, the MTC device which needs to send its data continuously listens to the broadcast information containing the updated  $p$  value for the next RA cycle until it

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



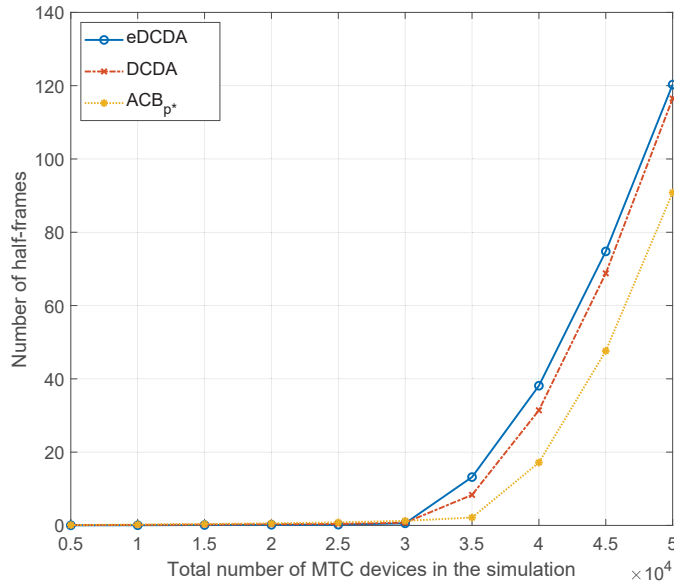
**Figure 3.16:** Connectivity Capability Loss in terms of number of failed communications, when  $N_{MTC} = 50000$ .



**Figure 3.17:** Connectivity Capability Loss in terms of number of failed communications, under different  $N_{MTC}$  values.

extracts  $q \in [0, 1]$  which is less or equal to  $p$ . In the case of failed transmission, the MTC device follows the random backoff, i.e., it waits for a time  $z$  uniformly extracted inside the backoff window before attempting the access procedure again.

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



**Figure 3.18:** Mean value of the service delay for succeeded communication, under different  $N_{MTC}$  values.

#### Connectivity Capability Loss and Delay Analysis

In Fig. 3.16 we show the Connectivity Capability Loss in terms of the average number of MTC devices which have failed their communication in the simulation time. We made a comparison between the static resource allocation schemes (i.e., the traditional *LTE* system,  $S_2$ ,  $S_3$ , and  $S_4$ ), the  $ACB_{p^*}$  scheme, the DCDA framework, and its enhanced version here proposed, called eDCDA, when  $N_{MTC} = 50000$ . It can be seen that the number of failed communications is huge for any static resource allocation schemes. On the other hand,  $ACB_{p^*}$  scheme (for any  $B_W$  values), DCDA and eDCDA, which jointly consider a dynamic uplink resource allocation and an ACB scheme, show a remarkable improvement over any static dimensioning schemes. Consequently, the following comparative evaluation, under different traffic loads (i.e.,  $N_{MTC}$  from 5000 to 50000), considers the latter three control schemes only.

Fig. 3.17 depicts the Connectivity Capability Loss under different values of  $N_{MTC}$ . It can be seen that number of failed communications is non-zero when  $N_{MTC} > 25000$  for the  $ACB_{p^*}$  and  $N_{MTC} > 30000$  for both DCDA and eDCDA.

In addition, the number of failed communications for  $ACB_{p^*}$  is, in the overload region, more than one order of magnitude higher than frame-

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

work and its improved version. Comparing the latter two schemes, the eDCDA framework shows an even more negligible loss. We underline that the figure shows the values averaged over a large number of simulations, thus the benefits of the new proposed algorithm when the exceptional cases occur cannot be so noticeable.

Fig. 3.18 presents the mean service delay for the succeeded communications, which is quantified, for each MTC device, as the amount of time elapsed from the time the data are ready to be transmitted to their actual transmission, expressed in units of RA cycle. The  $ACB_p^*$  scheme shows, in the overload region, slightly better performance in comparison with both the DCDA and eDCDA control frameworks. However, we note that our control framework achieves a considerable advantage in terms of number of served MTC devices (see Fig. 3.17), at the expense of a slight increase in delays (at most, a delay less than 150 ms), which are not significant for the considered delay-tolerant services. Comparing DCDA and its enhanced version, the new scheme presents a very slight increase in terms of delay (i.e., at most 10 ms). On the other hand, this insignificant additional delay makes possible to serve almost the whole amount of MTC devices, even if  $N_{MTC} = 50000$ .

#### Energy Consumption Model and Analysis

It is well known that the energy consumption related to the data transmission, i.e., the energy consumed for performing the RA procedure together with the ACB scheme, is the prevailing energy contribution compared to the energy consumption related to the collection and preparation of the data to be sent.

At this scope, we adopt the energy model depicted in Fig. 3.19, proposed in [92] and already used in [93, 94]. It defines several MTC states, i.e., Deep Sleep (DS), Light Sleep (LS) and Active (A), each one leading to different power consumption. The MTC device can stay only in one of these states at the same time and for changing state it consumes power and takes a certain time period. For instance, to change from the DS state to the LS state it takes 1 TTI and consumes 22 mW/TTI. In order to make more clear the following description, the states and the related changing operations are enumerated from 1 to 10 and the most significant are represented as energy blocks in Fig. 3.20a. The height of

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

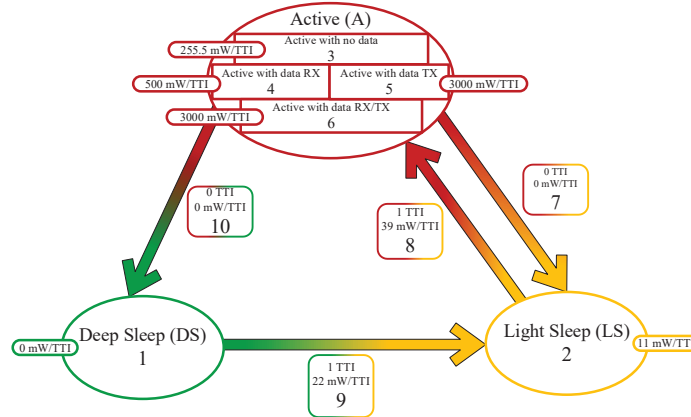
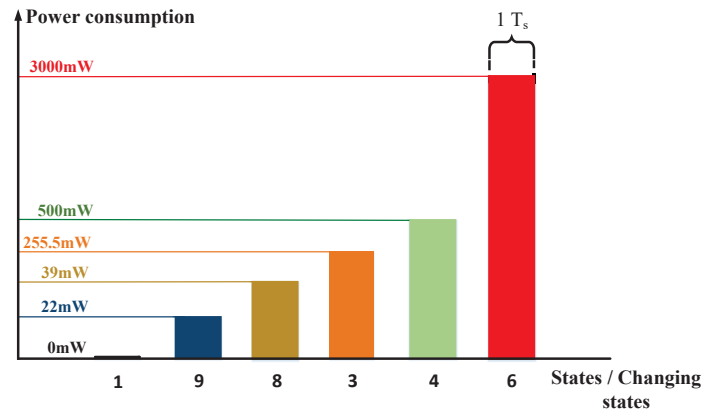
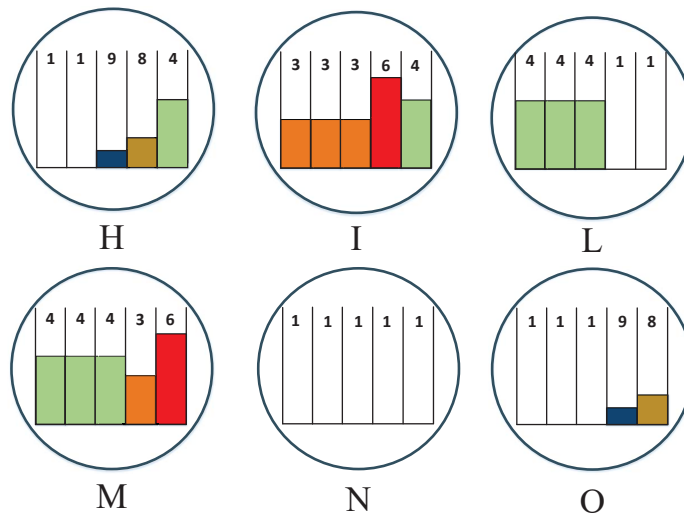


Figure 3.19: Adopted energy model.



(a) Energy blocks.



(b) Energy consumption events.

Figure 3.20: Energy blocks and energy consumption events.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

each energy block represents the power consumption per TTI while the width constitutes the block time duration. For instance, energy block 6 represents the Active state with data TX/RX (see Fig. 3.19), characterized by a power consumption of 3000mW and a time duration of 1 TTI. Moreover, analyzing the behavior of a MTC device in one RA cycle, we can derive different energy events (see Fig. 3.20b), each one represented by a letter from the Latin alphabet from H to O. For instance, energy consumption event H represents the changing state from DS to A with data RX. More specifically, during the first 2 TTIs the MTC device is in DS state and consumes 0 mW (energy blocks 1). Next, it consumes 11 mW for the time of 1 TTI to transit from the DS to the LS state (energy block 9). After that, it consumes 39mW for the time of 1 TTI to change from the LS state to the A with data RX state. Finally, it remains in the A with data RX state for 1 TTI with a consumption of 500 mW. So, the total energy consumption related to block H, indicated with  $\mathcal{H}$ , is 0.561 J.

In the following, we analytically prove that energy consumed by an MTC device which adopts the benchmark control scheme is at least equal to the one consumed by adopting our control framework.

*Proof.* Let us consider the flow diagram of the proposed joint control scheme in Fig. 3.10 and the one of  $ACB_{p^*}$  in Fig. 3.15. In particular, we focus on the blue circles inside the blocks which contain the letters related to the energy consumption events depicted in Fig. 3.20b. We analyze the worst case, i.e.,  $T_{pr} = 4$ , the preamble is transmitted in the last  $T_s$  of the PRACH, and the data transmission occurs in the last time slot of the PUSCH.

In the case where the communication has been succeeded, the energy consumption of our control framework depends almost exclusively on the number of access attempts made by the MTC device, i.e.,  $N_A$ . In fact, in our control scheme (see Fig. 3.10) the MTC device needing to send data reaches the Active state with data RX (Event H) for listening to the broadcast information. Next, the preamble is transmitted and the energy consumption follows the Event I. After that, the MTC device remains in the Active with data RX state and if it receives the RRC setup message within  $W_{RAR}$ , then it transmits its data in the scheduled PUSCH resource (Event M). In the generic case, the total energy consumption



### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

for a succeeded communication in the  $N_A$ th attempt,  $E_{eDCDA}$ , is derived considering the total sequence of events, as follows:

$$E_{eDCDA} = N_A(\mathcal{H} + \mathcal{I}) + (N_A - 1)(2\mathcal{N} + \mathcal{L} + \mathcal{O}) + \mathcal{M}. \quad (3.37)$$

As regards the benchmark scheme, to quantify the energy consumption, we introduce a parameter  $\chi$  which represents the total number of  $q$  extractions (see Fig. 3.15) for one succeeded communication, thus the total number of times it needs to listen to the updated broadcast information. The total energy,  $E_{ACB_{p^*}}$ , is:

$$E_{ACB_{p^*}} = \chi\mathcal{H} + N_A\mathcal{I} + (N_A - 1)(\mathcal{L} + \mathcal{N}) + \mathcal{M}. \quad (3.38)$$

So, for evaluating the energy consumption difference, we derive in which cases the following inequality is satisfied

$$E_{ACB_{p^*}} \geq E_{eDCDA}. \quad (3.39)$$

It follows:

$$\begin{aligned} \chi\mathcal{H} + N_A(\mathcal{I} + \mathcal{L}) - \mathcal{L} + \mathcal{M} &\geq \\ N_A(\mathcal{H} + \mathcal{I} + \mathcal{L}) + (N_A - 1)\mathcal{O} - \mathcal{L} + \mathcal{M}. & \end{aligned} \quad (3.40)$$

Then:

$$\chi\mathcal{H} \geq N_A\mathcal{H} + (N_A - 1)\mathcal{O}. \quad (3.41)$$

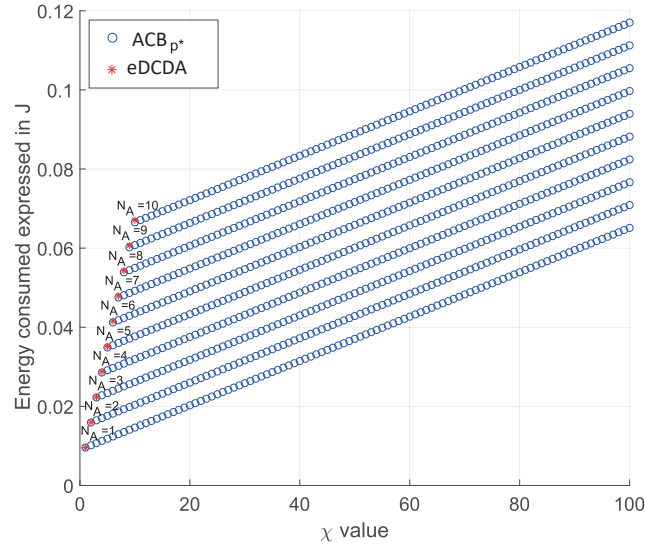
Finally, the last equation might be simplified by ignoring the  $\mathcal{O}$  event, whose energy consumption is negligible being  $\mathcal{H} \gg \mathcal{O}$ . So, it turns into

$$\chi \geq N_A. \quad (3.42)$$

Due to the fact that for each transmission in  $ACB_{p^*}$  there is at least one  $\chi$  extraction for each  $N_A$ , it follows that (3.42) is verified for each  $N_A$  value. Thus, the energy consumption by adopting the  $ACB_{p^*}$  scheme is always greater or equal to that of the proposed control framework.  $\square$

Eqs. (3.37) and (3.38) for different  $\chi$  and  $N_A$  values are depicted in Fig. 3.21. As shown, the energy consumption of  $ACB_{p^*}$  depends strongly on  $\chi$ . For instance, if  $N_A = 5$  and  $\chi = 70$ , i.e., on average 14 extractions per access attempt, the MTC device which adopts the

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.21:** Analytical representation of the energy consumed by adopting the proposed control framework and the  $ACB_{p^*}$  scheme.

$ACB_{p^*}$  scheme consumes twice as much energy than an MTC device which adopts our control framework.

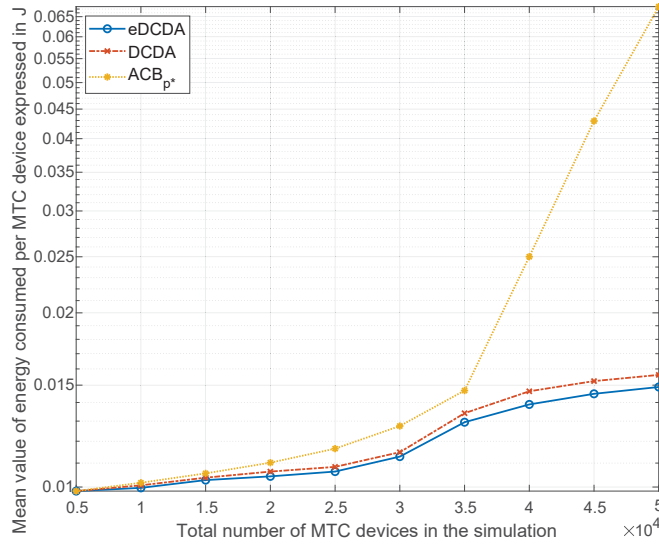
Now, we want to quantify the value of energy savings achieved. To this end, by simulations, we obtain the mean values of Energy Consumption per MTC device, expressed in J, under different  $N_{MTC}$  values, as shown in Fig. 3.22.

We note that, for any  $N_{MTC}$  value, the proposed control framework achieves better performance compared to those of  $ACB_{p^*}$  and DCDA schemes, thus MTC device battery life is increased.

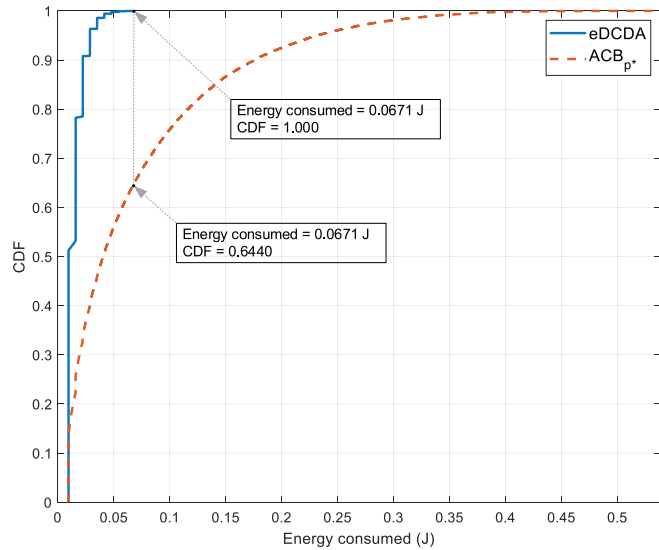
The performance difference between eDCDA and  $ACB_{p^*}$  becomes very significant in a massive MTC scenario, e.g.,  $N_{MTC} \geq 40,000$ .

Finally, due to the fact that the MTC device energy consumption is one of the most important Key Performance Indicators (KPIs), we focus on analyzing also the Cumulative Distribution Functions (CDFs) of the energy consumption. Fig. 3.23 depicts these CDFs calculated by adopting the  $ACB_{p^*}$  and the eDCDA approaches. As regards the eDCDA control framework, we note that, since the energy consumption is only a function of  $N_A$ , as reported in Section 3.3.5, the CDF assumes discrete values. On the other hand, the CDF of  $ACB_{p^*}$  assumes continuous values, since the energy consumption depends on both  $N_A$  and  $\chi$ . Comparing the two CDFs, we note that, when eDCDA is adopted, the

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA



**Figure 3.22:** Mean value of the energy consumed per MTC device for a single succeeded communication.



**Figure 3.23:** CDF of the MTC device energy consumption for one succeeded communication.

total amount of MTC devices consume at most 0.0671 J, and this value corresponds to the maximum value of energy consumption experienced by the 64.4% of devices when  $ACB_p^*$  is adopted. In addition, the maximum energy consumption of  $ACB_p^*$  is 0.535 J, which is almost one order of magnitude larger than the proposed joint control framework.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

Table 3.2: Simulation Test

System	$I_R$	$I_T$	$\sum_{j=1}^{G_{i-1}} \tilde{u}_{i-1}^j$	$\tilde{n}_{r_i}^1$	$G_i$	$\bar{t}_{r_i}^j$	$\bar{s}_i^j$	$\sum_{j=1}^{G_i} s_i^j$
DCDA	1	1693	1693	0	1693	1.61	1.60	627
DCDA $_{G_{max}}$	1	296	1693	0	296	6.46	6.30	1865
	2	49	49	0	49	1.00	1.00	49
eDCDA	1	43	1693	0	85	22.29	17.01	1446
	2	14	435	1	24	19.25	15.37	369
	3	3	93	0	5	19.8	16.4	82
	4	1	16	0	1	16	12	12
	5	1	4	4	0	1	4	4

#### eDCDA vs DCDA

In this subsection, we want to better capture the advantages of the eDCDA in comparison with the traditional DCDA. At this aim, as an example, we consider a single simulation test, and in Table 3.2 we report some results starting from a simulation snapshot at the 1608th RA cycle, when the estimated number of new access requests (column 5) is zero and the estimated number of re-attempting MTC devices (column 4) is equal to 1693. We underline that, given that  $T_{sim}$  is equal to 2000 RA cycles, the amount of remaining RA cycles is 392 and the table shows only the results that occurred within the simulation time. We denote  $I_R$  (column 2) as the index related to the number of times the algorithm runs, and  $I_T$  (column 3) as the number of iterations performed in each run.

As regards the DCDA, since it does not handle the exception  $\tilde{n}_{r_i}^1 = 0$ , the cluster size in output (column 6) is calculated on the basis of estimated values, as  $G_i = \sum_{j=1}^{G_{i-1}} \tilde{u}_{r_{i-1}}^j = 1693$ . It results that the average number of total requests per RA cycle (column 7) is 1.61, the average number of succeeded communications (column 8) is 1.60, and the total number of succeeded communications (column 9) is 627. So, this dimensioning  $G_i$  is erroneous because, although the ratio between the succeeded communications and the access attempts is almost 1, there are still 1287 pending MTC devices at the end of the simulation.

We also report the results obtained by considering a slightly modified version of the DCDA, denoted as DCDA $_{G_{max}}$  which only adds the maximum limit of the cluster size introduced in Section 3.3.1, i.e.,  $G_{max} = 296$ . In this case, the table shows two rows, that is, the algorithm runs

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

twice within the remaining simulation time. In particular, the first run creates a cluster with size  $G_{max}$ , with an average number of total access requests equal to 6.46. Not all MTC devices have been served, so the algorithm is re-run. However, in this case, the exception  $\tilde{n}_{r_i}^1 = 0$  occurs again, and cluster is oversized.

Finally, we report the results obtained by means of the eDCDA. As shown in the table, the algorithm runs five times and the amount of 5 cluster sizes is equal to 116 RA cycles. Within the time of 116 RA cycles all the MTC devices are served by adopting the eDCDA, while in the same time interval the DCDA $_{G_{max}}$  serves only 731 MTC devices and the DCDA only 186. This shows that this cluster sizing is more efficient than both the previous algorithms. In addition, the number of iterations per single run is further reduced thanks to the introduction of an appropriate initial value of  $G$ .

#### 3.3.6 Appendix

#### 3.3.7 Proof of Equation (3.77)

We define the random variables  $X_i$ , with  $i \in \{1, \dots, L\}$ , as:

$$X_i = \begin{cases} 1 & \text{if the } i\text{th preamble has been selected only by} \\ & \text{one MTC device out of the } M \\ 0 & \text{otherwise.} \end{cases} \quad (3.43)$$

Therefore:

$$N_{SA} = \sum_{i=1}^L X_i. \quad (3.44)$$

The mean value of  $N_{SA}$ ,  $\bar{N}_{SA}$ , is calculated as follows:

$$\bar{N}_{SA} = E\{N_{SA}\} = E\left\{\sum_{i=1}^L X_i\right\} = \sum_{i=1}^L E\{X_i\}. \quad (3.45)$$

Since  $X_i$  are  $L$  random variables identically distributed, the mean value  $E\{X_i\}$ , for each  $i$ , is equal to:

$$E\{X_i\} = \Pr(X_i = 1) = M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1}. \quad (3.46)$$

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

The last step has been calculated by applying the total probability theorem for  $M$  independent extractions. So:

$$\bar{N}_{SA} = M \left(1 - \frac{1}{L}\right)^{M-1} = M \left(1 - \frac{1}{L_0 T_{pr}^*}\right)^{M-1}. \quad (3.47)$$

#### Computation of $M_{max}$ value

Given the  $T_{pr}^*$  value calculated in (3.17), if curves  $\bar{N}_{SA}(T_{pr}^*, M)$  and  $N_{ST}(T_{pr}^*)$  have no intersection points, then  $M_{max}$  can be calculated by deriving the  $M$  value which maximizes  $\bar{N}_{SA}$ . At this aim, we impose the derivative of  $\bar{N}_{SA}$  equal to zero:

$$\begin{aligned} 1 + M \log \left(1 - \frac{1}{L_0 T_{pr}^*}\right) &= 0 \implies \\ M &= - \left[ \log \left(1 - \frac{1}{L_0 T_{pr}^*}\right) \right]^{-1}. \end{aligned} \quad (3.48)$$

So, we obtain the following single value:

$$M_{max} = \left[ \log \left( \frac{L_0 T_{pr}^*}{L_0 T_{pr}^* - 1} \right) \right]^{-1}. \quad (3.49)$$

Otherwise,  $M_{max}$  results as a range of values and (3.49) is a generic point belonging to that range. The minimum value of this range can be calculated as:

$$M_{max} = \frac{W \left[ N_{ST}(T_{pr}^*) \left(1 - \frac{1}{L_0 T_{pr}^*}\right) \log \left(1 - \frac{1}{L_0 T_{pr}^*}\right) \right]}{\log \left(1 - \frac{1}{L_0 T_{pr}^*}\right)}, \quad (3.50)$$

where  $W[x]$  is the Lambert W-Function.

#### Signaling Overhead

In this section, we derive the number of bits needed by the gNB for delivering to the MTC device the SCMA allocation information. We remind that in an SCMA system, the physical channel corresponds to

### 3.3. Joint Control of Random Access and Dynamic Uplink Resource Dimensioning for massive MTC in 5G NR based on SCMA

one or more transmission layers to be used inside the allocated SCMA blocks. In this paper, we assume that the radio resources are allocated in multiples of SBGs, each one consisting of  $Q$  subcarriers for the time of 1  $TTI$ . So, considering that the amount of subcarriers allocated to the PUSCH is 72 and that the maximum  $T_{pu}$  value is 4 time slot units, the maximum number of available SBGs in one RA cycle is  $\frac{72}{Q} \cdot \frac{4T_s}{TTI}$ . We also assume that, for each MTC device, the gNB can assign a range of consecutive SBGs and one layer to be used for each SCMA blocks allocated to it. Under this assumption, the total amount of bits required for delivering the SCMA allocation information is equal to:

$$bit_{req} = 2 \cdot \left\lceil \log_2 \left( \frac{72}{Q} \cdot \frac{4T_s}{TTI} \right) \right\rceil + \lceil \log_2 L_{SB} \rceil. \quad (3.51)$$

By considering (3.3) and that  $T_s = TTI = 1$  ms, it follows:

$$bit_{req} = 2 \cdot \left\lceil \log_2 \left( \frac{72}{Q} \cdot 4 \right) \right\rceil + \left\lceil \log_2 \binom{Q}{S} \right\rceil. \quad (3.52)$$

Since  $2 \cdot \lceil x \rceil \leq \lceil 2 \cdot x \rceil + 1$ , it follows:

$$bit_{req} \leq \left\lceil 2 \cdot \log_2 \left( \frac{72}{Q} \cdot 4 \right) \right\rceil + \left\lceil \log_2 \binom{Q}{S} \right\rceil + 1. \quad (3.53)$$

By applying the property  $\lceil x \rceil + \lceil y \rceil \leq \lceil x + y \rceil + 1$ , it results:

$$bit_{req} \leq \left\lceil 2 \cdot \log_2 \left( \frac{72}{Q} \cdot 4 \right) + \log_2 \binom{Q}{S} \right\rceil + 2. \quad (3.54)$$

After some algebraic manipulations, (3.54) can be written as:

$$\begin{aligned} bit_{req} &= \left\lceil \log_2 \left( \frac{72}{Q} \cdot 4 \right)^2 + \log_2 \left( \frac{Q!}{S! (Q-S)!} \right) \right\rceil + 2 \\ &= \left\lceil \log_2 \left( \frac{72}{Q} \cdot 4 \right)^2 + \log_2 \left( \frac{\prod_{i=0}^{S-1} (Q-i)}{S!} \right) \right\rceil + 2 \\ &= \left\lceil \log_2 \left( \frac{(72 \cdot 4)^2}{S!} \cdot \frac{\prod_{i=1}^{S-1} (Q-i)}{Q} \right) \right\rceil + 2. \end{aligned} \quad (3.55)$$

Finally, by applying the property  $\lceil x + y \rceil \leq \lceil x \rceil + \lceil y \rceil$ , the number of required bits as function of  $S$  and  $Q$  is equal to:

$$bit_{req} \leq \left\lceil \log_2 \left( \frac{(72 \cdot 4)^2}{S!} \right) \right\rceil + \left\lceil \log_2 \left( \frac{\prod_{i=1}^{S-1} (Q - i)}{Q} \right) \right\rceil + 2. \quad (3.56)$$

When  $S = 2$ , the number of required bits is given by:

$$bit_{req} \leq \left\lceil \log_2 \left( 1 - \frac{1}{Q} \right) \right\rceil + 18. \quad (3.57)$$

It is obvious that, for each value of  $Q > 2$ ,  $bit_{req} \leq 18$ . This means that the 18 bits reserved for time and frequency allocation in the conventional RAR, are enough to deliver the SCMA allocation information, regardless of the value of  $Q$ .

### **3.4 Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH**

---

#### **3.4.1 System Model and Background**

We consider a 5G network scenario with one Next Generation Node B (gNB) that supports  $N_{MTC}$  MTC devices using a licensed band of  $B = 1.08$  MHz exclusively dedicated to the mMTC scenario. Each device generates a single data at time  $t \in [0, T_{arrival}]$  according to the Beta distribution, as recommended by the 3GPP [67]. As regards the uplink radio interface, we adopt the smallest 5G NR numerology, i.e., subcarrier spacing of 15 kHz, time slot duration ( $T_s$ ) of 1 ms, containing 14 OFDM symbols. This choice has been made because delay tolerant services, such as MTC services, can benefit from small subcarrier spacing to reduce bandwidth consumption [95]. A typical RA cycle lasts  $T_{ra} = 5$  time slots and the uplink radio resources are divided into a PRACH subset and a PUSCH subset. For the preamble sequences in the PRACH we adopt the NR Preamble Format 0, i.e., the PRACH access resources consist of 64 orthogonal preamble sequences that are mapped to 839 subcarriers of 1.25 kHz, and each preamble lasts 1 ms [35]. This choice has been made because the 1.25 kHz numerology is the only option available to support the bandwidth considered, and this Preamble Format maximizes the number of orthogonal preambles available per time



### 3.4. Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH

slot. The total 64 preambles are divided into two groups: the first, consisting of  $L_0$  preambles, is dedicated for the contention based procedure, while the second one is reserved for the contention free procedure. In the time domain PRACH lasts  $T_{pr} \in \{1, 2, \dots, T_{ra} - 1\}$  time slots and PUSCH the remaining  $T_{pu}$  slots. In summary, for each RA cycle we have  $T_{ra} = T_{pr} + T_{pu}$ . Since each preamble lasts  $1 T_s$ , the total number of preamble sequences available for the contention based procedure in each RA cycle is  $L = L_0 T_{pr}$ .

In order to reduce the MTC device energy consumption, we consider the 2-step connectionless RA procedure proposed in [16]. In summary, during Step 1, each MTC device randomly selects one preamble out of the  $L$  available and transmits a tagged preamble sequence on the PRACH. Obviously, there is a non-zero collision probability, since the same preamble can be selected by more than one device. This tagged preamble sequence allows the gNB to detect whether a collision has been occurred immediately after the preamble reception. However, we underline that the gNB does not know how many MTC devices have selected each collided preamble. During Step 2, if the preamble is successfully transmitted, the gNB will send to the MTC a Random Access Response (RAR) message. More specifically, because for the MTC scenario a very small amount of data is expected, we set for each transmission request an upper bound value,  $\theta_{max}$  bits, and we assume that the gNB assigns to each successful access attempt, the  $R_\theta$  PUSCH resources, if any, enough to transmit  $\theta_{max}$  bits. Consequently, the MTC device which has received the RAR message from the gNB, transmits its data packet in the PUSCH of the next RA cycle together with an UL context containing all necessary information related to the device identity and PDN-ID. Conversely, the device reattempts the RA procedure inside the back-off window ( $B_W$ ). In addition, to further increase the transmission efficiency, we adopt the Sparse Code Multiple Access (SCMA) technique for PUSCH resources, that is a promising Non-Orthogonal Multiple Access (NOMA) technique to support massive MTC connectivity requirements with small-size data. By using SCMA scheme, multiple MTC devices can transmit on the same Resource Element (RE) with different sparse codebooks [96].

As regards the PRACH, we denote  $M$  as the number of MTC de-

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

vices which are performing the contention-based RA procedure in the considered RA cycle, and  $P_S$  as the number of succeeded preambles in the PRACH, i.e., the number of successful access attempts. In particular, given  $M$  and  $L$ ,  $P_S$  is a random variable whose expected value is

$$\bar{P}_S = M (1 - 1/L)^{M-1} = M (1 - 1/L_0 T_{pr})^{M-1}. \quad (3.58)$$

For the same  $M$  value, the higher  $T_{pr}$ , the higher  $\bar{P}_S$ .

As regards the PUSCH resources, according to the SCMA encoder, one  $SCMA_{block}$  is the minimum resource quantity which can be shared by different transmissions, called layers. The total number of data transmissions ( $DT_{max}$ ) which can be satisfied in the PUSCH is:

$$DT_{max} = \left\lfloor \frac{\lfloor 72/Q \cdot 14 \rfloor (QK_{max})/S}{\lceil \theta_{max}/\log_2(I) \rceil} \right\rfloor (T_{ra} - T_{pr}), \quad (3.59)$$

where  $Q$  is the the number of REs in one  $SCMA_{block}$ ,  $S$  is the number of REs which a layer occupies respect to  $Q$ ,  $K_{max}$  is the maximum number of overlapped layers in one RE, and  $\log_2(I)$  is the number of bits per symbol [68]. It is evident that, the higher  $T_{pr}$ , the lower  $DT_{max}$ . Finally, the average number of MTC devices which successfully transmit inside a RA cycle, i.e., the number of succeeded communications, is

$$\bar{C}_S = \min(\bar{P}_S, DT_{max}). \quad (3.60)$$

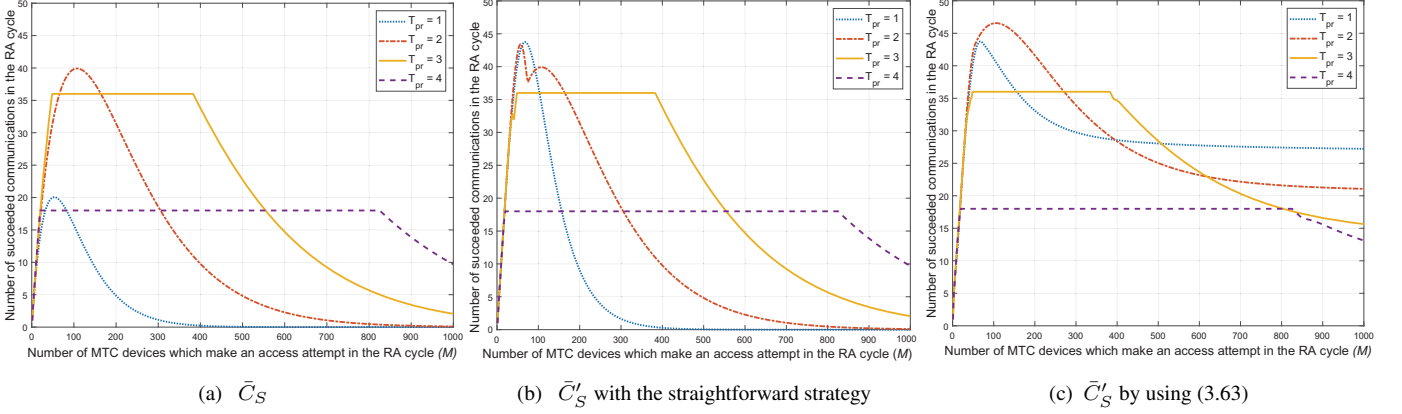
In Fig. 3.34a we plot  $\bar{C}_S$  versus  $M^4$  for any  $T_{pr}$  value. We observe that the straight lines are caused by the saturated PUSCH resources (e.g., with  $T_{pr} = 3$ ,  $\bar{C}_S = DT_{max} = 36$ , for  $M$  ranging from 50 to 380).

#### 3.4.2 Our proposal

Starting from (3.60), we note that when  $P_S < DT_{max}$ , i.e., the number succeeded access attempts is less than the amount of resources available for the data transmission,  $P_S$  out of  $DT_{max}$  data transmissions are scheduled in a collision free mode to the succeeded preambles, while  $DT_U = DT_{max} - P_S$  available data transmissions remain unused in the PUSCH.

<sup>4</sup>As shown in [68], the values assumed by  $M$  are consistent with the MTC arrival distribution suggested by the 3GPP in [67], with  $N_{MTC} = 100000$ .

### 3.4. Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH



**Figure 3.24:** Number of successful transmissions in an RA cycle vs the number of MTC devices ( $M$ ) which carry out the RA procedure for different  $T_{pr}$  values.

Our strategy aims to rise  $\bar{C}_S$  trying to serve also the collided preambles by assigning to them the  $DT_U$  resources in a proper manner. In this case, the data transmission resources cannot be allocated in a collision free mode, since there is a 1 to many relationship between each collided preamble and the related MTC devices. So, the scheduler needs to allocate, in a contention based mode, a pool of  $R_\theta$  resources consisting of  $DT_{P_C}$  data transmissions per collided preamble, with  $DT_{P_C} > 1$ . Specifically, the gNB reserves a single pool of  $DT_{P_C}$  resources to the MTC devices that have transmitted the same collided preamble. Then, each related MTC device draws independently a uniform random number  $g \in \{1, \dots, DT_{P_C}\}$  and sends its data in the  $g$ th data transmission resource. Similarly to the preamble transmission in the PRACH, there is a non-zero data transmission collision probability. We emphasize that, if the additional communication is successful, the gNB knows which device has transmitted data by inspecting the related UL context. Hence, the gNB is able to ACK the data to the related device and to further process and forward the data to the expected network entity. Conversely, if the additional communication is collided, the MTC device waits in vain the ACK from the gNB within the related waiting window. Then, it re-attempts the RA procedure during a randomly selected PRACH inside the Backoff Window ( $B_W$ ).

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

#### Analysis and Formulation of Proposed Strategy

The number of successful additional data transmission ( $A_S$ ) is a random variable that depends on  $DT_{P_C}$  and on the number of MTC devices that have transmitted the same collided preamble. We denoted it as  $K_C$  and it is termed as "collision coefficient". The average number of successful additional data transmissions per collided preamble ( $\bar{A}_{S,P_C}$ ) can be calculated as (3.58), with  $M = K_C$  and  $L = DT_{P_C}$ , as follows:

$$\bar{A}_{S,P_C} = \begin{cases} K_C (1 - 1/DT_{P_C})^{K_C-1} & \text{if } DT_{P_C} > 1 \\ 0 & \text{otherwise.} \end{cases} \quad (3.61)$$

Since  $\bar{A}_{S,P_C}$  and  $\bar{P}_C$  are independent random variables, the average total number of successful additional data transmissions is  $\bar{A}_S = \bar{A}_{S,P_C} \bar{P}_C$ . Finally, we redefine the average number of succeeded communications ( $\bar{C}'_S$ ) as:

$$\bar{C}'_S = \begin{cases} \bar{P}_S + \bar{A}_S & \text{if } \bar{P}_S < DT_{max} - 1 \\ \bar{C}_S & \text{otherwise.} \end{cases} \quad (3.62)$$

To calculate  $\bar{C}'_S$ , we need to derive  $K_C$  and  $DT_{P_C}$ . Let us start by calculating the total number of preambles out of  $L$  transmitted by  $M$  MTC devices. It is denoted as  $P_T$ . At this aim, we define  $Y_i$ , with  $i \in \{1, \dots, L\}$ , as Boolean random variables, each one equal to 1 if and only if preamble  $i$  has been selected by at least one MTC device. Since  $Y_i$  are  $L$  random variables identically distributed, the mean value of  $P_T$  can be calculated as  $\bar{P}_T = E \left\{ \sum_{i=1}^L Y_i \right\} = LE \{Y_i\}$ , where  $E \{Y_i\} = \Pr(Y_i = 1) = 1 - (1 - 1/L)^M$ , for each  $i$ .

Now, the average number of collided preambles can be easily calculated as  $\bar{P}_C = \bar{P}_T - \bar{P}_S$ , and the average number of MTC devices that have selected a collided preamble can be approximated as  $\bar{K}_C \simeq (M - \bar{P}_S) / \bar{P}_C$ . Then, having  $\bar{DT}_U$  unused data transmissions and  $\bar{P}_C$  collided preambles, the straightforward strategy is to equally distribute the  $\bar{DT}_U$  resources among the  $\bar{P}_C$  preambles, i.e.,  $\bar{DT}_{P_C} = \bar{DT}_U / \bar{P}_C$ .

Having derived  $\bar{K}_C$  and  $\bar{DT}_{P_C}$ , we can calculate  $\bar{A}_{S,P_C}$  by using (3.61), where we set  $K_C = \bar{K}_C$  and  $DT_{P_C} = \bar{DT}_{P_C}$ . This treatment is valid

### 3.4. Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH

when using the drift approximation [97]. Then, we calculate  $\bar{C}'_S$  by using (3.62) and plot the variation of  $\bar{C}'_S$  versus  $M$  for any  $T_{pr}$  value in Fig. 3.34b. As expected, the only work areas that permit to improve the number of successful transmissions are outside of straight lines, i.e., where the PUSCH resources are not saturated. More specifically, the best improvement in terms of  $\bar{C}'_S$  is obtained by the  $T_{pr} = 1$  dimensioning, thanks to the high availability of  $\overline{DT}_U$  PUSCH resources. As regards the  $T_{pr} = 2$  dimensioning, it exhibits a modest gain with respect to  $\bar{C}_S$  only up to  $M = 74$ , since for  $M \geq 75$ , it results  $\bar{A}_S = 0$ , as can be verified. Instead, the remaining dimensionings show a very slight improvement only before the related straight lines, that can be neglected. Moreover, all the dimensionings do not show improvements for high  $M$  values because the higher  $M$ , the higher  $\bar{P}_C$  and  $\bar{K}_C$ , the lower  $\overline{DT}_{P_C}$ . So, the straightforward strategy of serving all the collided preambles leads to  $\bar{A}_S = 0$ .

To overcome this issue, we propose an enhanced strategy to maximize the term  $\bar{C}'_S$ . We note that, given a fixed  $T_{pr}$  value,  $M$ , and  $\theta_{max}$ , the values of  $\bar{P}_S$  and  $DT_{max}$  are determined by (3.58) and (3.59), respectively. Then, in (3.62) the only term that can be maximized is  $\bar{A}_S$ . Unlike the straightforward strategy, we propose to divide, on average, the unused PUSCH resources among a proper subset of collided preambles (i.e.,  $\bar{P}_{C_S} \leq \bar{P}_C$ ). For this reason, we can formulate the goal as follows: *"Given  $M$  and  $T_{pr}$ , the aim is to maximize  $\bar{A}_S$  by varying the number of collided preambles ( $\bar{P}_{C_S}$ ) to which the unused PUSCH resources should be allocated in the contention way:*

$$\max_{\bar{P}_{C_S} \in [0, \min(\overline{DT}_U/2, \bar{P}_C)]} \{ \bar{A}_S \} .'' \quad (3.63)$$

We calculated  $\bar{A}_S$  numerically and in Fig. 3.24c, we plot the variation of  $\bar{C}'_S$  versus  $M$  for any  $T_{pr}$  value. The  $T_{pr} = 1$  and  $T_{pr} = 2$  dimensionings achieve a remarkable improvement for any value of  $M$  in comparison with Figs. 3.24a and 3.24b. As regards  $T_{pr} = 3$ , a light improvement occurs when  $M \geq 500$ .  $T_{pr} = 4$  has no significant room for improvement. The lower the  $T_{pr}$  value, the larger  $\overline{DT}_U$ , the higher the improvement.

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

### The PUSCH Resource Reallocation Algorithm (PRRA)

In light of the above theoretical average advantages, we implement the proposed strategy by presenting the PUSCH Resource Reallocation Algorithm (PRRA). It allows the scheduler to re-allocate the PUSCH resources efficiently, inside a given RA cycle, known  $M$ ,  $P_S$ ,  $P_C$  and  $T_{pr}$ . The purpose is to find the optimal  $P_{C_S}^* \leq P_C$  value that maximizes the expected value of  $A_S$  inside that RA cycle, and the relative  $DT_{P_C}^*$  value. Then, the  $P_{C_S}^*$  collided preambles will be chosen randomly. The algorithm is described in pseudo-code 3, and its time complexity is  $O(n)$ , where  $n = \min\{DT_U, P_C\}$ .

---

#### Pseudo-code 3 PUSCH Resources Reallocation Algorithm

---

**Inputs:**  $M, Q, K_{max}, S, \theta_{max}, T_{pr}, T_{ra}, P_S$ , and  $P_C$

**Iteration:**

- 1: calculate  $DT_{max}$  by using (3.88),  $K_C = (M - P_S)/P_C$ , and  $DT_U = DT_{max} - P_S$ ;
  - 2: create an auxiliary vector  $\mathbf{A} = [1, \dots, \min\{DT_U, P_C\}]$ ;
  - 3: calculate  $\mathbf{B} = \lfloor DT_U \mathbf{A}^{\circ(-1)} \rfloor$ , where  $(\cdot)^\circ$  is the Hadamard power operator;
  - 4: calculate  $\mathbf{C} = K_C \left( \mathbf{J}_{1,|\mathbf{B}|} - \mathbf{B}^{\circ(-1)} \right)^{\circ(K_C - 1)}$  where  $\mathbf{J}_{1,|\mathbf{B}|}$  is the all-ones matrix of size  $1 \times |\mathbf{B}|$ ;
  - 5: calculate  $\mathbf{D} = \mathbf{A} \circ \mathbf{C}$ , where  $\circ$  is the Hadamard product;
  - 6: calculate  $\bar{A}_S = \max\{\mathbf{D}\}$  and  $P_{C_S}^* = \arg \max\{\mathbf{D}\}$ ;
  - 7: calculate  $DT_{P_C}^* = \mathbf{B}[P_{C_S}^*]$ .
- 

However, in order to apply the PRRA, the gNB should know for each RA cycle, inter alia, the total amount of MTC devices  $M$  attempting to access. Clearly,  $M$  is not known at the base station a priori, and therefore this parameter should be estimated. From an empirical study, we derived in [68] the estimated number of attempting devices ( $\tilde{M}$ ) as function of  $P_S$  and  $P_C$ , which the gNB knows at Step 2 of the RA procedure:

$$\tilde{M} = P_S + 2 \cdot 1.97^{P_C/L} P_C. \quad (3.64)$$

This estimate exploits the fact that the process of total access attempts in two consecutive RA cycles is strongly correlated, as reported in [68] and proved in [16]. The accuracy of (3.64) is verified by means of the simulation results in Section 3.4.3.

### A New Enhanced Dynamic Uplink Resource Dimensioning

Observing Fig. 3.24c, we note that there is not a unique value of  $T_{pr}$  which optimizes  $\bar{C}'_S$  for any number of attempting devices  $M$ . In [16,68]

### 3.4. Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH

we showed, with a long time evaluation (i.e., considering thousands of successive RA cycles), that the optimal PRACH/PUSCH resource allocation can be achieved only if a dynamic load-aware control is applied. More specifically, we proposed to dynamically set the  $T_{pr}$  dimensioning for each RA cycle, on the basis of  $M$ . In [68] the proposed Dynamic Uplink Resource Dimensioning (DURD) was applied to the  $\bar{C}_S$  values derived from (3.87) and reported in Fig. 3.24a. In this letter, we apply the same approach, considering the new curves derived from (3.62) and (3.63), and shown in Fig. 3.24c.

"Given  $M$ , the goal is to find the optimal  $T_{pr}$  value so that

$$T_{pr}^* = \arg \max_{T_{pr} \in \{1, \dots, T_{ra}-1\}} \{\bar{C}'_S\}, \quad (3.65)$$

where the dependence on  $M$  is contained in  $\bar{C}'_S$ ." Choosing  $T_{pr} = T_{pr}^*$ , we have the optimal resource dimensioning in the RA cycle. In this way, we obtain an Enhanced Dynamic Uplink Resource Dimensioning system, termed EDURD and, as shown in Fig. 3.34c, the optimal dimensioning results  $T_{pr}^* = 2$  for  $M \leq 271$ ,  $T_{pr}^* = 3$  for  $272 \leq M \leq 507$ , and  $T_{pr}^* = 1$  for  $M \geq 508$ . Compared to Fig. 3.34a, the high  $M$  values correspond to the lowest  $T_{pr}^*$  value because for  $T_{pr} = 1$  a large portion of the preambles in the PRACH are collided and, consequently, the succeeded transmissions are related to the reallocation in the PUSCH, that works better with  $T_{pr} = 1$ . In summary, adopting the EDURD system, for each RA cycle, on the basis of  $M$ , the gNB sends in broadcast the optimal  $T_{pr}^*$  value for the next RA cycle. Then, on basis of  $P_S$  and  $P_C$  detected, the gNB scheduler assigns one  $R_\theta$  resource to each of the  $P_S$  MTC devices in a collision free mode, and  $DT_{P_C}^*$  resources to  $P_{C_S}^* \leq P_C$  preambles in a contention-based mode, that were estimated as output of the PRRA.

#### 3.4.3 Performance Evaluation

In this section, we consider the following uplink dimensioning systems, where the PUSCH resources are assigned on basis of the SCMA technique. We denote them as:

- $S_i$ , with  $i \in \{1, \dots, 4\}$ , as the system with static uplink dimensioning  $T_{pr} = i$ ;

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

**Table 3.3:** *Simulation Parameters*

Parameter	Symbol	Value
Preambles reserved for contention-based procedure	$L_0$	54
Number of total MTC devices	$N_{MTC}$	40000, 50000, 75000, 100000
Number of data transmission requests per MTC device		1
Maximum number of RA attempts	$M_A$	10
Data transmission size	$\theta_{max}$	160 bits
Arrival distribution	Beta	$\alpha = 3, \beta = 4$ $T_{Arrival} = 10s$
RA cycle duration	$T_{ra}$	5 time slots
Simulation Length	$T_{sim}$	10s
SCMA parameters	$Q, K_{max}, S$	4, 3, 2
Backoff Window	$B_W$	20ms
Constellation Points	$I$	4

- *DURD* [68] as the system with optimal dynamic uplink dimensioning based on Fig. 3.34a;
- *eDURD* [68] as the *DURD* system based on the predictive estimation of  $M$  by using (3.64).
- $Ex$ , with  $x = \{S_i, DURD, eDURD\}$ , as the enhanced system  $x$  with the contention-based additional transmissions in the PUSCH, by using the PRRA.

We compare the performance of above systems by simulation in MATLAB environment. The simulation parameters are provided in Table 3.6, and results are averaged over 50 independent simulations. We introduce two metrics.

1) *System Throughput Gain*. It is the percentage throughput gain compared to  $S_1$ , i.e.,  $\Delta Th_x = [(Th_x - Th_{S_1}) / Th_{S_1}] \cdot 100$ , where  $Th_x$  is the system throughput in terms of number of successful communications achieved by means of system  $x$ .

2) *Energy Consumption Indicator*. It is the average number of times in which an MTC device enters in the RX/TX active state [98], denoted as  $\bar{E}_I$ . For each MTC device,  $E_I$  is calculated as  $E_I = E_{RA} + E_T$ , where  $E_{RA} \in \{1, \dots, M_A\}$  is the number of access attempts, and  $E_T$  counts the total number of contention-based and contention-free transmissions



### 3.4. Enhancement of the Resource Allocation in mMTC Scenarios through a Contention-Based transmission in the PUSCH

in the PUSCH. In particular,  $E_T \in \{0, 1\}$  for the non-Enhanced systems, while  $E_T \in \{0, \dots, M_A\}$  for the other ones.

Figs. 3.25 and 3.26 show  $\Delta Th_x(\%)$  and the Energy Consumption Indicator  $\bar{E}_I$ , respectively. Let us start considering a moderate load, i.e.,  $N_{MTC} = 40000$ . Among all the static dimensioning systems, only  $ES_1$  and  $ES_2$  show a throughput gain with respect to the related standard systems. Instead, as expected,  $ES_3$  does not show improvement with respect to  $S_3$  because, in this light traffic conditions,  $M \leq 500$  for each RA cycle. By comparing  $ES_2$  with DURD, we underline that the PRRA improves the performance of a static system also compared to an optimized dynamic solution. As regards the EDURD, the results are the same of the ones obtained by  $ES_2$ , being the optimal dimensioning  $T_{pr}^* = 2$  for any RA cycle. Finally, EeDURD obtains the same results as EDURD, confirming that the predictive estimation of  $M$  is efficient in the considered working zone. As regards a medium-high load ( $N_{MTC} = 50000$ ), compared to the previous case, also  $ES_3$  shows a gain with respect to  $S_3$ . As regards the dynamic systems, EDURD shows the similar performance of the  $ES_3$ , because the optimal dimensioning is equal to  $T_{pr}^* = 3$  for almost the totality of the RA cycles. In addition, it achieves the best performance in terms of energy consumption. Also in this case, the predictive estimation of  $M$  is effective. In the high load scenarios ( $N_{MTC} \in \{75000, 100000\}$ ), among the static systems the best performances in throughput are obtained with  $ES_1$ . The EDURD system has a good gain with respect to any static system. In particular, the highest gain is achieved for  $N_{MTC} = 75000$  because the optimal dimensioning ranges from  $T_{pr}^* = 3$  to  $T_{pr}^* = 1$ , while for  $N_{MTC} = 100000$  the gain is lower, because the optimal dimensioning is  $T_{pr}^* = 1$  for many RA cycles. The energy consumption is slightly higher with respect to the other dynamic systems, but in line with the best static energy consumption. As regards the estimated versions under very high load, for EeDURD the predictive estimation does not work so well, because the optimal dimensioning is  $T_{pr}^* = 1$  and only  $L = 54$  preambles are available. The  $L$  value is too much low with respect to the number of attempting MTC devices and the estimated value of  $K_C$ , as a function of  $P_C$  and  $L$ , is more approximated. In stead, the predictive estimation in the eDURD system shows good performance because the optimal dimensioning is  $T_{pr}^* = 4$ ,

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

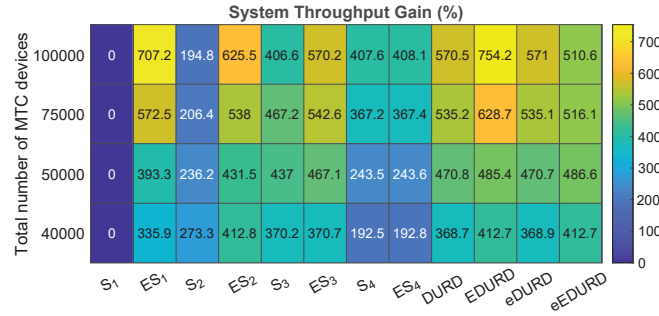


Figure 3.25: Throughput Gain under different systems and  $M$

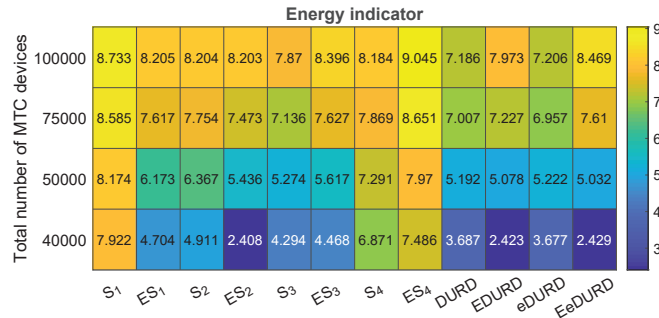


Figure 3.26: Energy Consumption Indicator under different systems and  $M$

see Fig. 3.34a.

## 3.5 A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

### 3.5.1 Motivations and related work

In order to work well, all the works mentioned above and many other ones available in literature need to know the current PRACH traffic load<sup>5</sup>, denoted with  $M$ . It represents the total access requests per RA cycle, i.e., it is the sum of new access requests and the re-transmitting ones. We underline that these works adopt either the 4-step or the 2-step connectionless RA procedure. However, in both the RA procedures, for each RA cycle the only information known to the BS are: the amount of collided preambles ( $P_C$ ) and the amount of succeeded preambles ( $P_S$ ), out of the

<sup>5</sup>A conventional analysis of the traffic load in the PUSCH, as largely available in literature for enhanced Mobile BroadBand (eMBB) scenarios, is not significant of the congestion in an mMTC scenario, where the bottleneck is observed at the connection setup phase (i.e., in the PRACH). In fact, a light traffic load in the PUSCH does not imply that the system is far from congestion, because this condition could occur either when a small number of access requests are carried out in the PRACH (i.e., light PRACH traffic load) or due to a high number of access request collisions (i.e., when the PRACH is overloaded).

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

$L$  available ones. The BS cannot know how many MTC devices have selected each collided preamble. So, the number of attempting MTC devices, i.e.,  $M$ , in the RA cycle is unknown and should be properly estimated.

Several works (e.g., [52,54,55,62]) ideally assume the perfect knowledge of the PRACH traffic load. Some works [53, 56, 68] propose conventional ways to correctly estimate the current PRACH traffic load  $M$ , taking into account the information available at the BS. Other works (e.g., [58]) exploit AI techniques to predict future traffic statistics, based on the knowledge of the current PRACH traffic load, but still estimated with conventional methods available in the literature.

More specifically, in [53], the authors propose an algorithm to determine the ACB factors. They analyze only the case when  $M \geq L$ , since when  $M < L$ , no ACB check is performed, and all the MTC devices will attempt random access upon activation. In particular, their approach is based only on the analysis of the number of collided preambles  $P_C$ , and does not take into account the number of succeeded preambles  $P_S$ . They apply an ACB factor  $r$  in the RA cycle  $j$ , and then, on the basis of the observed collisions, they estimate the proper ACB factor  $\tilde{r}$  that should have been applied in the same RA cycle  $j$ . This value is calculated as:

$$\tilde{r} = \min \left\{ 1, r \left( 1 + \frac{[P_C - L(1 - 2e^{-1})]e}{2L} \right)^{-1} \right\}. \quad (3.67)$$

Finally, on the basis of  $\tilde{r}$ , they estimate the  $M$  value as

$$\tilde{M} = \frac{L}{\tilde{r}}. \quad (3.68)$$

Then, on basis of  $\tilde{M}$ , they predict the ACB factor  $\tilde{r}$  that will be used in RA cycle  $j + 1$ . This procedure is cyclically applied for each RA cycle.

---


$$Pr\{P_S, P_C \mid M, L\} =$$

$$\frac{1}{L^M} \binom{M}{P_S} \binom{L}{P_S} P_S! \binom{L - P_S}{L - P_S - P_C} \sum_{i=0}^{P_C} (-1)^i \binom{P_C}{i} \sum_{w=0}^i \binom{M - P_S}{w} \binom{i}{w} w! (P_C - i)^{M - P_S - w} \quad (3.66)$$

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

We underline that this simple estimate works properly only after several RA cycles in which the arrival process does not fluctuate.

In [56], the authors use a maximum likelihood estimation method to calculate  $M$ . In detail, they derive  $Pr\{P_S, P_C \mid M, L\}$ , which is the probability that  $P_S$  preambles are successfully transmitted and  $P_C$  preambles are collided given that each of the  $M$  devices sends one of the  $L$  available preambles. This probability is shown in (3.66) at the bottom of the page. Finally, the BS can estimate the number of PRACH attempts in the  $j$ th RA cycle as

$$\tilde{M} = \arg \max_{M'} (Pr\{P_S, P_C \mid M', L\}). \quad (3.69)$$

We note that the computational complexity is very high, since for each possible  $M$  value, it is needed the calculation of the product of several binomial coefficients, inside a nested summation, where the upper limit is  $P_C$ . Clearly, when the number of collided preambles is high, this approach becomes computationally expensive in both time and space.

Also, a very simple method for estimating  $M$  for each RA cycle is applied in [16, 68]. It is derived from an empirical study and is equal to:

$$\tilde{M} = P_S + 2 \cdot 1.97^{\frac{P_C}{L}} P_C. \quad (3.70)$$

In particular, this empirical method was proved to be sufficiently useful in [68] to determine the optimal Dynamic Uplink Resource Dimensioning (DURD).

Finally, in [58] the authors propose Deep Reinforcement Learning (DRL) algorithms to jointly and dynamically adapt the parameters required by their hybrid RA scheme. Their strategy uses a Recurrent Neural Network (RNN) model to predict the future PRACH traffic statistic. In order to enable online updating, they implement the supervised learning method proposed in [99]. Specifically, the weights of the RNN are updated by using the current  $\tilde{M}$  value obtained through a state-of-the-art estimation method [53, 56], without the introduction of any AI-aided technique.

#### 3.5.2 Contributions

In light of what described, it is clear how important is to accurately estimate the traffic load for several control schemes, which deal with reg-

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

ulating access in massive IoT scenarios. Moreover, since the estimation of  $M$  should be carried out for each RA cycle, it must require a low time computational complexity, for making these schemes viable. These needs will be increasingly stringent in the prospect of the future 6G networks, where  $10^7$  MTC devices per square kilometer are expected [100]. Therefore, in this paper we address the problem of estimating the traffic load  $M$  when the radio access follows a grant-based approach in an accurate and realistic way. We propose a new AI-based traffic load estimation method which exploits only the information related to the PRACH, available at the gNB. We chose a conventional Deep Neural Network (DNN) architecture used for regression tasks and we properly derived its topology for achieving the best accuracy with the lowest computational complexity.

#### 3.5.3 System Model for the PRACH traffic load prediction

In this work, for defining our DNN-based estimation method and then to compare its performance against certain benchmark ones, we consider an actual and viable 5G scenario with one gNB which supports a group of  $N_{MTC}$  MTC devices using a licensed band of bandwidth  $B$  exclusively dedicated to the mMTC service. In particular, we consider  $B = 1.08$  MHz, that is, the the minimum amount of bandwidth required to support the mMTC service [35].

As regards the uplink radio interface, we adopt the smallest 5G NR numerology, i.e., subcarrier spacing of 15 kHz with normal cyclic prefix, which corresponds to a time slot duration ( $T_s$ ) of 1 ms, containing 14 OFDM symbols. This choice has been made because small subcarrier spacing results in longer symbol duration and lower overhead. Therefore, delay tolerant services, such as several MTC services, can benefit from small subcarrier spacing to reduce bandwidth consumption [95].

In 5G NR there are 64 orthogonal preambles defined in each time-frequency PRACH occasion, and  $L_0$  are available for the contention-based requests. We set  $L_0 = 54$ , that is typically reserved for contention-based RA [101]. In this work, we adopt the NR preambles Format 0, i.e., preambles are mapped to 839 subcarriers of 1.25 kHz (i.e., 1.05 MHz), and each preamble lasts 1 ms [35].

A typical RA cycle lasts  $T_{ra} = 5T_s$  and the uplink radio resources

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

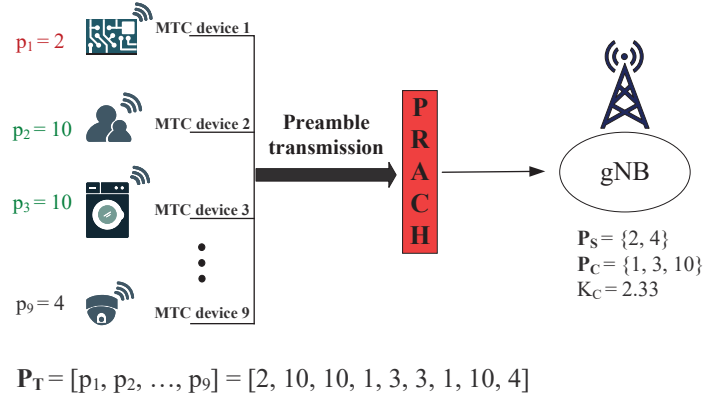
are divided into a PRACH subset and a PUSCH subset, as shown in Fig. 3.1. The PUSCH is used to transmit data and occupies, in the frequency domain, the entire bandwidth  $B$ . In the time domain PRACH lasts  $T_{pr} \in \{1, 2, \dots, T_{ra} - 1\}$  time slots and PUSCH the remaining  $T_{pu}$  slots. In summary, for each RA cycle we have  $T_{ra} = T_{pr} + T_{pu}$ . So, the total number of available preamble sequences in a RA cycle is  $L = L_0 T_{pr}$ .

As regards the RA procedure, we consider the 2-step connectionless RA procedure proposed in [16]. Also, for the PUSCH resources, we adopt the Sparse Code Multiple Access (SCMA) technique on the OFDMA grid. The SCMA technique appeared to be the most promising Non Orthogonal Multiple Access (NOMA) technique to multiplex a large number of small-size data [96]. It is a code division multiple access scheme, characterized by sparse codebooks built on multidimensional constellations. By using SCMA scheme multiple users can transmit on the same Resource Element (RE) with different codebooks. According to the SCMA encoder defined in [77], one  $SCMA_{block}$  is the minimum resource quantity which can be shared by different transmissions, called layers. In addition, since for the MTC devices a very small amount of data is expected, we set for each transmission request an upper bound value,  $\theta_{max}$  bits, and we assume that the gNB assigns to each successful access attempt, the PUSCH resources, if any, enough to transmit  $\theta_{max}$  bits.

#### 3.5.4 The proposed DNN-based traffic load estimation method

In this section we propose a new DNN-based method to estimate  $M$ . In this regard, we describe the Step 1 of the RA procedure by means of an example reported in Fig. 3.27. As shown,  $M = 9$  MTC devices transmit one randomly chosen preamble out of  $L$  in the PRACH. Let  $\mathbf{P}_T = [p_1, \dots, p_M]$  denote the vector containing the preambles transmitted by the MTC devices. Specifically, in the example  $\mathbf{P}_T = [2, 10, 10, 1, 3, 3, 1, 10, 4]$ , i.e., preambles 2 and 4 have been transmitted by only one MTC device, whereas preambles 1, 3, and 10 by two or more MTC devices. At the end of the Step 3 of the conventional 4-step RA procedure, or at the end of the Step 1 of the connectionless 2-step RA procedure, the gNB detects the preambles 2 and 4 as succeeded, whereas preambles 1, 3, and 10 as collided. We denote  $\mathbf{P}_S$  as the set of succeeded preambles with cardinal-

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond



**Figure 3.27:** Example of Step 1 of the RA procedure.

ity  $P_S$ , and  $\mathbf{P}_C$  as the set of collided preambles with cardinality  $P_C$ . In the example, we have  $\mathbf{P}_S = \{2, 4\}$  and  $\mathbf{P}_C = \{1, 3, 10\}$ . The value of  $M$  is unknown to the gNB, since it does not know the cardinality of each collided preamble, i.e., how many MTC devices have transmitted it. Let  $M_p$  be the cardinality of preamble  $p$ . Obviously,  $M_p \in \{0, 1, \dots, M\}$ . We define  $K_C$  as the "collision coefficient", and it is calculated as follows

$$K_C = \frac{1}{P_C} \sum_{p \in \mathbf{P}_C} M_p. \quad (3.71)$$

In the example,  $K_C \approx 2.33$ . Clearly,  $K_C = f(P_S, P_C, L)$ . It is straightforward to derive that  $M = P_S + K_C P_C$ .

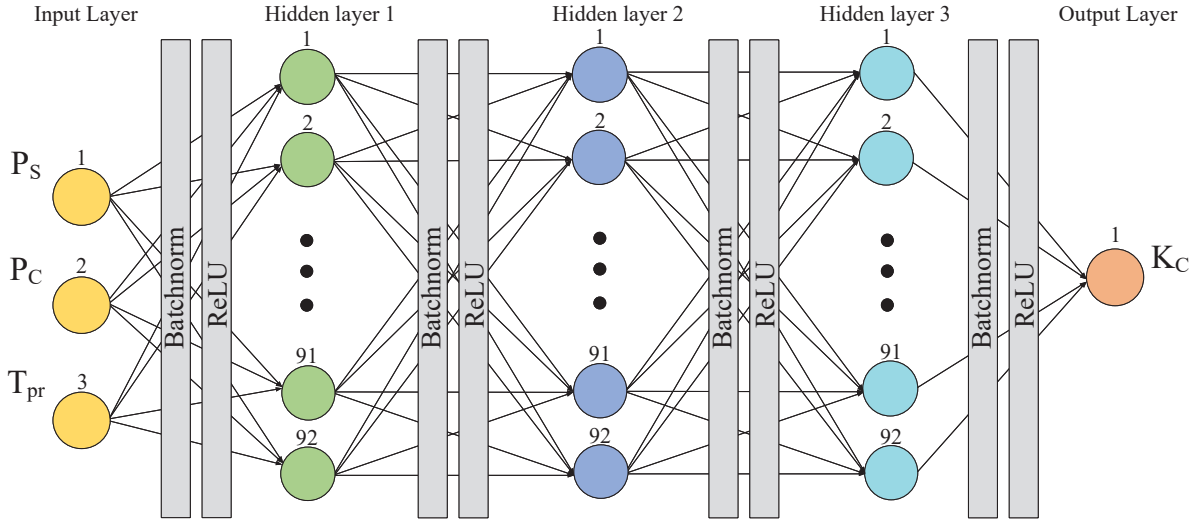
In the following we denote with  $\tilde{x}$  the estimated value of  $x$ . Since  $L = L_0 T_{pr}$ , we propose to estimate  $\tilde{K}_C$  by means of a DNN, denote as  $DNN_K$ , with input values  $P_S$ ,  $P_C$ , and  $T_{pr}$ . Then, the estimated  $M$  value is calculated as

$$\tilde{M} = P_S + \tilde{K}_C P_C. \quad (3.72)$$

#### $DNN_K$ architecture

The  $DNN_K$  architecture is depicted in Fig. 3.28. As shown, it is composed of  $L = 5$  layers, where the first layer is called input layer, the last layer is called output layer, and the other three layers are called hidden layers. Each layer  $l$  is made of  $N_l$  nodes called "neurons", and we have  $N_1 = 3$ ,  $N_2 = N_3 = N_4 = 92$  and  $N_5 = 1$ . For each layer  $l$ , with  $l > 1$ , the input data  $\mathbf{x}_l \in \mathbb{R}^{N_{l-1}}$  is multiplied by weight  $\mathbf{w}_l \in \mathbb{R}^{N_{l-1} \times N_l}$  and all the

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.28:** Architecture of  $DNN_K$

multiplied data plus the bias  $\mathbf{a}_l \in \mathbb{R}^{N_l}$  are added up and the sum is passed through an activation function  $\phi_l$  to generate the output  $\mathbf{y}_l \in \mathbb{R}^{N_l}$ . So, this computation related to layer  $l$  can be described by:

$$\mathbf{y}_l = \phi_l(\mathbf{w}_l^T \mathbf{x}_l + \mathbf{a}_l). \quad (3.73)$$

In our network, we adopt as input vector  $\mathbf{x}_1 = [P_S, P_C, T_{pr}]^T$ , as output  $y_5 = \tilde{K}_C$ , and as  $\phi_l$  the Rectified Linear Unit (ReLU).

We underline that this DNN architecture is tailored for the considered case study. In fact, the topology was not arbitrary chosen but determined during the training procedure through a large number of experiments, by varying the number of hidden layers and neurons per layer.

#### Training and Testing Procedure

The proposed model needs to fit on a training set  $\mathbf{S}_{TR}$ , which is a set containing labeled data that are used in the supervised learning to derive the parameters of the model. In practice, the training set contains  $S_{TR}$  examples, where the  $i$ th example consists of an input vector  $\mathbf{x}_{1_i}$  and the corresponding output value  $K_{C_i}$ , commonly denoted as ground truth data or label. Using the training set, the proposed DNN structure needs to be trained such that the differences between the estimated values  $\tilde{K}_{C_i}$  and the related ground truth values  $K_{C_i}$  are minimized, i.e., the most accurate reconstructed output values are obtained. Since we are addressing a



### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

regression problem, we adopt as loss function the Mean Squared Error (MSE), which can be defined as:

$$Loss = \frac{1}{S_{TR}} \sum_{i=1}^{S_{TR}} (K_{C_i} - \tilde{K}_{C_i})^2. \quad (3.74)$$

So, in order to minimize the loss function, each element  $w_{i,j} \in \mathbf{w}_l$ , for each layer  $l$  is updated using the Gradient Descent (GD) method until the loss function reaches the lowest value. Specifically, by using the back propagation method, the weights and biases of the network are updated according to their derivatives of the loss function, which is expressed as follows:

$$w'_{i,j} = w_{i,j} - \alpha \frac{\partial Loss}{\partial w_{i,j}} \quad (3.75)$$

$$a'_j = a_j - \alpha \frac{\partial Loss}{\partial a_j} \quad (3.76)$$

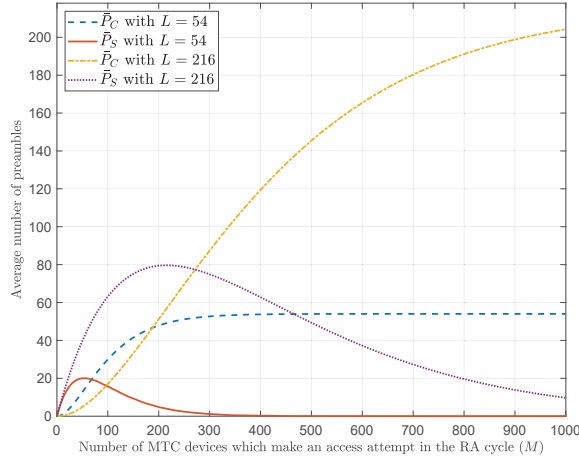
where  $i \in \{1, 2, \dots, N_{l-1}\}$ ,  $j \in \{1, 2, \dots, N_l\}$ , and  $\alpha$  is the learning rate. To gain better training performance and decrease computational complexity, in our networks we adopt the GD training optimizer ADAM with the  $\alpha$  value set to 0.0003. Moreover, for avoiding vanishing and exploding gradient problems, variable initialization and batch normalization are also taken into consideration. Regarding the variable initialization, we adopted the Xavier initialization [102], while as batch normalization, the approach adopted is the normalization via mini batch statistics [103].

Finally, once the model fits on the training set, a test set  $\mathbf{S}_{TE}$  is used only to assess the performance (i.e., generalization) of the model.

#### Dataset creation

The use of proper training and test sets is a very important part of building a supervised model with good generalization performance. In fact, the performance achievable depends on the availability of a large quantity of data and on their quality. As regards the quantity, both the sets should be representative of the entire population. As concerns the quality, the neural network should be trained and tested in a proper working area where it is possible to obtain a unique input-output relationship, i.e.,

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.29:** Average number of succeeded and collided preambles vs  $M$  for different  $L$  values.

the input data must not lead to an indeterminate output value. However, due to the unavailability in literature of a dataset with these characteristics for our task, we provide a new dataset created by simulations.

With the aim of ensuring the quality of the data, we first examined whether there exist areas of indeterminateness and, accordingly, determined the appropriate simulation campaign to carry out. More specifically, we analyzed the limits of the univocal estimate of  $M$  as a function of the three input values  $P_S$ ,  $P_C$ , and  $T_{pr}$  (i.e.,  $L$ ), performing an *a priori* analysis. Given  $M$  as the total number of access attempts in a RA cycle,  $P_S$  is a random variable and its expected value is:

$$\bar{P}_S = M \left(1 - \frac{1}{L}\right)^{M-1}. \quad (3.77)$$

As regards the number of collided preambles ( $P_C$ ), its expected value is equal to

$$\bar{P}_C = L \left[ 1 - \left(1 - \frac{1}{L}\right)^M - M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1} \right]. \quad (3.78)$$

The proofs of (3.77) and (3.78) are reported in 3.5.8. In Fig. 3.29 we plot  $\bar{P}_C$ , and  $\bar{P}_S$  vs  $M$ , when  $T_{pr} = 1$  and  $T_{pr} = 4$ , i.e., the available preambles  $L$  are  $L = L_0 T_{pr} = 54$  and  $L = 216$ , respectively.

As shown, when  $T_{pr} = 1$ , for high traffic load values (about  $M > 300$ )  $\bar{P}_C$  tends to the maximum and constant value equal to the number

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

of available preambles,  $\bar{P}_C = L = 54$ , and  $\bar{P}_S$  results 0. So, the value of  $M$  cannot be determined as function of the above parameters, since it is not unique. Conversely, for  $L = 216$ , given a pair of values  $\bar{P}_C$ , and  $\bar{P}_S$  we can determine a single value of  $M$ .

Since the DNN should not been trained in areas of indeterminateness, we need to determine the correct range of values of  $M$ , depending on  $T_{pr}$ , so that  $\Pr\{P_C = L\}$  can be neglected. For this reason, we derive also in 3.5.8 the probability mass function (pmf) of  $P_C$ . It results:

$$p_{P_C}[k] = \binom{L}{k} \left[ 1 - \left(1 - \frac{1}{L}\right)^M - M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1} \right]^k \left[ 1 - \left(\frac{1}{L}\right)^M + M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1} \right]^{L-k} \quad (3.79)$$

where  $k = 0, 1, \dots, L$ . Given the pmf (3.79), we set the probability that  $P_C = L$  can be neglected at  $10^{-6}$ , i.e.,

$$\Pr\{P_C = L_0 T_{pr}\} \leq 10^{-6}. \quad (3.80)$$

Then, for each value of  $T_{pr} = h$ , with  $h = \{1, 2, 3, 4\}$ , we found the related maximum value of  $M$ , denoted as  $M_{max}^h$ . We obtain,  $M_{max}^1 = 151$ ,  $M_{max}^2 = 393$ ,  $M_{max}^3 = 669$ , and  $M_{max}^4 = 966$ . Therefore, as analytically derived, given  $T_{pr} = h$ , the working area for properly training the model concerns  $M \in \{1, 2, \dots, M_{max}^h\}$ .

In light of the obtained results, we divided the simulation campaign into 4 groups, where the  $h$ th group is related to the dimensioning  $T_{pr} = h$ . For each group, we fixed one value of  $M \in \{1, 2, \dots, M_{max}^h\}$  and simulated the first step of the RA procedure described at the beginning of this Section and drawn in Fig. 3.27. As a result, we obtained  $P_S$ ,  $P_C$  and  $K_C$ , and we stored in the dataset the obtained values, together with the adopted dimensioning  $h$ . Globally, for each dimensioning, we carried out 10000 different and independent experiments for each value of  $M \in \{1, 2, \dots, M_{max}^h\}$ . Thus, the whole dataset is composed by 21790000 multi-dimensional data points. The adopted simulation campaign allows us to create a dataset representative of the population, because it is obtained considering 10000 independent experiments for each

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

possible value of  $M$ , and the data are of quality, since it does not contain areas of indeterminateness.

The obtained dataset was divided into the training set  $\mathbf{S}_{\text{TR}}$  containing the 90% of the generated points, and the test set  $\mathbf{S}_{\text{TE}}$ , containing the 10% of the points. Moreover, we denote with  $\mathbf{S}_{\text{TE}}^h \subset \mathbf{S}_{\text{TE}}$  the subset containing the points related to group  $h$ , and with  $S_{\text{TE}}^h$  its cardinality. This dataset and the python code to train and test our proposed DNN are available at the following link: <https://figshare.com/s/01107999ec85158d966e>.

#### 3.5.5 Comparison and assessment of estimation methods in a stand-alone RA cycle

The analysis developed in this section considers a stand-alone RA cycle. We first evaluate the performance of the proposed DNN in terms of accuracy in regression for the estimation of the output  $K_C$ . Then, we assess the accuracy of the proposed DNN-based method in estimating  $M$ .

Since it is important to evaluate the performance for each  $T_{pr}$  value, we introduce the following performance metrics related to each of the  $h$ th test set group.

- The Root Mean Square Error (RMSE) between the estimated  $\tilde{y}$  values and the ground truth  $y$  values, for the set  $\mathbf{S}_{\text{TE}}^h$ .

$$RMSE_h = \sqrt{\frac{1}{S_{\text{TE}}^h} \sum_{i=1}^{S_{\text{TE}}^h} (\tilde{y}_i^h - y_i^h)^2}, \quad (3.81)$$

where  $\tilde{y}_i^h$  and  $y_i^h$  are the  $i$ th estimated point and the  $i$ th ground truth of the set  $\mathbf{S}_{\text{TE}}^h$ , respectively. The lower  $RMSE_h$ , the better the accuracy of the estimation.

- The Coefficient of Determination, denoted as  $R_h^2$ , and defined as

$$R_h^2 = 1 - \frac{SS_{res}}{SS_{tot}}, \quad (3.82)$$

where  $SS_{res}$  is the residual sum of squares:

$$SS_{res} = \sum_{i=1}^{S_{\text{TE}}^h} (\tilde{y}_i^h - y_i^h)^2, \quad (3.83)$$

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

$SS_{tot}$  is the total sum of squares:

$$SS_{tot} = \sum_{i=1}^{S_{TE}^h} (y_i^h - \bar{y}^h)^2, \quad (3.84)$$

and  $\bar{y}^h$  is the mean of the observed data:

$$\bar{y}^h = \frac{1}{S_{TE}^h} \sum_{i=1}^{S_{TE}^h} y_i^h. \quad (3.85)$$

The better the regression fits the data in comparison to the horizontal straight line  $\tilde{y}^h = \bar{y}^h$  (the null hypothesis), the closer the value of  $R_h^2$  is to 1. We note that  $R_h^2$  can be negative when the chosen model does not follow the trend of the data and fits worse than the horizontal line  $\tilde{y}^h = \bar{y}^h$ .

**Table 3.4:** *RMSE and  $R^2$  values achieved in the estimation of  $K_C$  for the considered methods*

	DNN <sub>K</sub> -based		Proposed in [53]		Proposed in [56]		Empirical method	
	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$
$T_{pr} = 1$	0.2180	0.7941	0.9302	-1.8897	0.2369	0.7548	0.2286	0.7770
$T_{pr} = 2$	0.1803	0.9202	1.2512	-5.1714	/	/	0.2017	0.8779
$T_{pr} = 3$	0.1651	0.9507	1.4521	-4.6340	/	/	0.2589	0.8402
$T_{pr} = 4$	0.1556	0.9644	1.5999	-4.3127	/	/	0.3328	0.7882

**Table 3.5:** *RMSE values achieved in the estimation of  $M$  for the considered methods*

	DNN <sub>K</sub> -based	Proposed in [53]	Proposed in [56]	Empirical method
$T_{pr} = 1$	5.80	32.64	6.47	6.23
$T_{pr} = 2$	12.60	104.41	/	14.60
$T_{pr} = 3$	19.30	193.70	/	34.15
$T_{pr} = 4$	25.81	294.41	/	62.93

We evaluate the accuracy of the neural network in estimating  $K_C$  and compare it with other proposals available in literature. We consider the estimation methods described in subsection 3.5.1: those proposed in [53, 56] and the empirical formula (3.70) applied in [68]. In effect, these methods estimate  $M$ , so, for each value of  $\tilde{M}$ , we derived the relative estimated  $\tilde{K}_C$  value. Thus, for each estimation method we calculated the

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

$\tilde{M}_i^h$  value related to the  $i$ th test set point of  $\mathbf{S}_{\text{TE}}^h$ . Then, we derived  $\tilde{K}_{C_i}^h$  as

$$\tilde{K}_{C_i}^h = \begin{cases} \frac{\tilde{M}_i^h - P_{S_i}^h}{P_{C_i}^h} & \text{if } P_{C_i}^h > 0 \\ 2 & \text{otherwise.} \end{cases} \quad (3.86)$$

So, we calculated (3.81) and (3.82) with  $\tilde{y}_i^h = \tilde{K}_{C_i}^h$ , and  $y_i^h = K_{C_i}^h$  for the proposed DNN and the benchmark methods. The software environment for all simulations include Python 3.7.7, the PyTorch 1.7.1 tensor library for deep learning, and the NVIDIA CUDA Toolkit version 11.0 for GPU-accelerated training and testing. The trained model and the related weights have been saved in the Open Neural Network Exchange (ONNX) standard format.

The results of the comparison are reported in Table 3.4. As regards the metric values obtained by using the iterative method proposed in [56], we report only the results achieved applying (3.69) with  $T_{pr} = 1$ , i.e.,  $L = 54$ . We did that since the time and space computational complexity, with  $T_{pr} > 1$ , becomes so large that it resulted infeasible to be computed. In fact, already with  $T_{pr} = 2$  the number of collided preamble can reach the value  $P_C = L = 108$ , which considerably increases the complexity of (3.66). As shown in Table 3.4, the proposed DNN-based estimation method outperforms the other ones, both in terms of  $RMSE$  and  $R^2$ , and the higher  $T_{pr}$ , the higher the advantage. Nevertheless, with  $T_{pr} = 1$ , the  $RMSE$  and the  $R^2$  values both for the proposal in [56] and for the empirical method are comparable with those obtained with our method. Conversely, the performance of the estimation method proposed in [53] is very poor; in this regard, it is right to remember that the assessment is carried out considering a stand-alone RA cycle.

Once the accuracy of the proposed DNN has been proved, let us evaluate the performance of the DNN-based method in estimating  $M$ , applying the output of the DNN to formula (3.72). At this aim, in Table 3.5 we report the related RMSE values obtained by the different methods in the reconstruction of  $M$ . Therefore, we calculated (3.81) with  $\tilde{y}_i^h = \tilde{M}_i^h$  related to the considered estimation methods, and  $y_i^h = M_i^h$ , derived as  $P_{S_i}^h + K_{C_i}^h P_{C_i}^h$ . As shown, the values reported in Table 3.5 conform to those of Table 3.4.

It is also useful to evaluate if and how the accuracy of the estimation

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

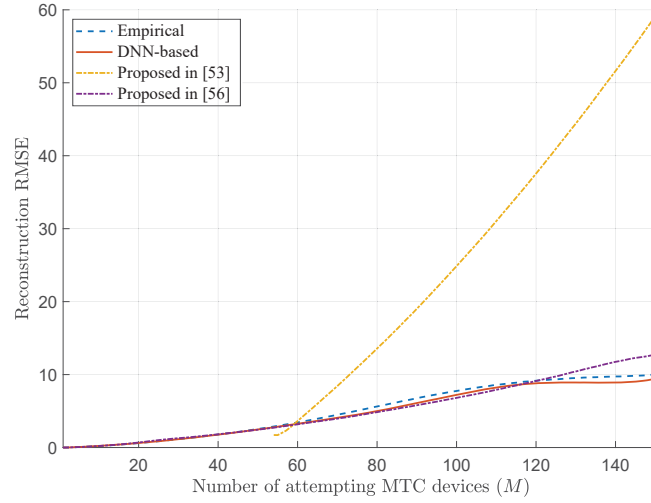
methods depends on the  $M$  value. For this reason, given a group  $h$ , for each value of  $M \in \{1, \dots, M_{max}^h\}$ , we considered the subset  $\mathbf{S}_{\text{TE},M}^h \subset \mathbf{S}_{\text{TE}}^h$ , so that  $y_i^h = M$ . Then, we calculated the related RMSE values by using (3.81) for each point belonging to  $\mathbf{S}_{\text{TE},M}^h$ .

In Fig. 3.30 we show the results obtained for the proposed and the benchmark schemes, when  $h = 1$ , i.e., the  $T_{pr} = 1$  dimensioning is adopted. As shown, the estimation method of [53] is effective only in the working region of an optimal ACB scheme, i.e., when  $M \approx L = 54$ . Otherwise, for  $M < L$  no ACB is performed and, consequently, the estimate cannot be carried out. For  $M > L$ , the higher  $M$  with respect to  $L$ , the lower the accuracy of the estimate. When  $M \gg L$ , the accuracy is very poor. As regards the comparison between the other estimation methods, it is noted that for all the methods, the accuracy degrades as  $M$  increases, but for high values of  $M$  the degradation is higher for the iterative method [56]. To better analyze the difference between these three methods, we report in Fig. 3.31 the related density scatter plots for the  $T_{pr} = 1$  group of the test set. It consists of about 151000 points, the black segment bisector represents the ground truth, the yellow color indicates higher density of points, while blue color lower density. For the DNN-based method and the one in [56], the high density points are symmetrically distributed around the ground truth, but it is confirmed that at high values of  $M$ , the method [56] is less accurate, in fact there are no high density points. As regards the method in [68], it tends to overestimate  $M$ , in fact the points are mainly concentrated above the ideal values.

Furthermore, we show in Fig.3.32, the performance behavior with any other  $T_{pr}$  dimensioning. We compare the DNN-based method with only the empirical method, since the performance achieved by [53] are significantly worse. As shown, the trend is similar for the three dimensionings. In particular, for small-medium values of  $M$ , the method based on DNN achieves a slightly better precision than the empirical method, while the inverse is valid for a small interval of  $M$  (medium-high values). However, for high values of  $M$ , the accuracy of the proposed method is considerably better, especially for the dimensioning  $T_{pr} = 4$ .

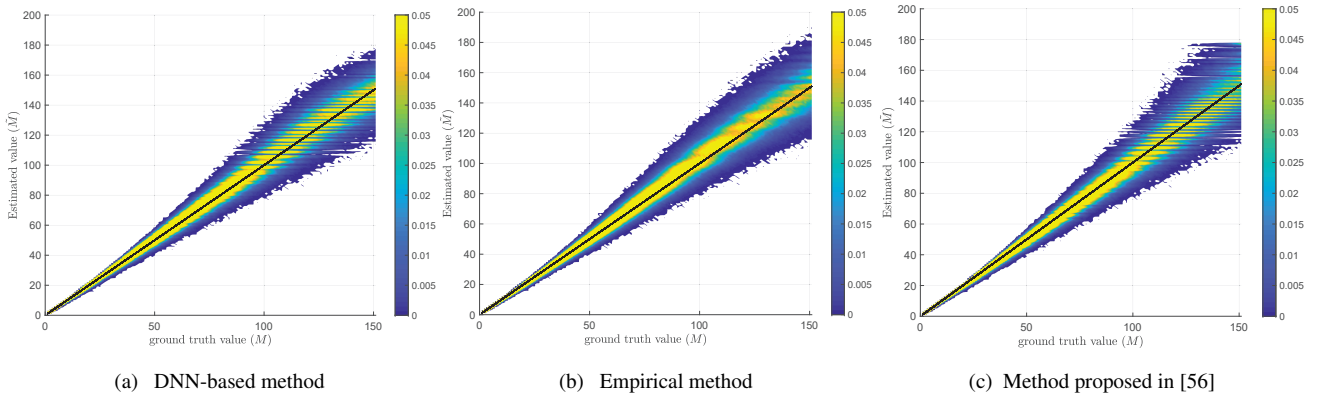
Finally, as done for the  $T_{pr} = 1$  dimensioning, we report in Fig. 3.33 the density scatter plots for the 4th group of the test set. It con-

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.30:** RMSE in the estimation of  $M$  vs  $M$  values with the  $T_{pr} = 1$  dimensioning

sists of about 966000 points, and the density depth is represented by means of the same color bar. As shown, all the points of the DNN-based method are symmetrically distributed around the ground truth (see Fig. 3.33a), while the empirical method presents lightly overestimated values for medium traffic load and strongly underestimated values for high and very high traffic load (see Fig. 3.33b).



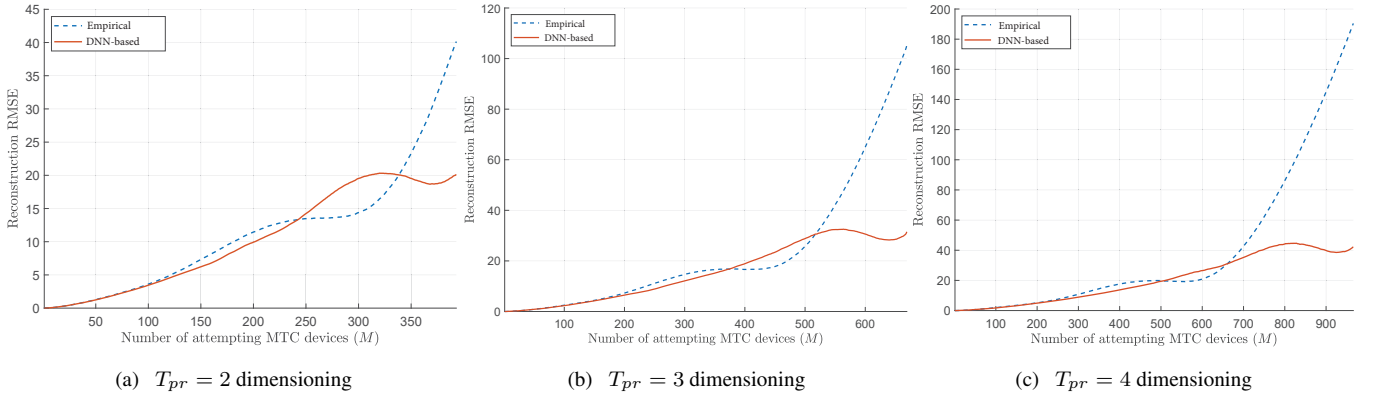
**Figure 3.31:** Density scatter plots in the estimation of  $M$  vs  $M$  values with the  $T_{pr} = 1$  dimensioning.

### 3.5.6 Comparison and assessment of the estimation methods in a long-term analysis

In the static analysis carried out so far, we assessed and compared the performance in estimating the traffic load given a single RA cycle.



### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond



**Figure 3.32:** RMSE in the estimation of  $M$  vs  $M$  values

Now, we want to evaluate the performance achieved by the proposed estimation method on a long-term analysis with time-varying arrival rates. We note that each  $M$  value includes both the new access attempts and the access reattempts, and the latter ones depend on the access control strategy adopted. Since in the next B5G networks it is expected a very high density of connection requests, it will be necessary to apply innovative load-aware access controls and/or a dynamic uplink resource dimensioning at the aim of minimizing congestion.

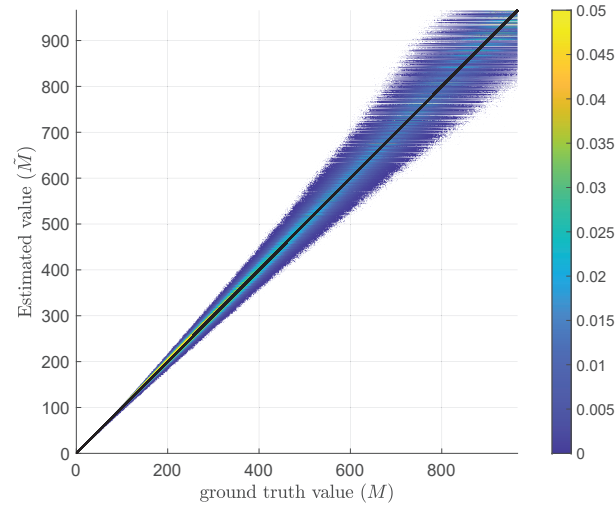
In this regard, in the following long-term analysis, we adopt the Dynamic Uplink Resource Dimensioning (DURD) scheme because it has shown good performance even in high traffic load condition [68]. The DURD control scheme will work in three different conditions: assuming an exact knowledge of the traffic load  $M$ ,  $DURD_I$ ; using an estimate derived from the DNN-based method,  $DURD_{DNN}$ ; or, applying an estimate derived using the empirical method based on formula (3.70),  $DURD_E$ .

#### DURD overview

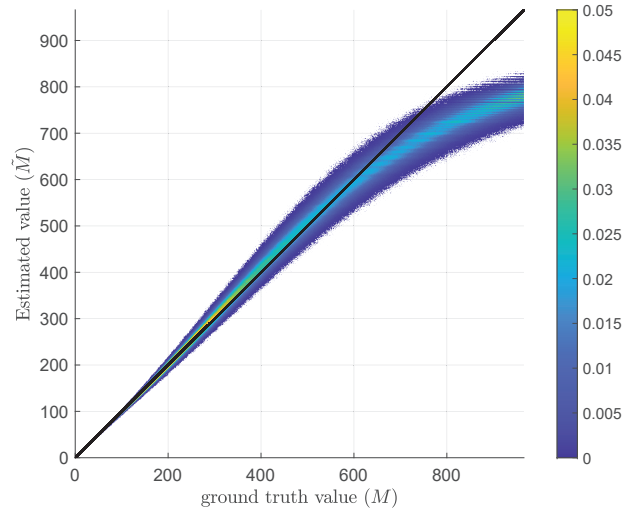
In this subsection, we briefly describe the DURD scheme proposed in [68], which aims to appropriately split the uplink radio resources between PRACH and PUSCH based on the current value of  $M$ .

Firstly, we observe that a communication is succeeded if both the following conditions occur: the MTC device transmits with success its preamble sequence in the PRACH; and there are resources available in the PUSCH for its data transmission.

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



(a) *DNN-based method*



(b) *Empirical method*

**Figure 3.33:** Density scatter plots in the estimation of  $M$  vs  $\bar{M}$  values with the  $T_{pr} = 4$  dimensioning.

So, given a RA cycle the average number of succeeded communications is:

$$\bar{C}_S = \min(\bar{P}_S, DT_{max}), \quad (3.87)$$

where  $DT_{max}$  is the total number of data transmissions which can be satisfied in the PUSCH. We remind that the data transmissions in the PUSCH are allocated by means of the SCMA technique, and the related

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

total number can be derived as:

$$DT_{max} = \left\lfloor \frac{\lfloor 72/Q \cdot 14 \rfloor (QK_{max})/S}{\lceil \theta_{max}/\log_2(I) \rceil} \right\rfloor (T_{ra} - T_{pr}), \quad (3.88)$$

where  $Q$  is the number of REs in one  $SCMA_{block}$ ,  $S$  is the number of REs which a layer occupies respect to  $Q$ ,  $K_{max}$  is the maximum number of overlapped layers with different codebooks in one RE, and  $\log_2(I)$  is the number of information bits sent for each symbol of the constellation.

We analytically derived the average number of successful communications ( $\bar{C}_S$ ) with respect to  $M$ , and plot it in Fig. 3.34 for different values of  $T_{pr}$  value. As shown, there is no a single optimal  $T_{pr}$  value, since it depends strongly on the offered load, i.e., on the value of  $M$ . For this reason, we proposed the DURD algorithm for dynamically dimensioning the uplink radio resources on the basis of the current value of  $M$ . The DURD scheme works as follows.

During RA cycle  $j$ , the following information are available at the gNB:  $L^j = L_0 T_{pr}^j$ , i.e., the number of available preambles in the PRACH of RA cycle  $j$ ;  $P_S^j$ , i.e., the number of succeeded preamble in RA cycle  $j$ ; and  $P_C^j$ , the number of collided preambles in RA cycle  $j$ . On the basis of these information, the value of  $M^j$  can be estimated by means of either the empirical estimation method or the DNN-based one proposed here. In [16] we showed that the process of total access attempts in two consecutive RA cycles was correlated. Hence, we assumed that

$$M^{j+1} \cong \tilde{M}^j. \quad (3.89)$$

So, the DURD algorithm derives the optimal  $T_{pr}$  dimensioning for the next RA cycle  $j + 1$  as:

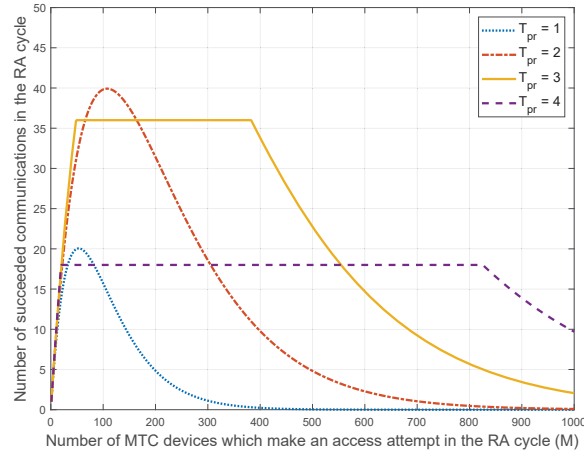
$$T_{pr}^{j+1} = \arg \max_{T_{pr} \in \{1, \dots, T_{ra}-1\}} \left\{ \bar{C}_S(\tilde{M}^j) \right\}. \quad (3.90)$$

Specifically, as shown in Fig. 3.34, the optimal dimensioning is  $T_{pr} = 3$  from a light load ( $M = 1$ ) to medium-high load ( $M = 554$ ), except for the interval  $M \in [66, \dots, 163]$ , where the optimal dimensioning is  $T_{pr} = 2$ . Conversely, in the presence of high load ( $M > 555$ ), the optimal  $T_{pr}$  dimensioning becomes  $T_{pr} = 4$ . It is worth pointing out that for each optimal  $T_{pr}$  dimensioning, the corresponding  $M$  values conform

## Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

**Table 3.6:** *Simulation Parameters*

Symbol	Value
$L_0$	54
$B$	1.08 MHz
$M_A$	10
$\theta_{max}$	160 bits
$T_{ra}$	5 time slots
$T_{sim}$	10s (2000 RA cycles)
$Q$	4
$K_{max}$	3
$S$	2
$B_W$	20ms
$I$	4



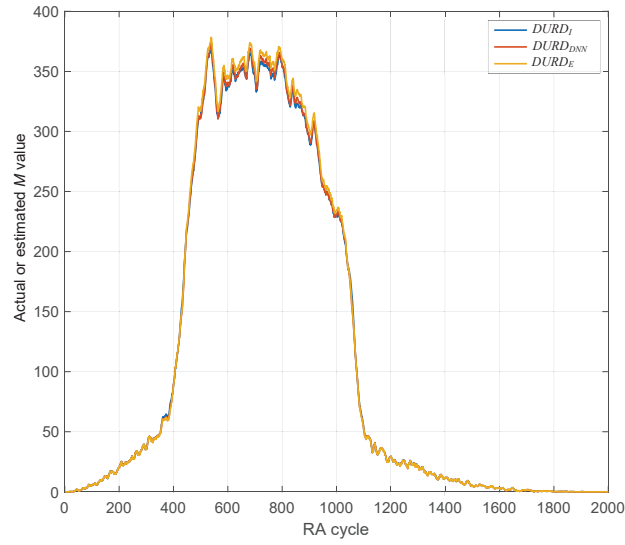
**Figure 3.34:** *Expected number of successful communications vs the number of MTC devices ( $M$ ) which attempts access in a given RA cycle, assuming different PRACH dimensionings.*

to the range of  $M$  values for which the DNN has been trained and tested. For example, the optimal dimensioning  $T_{pr} = 2$  is related to the interval  $M \in [66, \dots, 163]$  contained in the subset  $[1, \dots, M_{max}^2 = 393]$ .

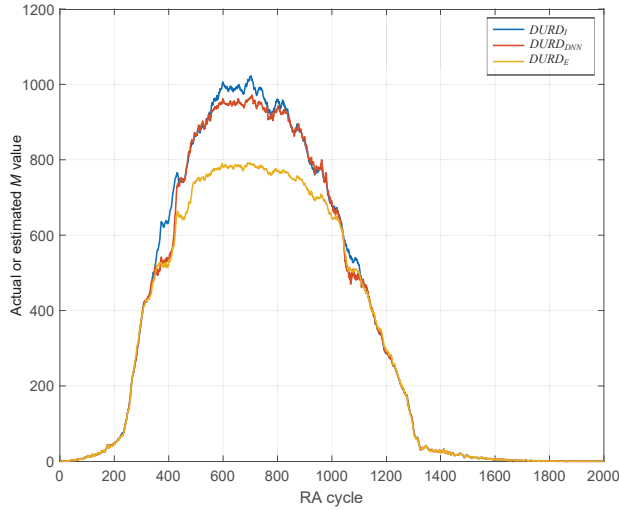
### Simulation setup

We consider two different arrival distributions for the new MTC access attempts. The first one is created through simulations by following the 3GPP guidelines [67]. The second one, at the aim of evaluating how well the DNN-based approach performs by using real traffic data traces, adopts an arrival distribution derived from the dataset presented in [104],

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond



(a)  $N_{MTC} = 50000$



(b)  $N_{MTC} = 100000$

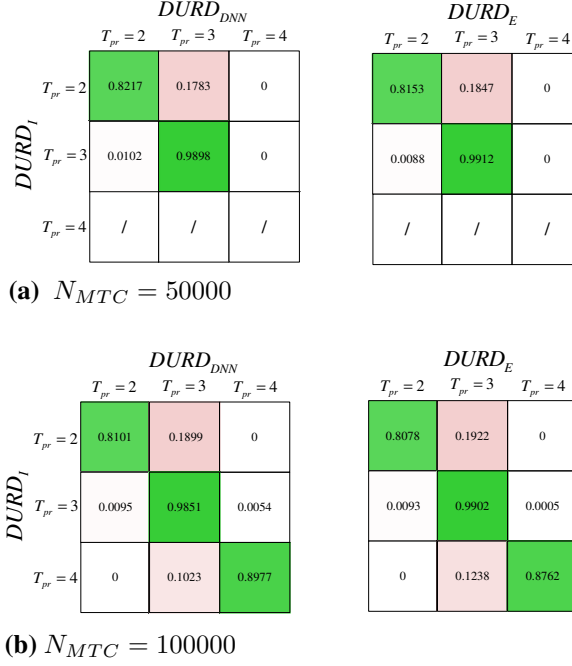
**Figure 3.35:** The temporal distributions of actual or estimated  $M$  values with different control schemes for the simulated Beta arrival distribution in the space of 10 s.

which is collected from a real IoT setup.

As regards the first arrival distribution, we adopt the Beta distribution, which is related to a scenario where a large number of MTC devices access the network in a highly synchronized manner. Each MTC device generates a single data at time  $t \in [0, T_{arrival}]$  following the Probability Density Function (PDF):

$$f(t) = \frac{t^{\alpha-1} (T_{arrival} - t)^{\beta-1}}{T_{arrival}^{\alpha+\beta-1} \text{Beta}(\alpha, \beta)}, t \in [0, T_{arrival}], \quad (3.91)$$

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.36:** Confusion Matrices with the simulated Beta arrival distribution.

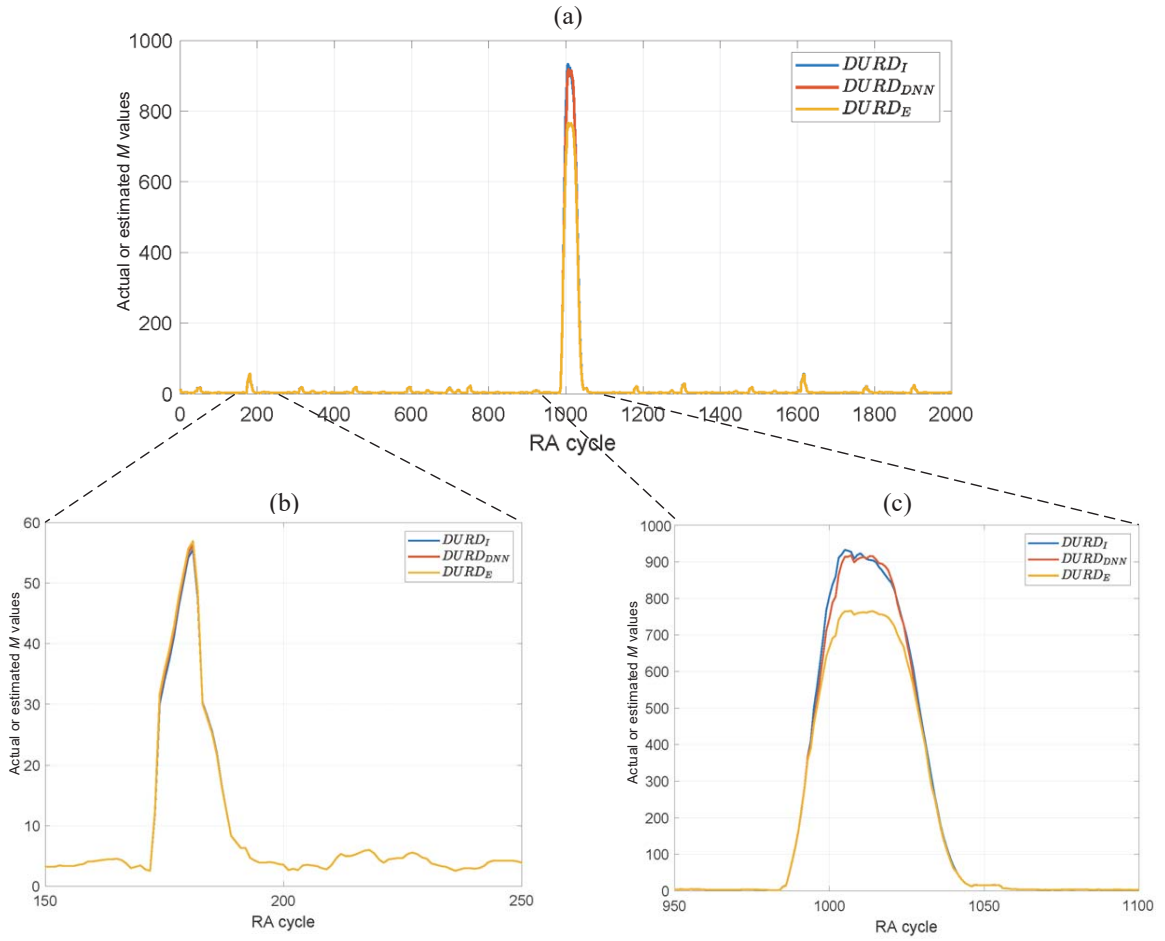
where

$$Beta(\alpha, \beta) = \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt. \quad (3.92)$$

Specifically, given  $N_{MTC}$  devices, we generated a realization of the random process by extracting  $N_{MTC}$  numbers from the PDF (3.91), with  $\alpha = 3$ ,  $\beta = 4$ , and  $T_{arrival} = 10s$  [67]. Then, we made the Beta arrival vector ( $\mathbf{A}_{Beta}$ ) containing the overall number of new arrivals for each RA cycle  $j$ , with  $j = 1, 2, \dots, \left\lfloor \frac{T_{arrival}}{T_{ra}} \right\rfloor$ , i.e., each element  $\mathbf{A}_{Beta}[j]$  is the sum of the new access attempts in the interval time  $[(j-1)T_{ra}, jT_{ra}]$ .

As concerns the second arrival distribution, we exploited the trace data openly available online in [104] containing daily logs from 28 unique IoT devices (e.g., motion and smoke sensors) for a total of two weeks. Among all the open datasets observed, it resulted one of the best data collection regarding the IoT traffic. For each day, the log contains, *inter alia*, the 'TIME' column representing the time points at which the access requests were captured with a recording frequency of 1 Hz. We made the real arrival vector ( $\mathbf{A}_{Real}$ ) whose element  $j$  is the number of the access requests collected in the  $j$ th second of the day, starting from midnight. However, this real dataset is related to a legacy scenario, where

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond



**Figure 3.37:** The temporal distributions of actual or estimated  $M$  values with different control schemes for the real arrival distribution in the space of 10 s.

		$DURD_{DNN}$			$DURD_E$		
		$T_{pr}=2$	$T_{pr}=3$	$T_{pr}=4$	$T_{pr}=2$	$T_{pr}=3$	$T_{pr}=4$
$DURD_I$	$T_{pr}=2$	0.9412	0.0588	0	0.8824	0.1176	0
	$T_{pr}=3$	0	0.9995	0.0005	0	0.9995	0.0005
	$T_{pr}=4$	0	0.0312	0.9688	0	0.0312	0.9688

**Figure 3.38:** Confusion Matrices for the real arrival distribution.

the maximum allowed connection density is  $10^5$  devices per  $Km^2$  [2]. So, in order to assess the estimation performance of the proposed DNN-based method in a futuristic 6G scenario, characterized by a maximum requirement of  $10^7$  devices per  $Km^2$ , we assume that the access requests contained in  $\mathbf{A}_{Real}$  are performed with a periodicity of one RA cycle in-

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

stead of one second. We underline that this assumption does not affect the arrival pattern, which remains in line with the real IoT setup, and introduces the burstiness of arrival requests, typical of the mIoT scenario. Out of all the available logs, we show the performance of an extract lasting 10 seconds, as same as the simulated distribution time<sup>6</sup>.

For both the considered arrival distributions, we made a simulation campaign in MATLAB environment. Moreover, for estimating  $M$  with the DNN-based method, the pre-trained DNN and its weights have been imported by means of the `importONNXNetwork` function of the Deep Learning Toolbox. The other simulation parameters are reported in Table 3.6.

#### Performance evaluation

In this subsection, we compare the DNN-based and the empirical estimation method based on formula (3.70) in the considered simulation setup. We underline that we do not perform a comparison neither with the estimation method [56], since this estimate supports only the  $T_{pr} = 1$  dimensioning, nor with [53] since it supports only ACB-based schemes.

We compare the performance of the above methods by evaluating in the  $j$ th RA cycle the  $M$  value, consisting of the number of new access attempts (in either  $\mathbf{A}_{\text{Beta}}[j]$  or  $\mathbf{A}_{\text{Real}}[j]$ ) and the access reattempts scheduled in the RA cycle considered.

Let us start considering the simulated arrival distribution. In Fig. 3.35a and 3.35b we show the temporal variation of  $M$ , with  $N_{MTC} = 50000$  and  $N_{MTC} = 100000$ , respectively. Each curve represents the average value over 10 simulations for the control system considered. In particular, the blue line is relative to the  $DURD_I$  system, so, for each RA cycle  $j$ , the perfectly known  $M^j$  value is used to provide the optimal dimensioning for the next RA cycle  $j + 1$  by using (3.90). The red line is relative to the  $DURD_{DNN}$  system, where for each RA cycle  $j$ , the number of access attempts is estimated by means of the DNN-based estimation method, and the next optimal dimensioning is calculated on the basis of this estimate. Similarly, the yellow line is related to the  $DURD_E$  system. Moreover, at the aim of evaluating the impact of the estimate methods for the specific control scheme in the uplink

---

<sup>6</sup>We adopt the  $\mathbf{A}_{\text{Real}}$  vector related to September 24th, 2016, with  $j \in \{16901, \dots, 18900\}$ .



### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

dimensionings, we present in Fig. 3.36, two Confusion Matrices, when  $N_{MTC} = 50000$  and  $N_{MTC} = 100000$ , respectively. Each row of the matrix represents the instances of the ideal dimensioning (i.e., the optimal dimensioning obtained with exact knowledge of  $M$ ), while each column represents the instances provided by the estimation methods. The diagonal elements represent the percentage of times for which the estimated dimensioning is equal to the ideal one, while off-diagonal elements are those that are wrongly estimated. Clearly, the higher the diagonal values of the confusion matrix the better the predicted  $T_{pr}$  dimensioning. Let us start analyzing the case  $N_{MTC} = 50000$ . As shown in Fig. 3.36a, the performances achieved by both systems are practically comparable. Nevertheless, we want to carry out a thorough analysis and to show the effect produced by a wrong dimensioning on the subsequent estimation of  $M$ . We focus on the most relevant error reported in Fig. 3.36a: when the ideal dimensioning is  $T_{pr} = 2$  (i.e.,  $M \in \{66, \dots, 163\}$ , see Fig. 3.34) but it has been wrongly estimated as  $T_{pr} = 3$  (17.83% and 18.47% error percentage for  $DURD_{DNN}$  and  $DURD_E$ , respectively). However, in the involved range of  $M$  values, the estimate of the traffic load  $\tilde{M}$  is still very accurate for both the two dimensionings (see Figs. 3.32a and 3.32b). So, the effects of these errors are practically negligible. In fact, in Fig. 3.35a, in the transition zones between the  $T_{pr} = 2$  and  $T_{pr} = 3$  dimensionings, the estimated  $\tilde{M}$  values are very accurate. In addition, during the whole simulation time, the number of attempting devices does not reach high values, in fact the peak value of  $M$  is about 368. As shown in Fig. 3.32b, the accuracies achieved by both the estimation methods adopting the  $T_{pr} = 3$  dimensioning with this low and medium traffic load are comparable. For this reason, both the estimation methods show in Fig. 3.35a the same good performance during the whole simulation time. As regards the case  $N_{MTC} = 100000$  (see Fig. 3.36b), the performances exhibited by the two schemes in the dimensioning are again comparable and substantially the same as before, but now the impact on the estimate of  $M$  is evidently different, as shown in Fig. 3.35b. Unlike the previous analysis, there is also the case in which the ideal dimensioning  $T_{pr} = 4$  is wrongly estimated as  $T_{pr} = 3$ . This wrong uplink dimensioning impacts the estimate of  $M$  for both the estimation methods and it is due to the assumption (3.89) taken for the DURD con-

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

trol system. In fact, since the arrival process is not always so strongly correlated, the DURD algorithm, based on an accurate estimate of  $M^j$ , could decide the dimensioning  $T_{pr} = 3$ , but then the actual traffic load  $M^{j+1}$  can be much greater than 555. This involves that the successive estimates are poor, because when  $T_{pr} = 3$  for  $M > 555$  the formula (3.70) is very inaccurate, as shown in Fig 3.33b, and the DNN-based method works for values of  $M > M_{max}^3$ . In Fig. 3.35b, this effect is noticeable from RA cycle 345 to 420; less, from 1050 to 1100. Furthermore, the difference of error in the estimate of  $M$  is very evident in the central range of Fig. 3.35b, from RA cycle 420 to 1050, where the optimal  $T_{pr}$  dimensioning is  $T_{pr} = 4$  for both the estimation methods. Under very high load condition, high values of  $M$ , the  $DURD_{DNN}$  is centered around the ground truth, while the  $DURD_E$  provides a serious underestimate, as shown in Figs. 3.33a and 3.33b, respectively. This implies, in the analyzed range, a very marked improvement in the estimate of the traffic load for the DNN-based method compared to the competitor one, as shown in Fig. 3.35b. We also underline that, in the interval where the blue curve reaches the peak value (around 1000), the estimate performed by the DNN is slightly less accurate since the  $M$  values to be estimated are greater than  $M_{max}^4$ .

Finally, let us focus on the real arrival distribution. Similarly to the simulated one, we show the temporal variation of  $M$  in Fig. 3.37a. As shown, the PRACH load  $M$  does not follow a Beta distribution throughout the entire simulation length (i.e., 10s), but it is quite low and about constant for almost the entire simulation time, and increases and decreases rapidly in only about 70 RA cycles (from RA cycle 980 to 1050). In the light PRACH load region, the performance of the  $DURD_{DNN}$  and  $DURD_E$  are very similar, as shown in the zoom reported in Fig. 3.37b. This result is in line with Fig.3.35a. In stead, in the central zone, the difference in the estimate of  $M$  by adopting the DNN-based method compared to the estimate one is very evident (see Fig. 3.37c), and this trend is similar to Fig.3.35b. As regards the impacts of the accuracy in the estimate on the functioning of the DURD control scheme, we show also the related Confusion Matrices in Fig. 3.38. As shown, the results obtained with the two control schemes,  $DURD_E$  and  $DURD_{DNN}$ , are very similar. This is due to the fact that, as shown in Fig. 3.34, for a

### 3.5. A DNN-based estimate of the PRACH traffic load for massive mMTC scenarios in 5G networks and beyond

wide range of very high traffic load values (i.e.,  $M > 555$ ) the optimal dimensioning adopted by the DURD scheme is the same and equal to  $T_{pr} = 4$ . Therefore, the considered resource dimensioning algorithm is not affected too much by the inaccuracy of the estimate, but this is not always valid for any other control scheme.

In conclusion, the performance achieved by adopting a real daily track proved to be in line with the one derived by adopting synthetic data. In both cases, we have demonstrated the best accuracy of the proposed DNN-based estimation method, especially under high traffic load condition, that is expected in the future B5G and 6G networks.

#### 3.5.7 Appendix

#### 3.5.8 Calculation of $\bar{P}_S$ , $\bar{P}_C$ and the pmf of $P_C$

Let  $|\mathcal{M}_p|$  be the cardinality of preamble  $p$ , that is, the number of devices that have chosen preamble index  $p$ . Firstly, we derive the probability that the preamble  $p$  has been selected by exactly  $h$  devices, with  $h = \{1, \dots, M\}$ . Since each device performs an independent extraction of one preamble out of  $L$  available ones, this probability follows the binomial distribution:

$$\Pr\{|\mathcal{M}_p| = h\} = \binom{M}{h} \left(\frac{1}{L}\right)^h \left(1 - \frac{1}{L}\right)^{M-h}. \quad (3.93)$$

The probability of success corresponds to (3.93) with  $h = 1$ . Let us define the following random variables  $X_p$ , with  $p \in \{1, \dots, L\}$ , as:

$$X_p = \begin{cases} 1 & \text{if } |\mathcal{M}_p| = 1 \\ 0 & \text{otherwise} \end{cases} \quad (3.94)$$

Therefore, the number of succeeded preambles results:

$$P_S = \sum_{p=1}^L X_p. \quad (3.95)$$

Since  $X_p$  are independent and identically distributed (i.i.d.) random variables, the mean value of  $P_S$  is:

$$\begin{aligned}\bar{P}_S &= \sum_{p=1}^L E\{X_p\} = L \Pr\{X_p = 1\} = \\ &L \Pr\{|\mathcal{M}_p| = 1\} = M \left(1 - \frac{1}{L}\right)^{M-1}.\end{aligned}\quad (3.96)$$

At the aim of deriving the average number of collided preambles,  $\bar{P}_C$ , we introduce the probability that the preamble  $p$  has been selected at least by  $i$  MTC devices, with  $i = \{1, \dots, M\}$ . It results:

$$\Pr\{|\mathcal{M}_p| \geq i\} = 1 - \sum_{j=0}^{i-1} \binom{M}{j} \left(\frac{1}{L}\right)^j \left(1 - \frac{1}{L}\right)^{M-j}.\quad (3.97)$$

Clearly, the probability of collision corresponds to (3.97) with  $i = 2$ . Let us define the following random variables  $Y_p$ , with  $p \in \{1, \dots, L\}$ , as:

$$Y_p = \begin{cases} 1 & \text{if } |\mathcal{M}_p| \geq 2 \\ 0 & \text{otherwise} \end{cases}\quad (3.98)$$

Therefore:

$$P_C = \sum_{p=1}^L Y_p.\quad (3.99)$$

$$\begin{aligned}\bar{P}_C &= L \Pr\{Y_p = 1\} = L \Pr\{|\mathcal{M}_p| \geq 2\} = \\ &L \left[1 - \left(1 - \frac{1}{L}\right)^M - M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1}\right].\end{aligned}\quad (3.100)$$

Finally, we derive the pmf for the random variable  $P_C$ . It is defined as:

$$p_{P_C}[k] = \Pr\{P_C = k\} = \Pr\left(\sum_{p=1}^L Y_p = k\right),\quad (3.101)$$

where  $k = 0, 1, \dots, L$ . Given  $k$ , the probability  $\Pr\{P_C = k\}$  follow the binomial distribution:

$$p_{P_C}[k] = \binom{L}{k} [\Pr\{Y_p = 1\}]^k [\Pr\{Y_p = 0\}]^{L-k},\quad (3.102)$$

### 3.6. A Wasserstein GAN autoencoder for SCMA networks

where  $\Pr\{Y_p = 0\} = 1 - \Pr\{|\mathcal{M}_p| \geq 2\}$ . It results:

$$p_{PC}[k] = \binom{L}{k} \left[ 1 - \left(1 - \frac{1}{L}\right)^M - M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1} \right]^k \left[ 1 - \left(\frac{1}{L}\right)^M + M \left(\frac{1}{L}\right) \left(1 - \frac{1}{L}\right)^{M-1} \right]^{L-k}. \quad (3.103)$$

## 3.6 A Wasserstein GAN autoencoder for SCMA networks

### 3.6.1 System model and Problem Formulation

We consider  $J$  SCMA layers, each one transmitting data symbols out of  $M$ , over  $K$  Resource Elements (REs), where  $K < J$ . Each layer  $j \in \{1, \dots, J\}$  adopts a unique codebook  $\mathbf{c}_j$ , composed of  $M$   $K$ -dimensional codewords (i.e.,  $\mathbf{c}_j = [\mathbf{c}_{1j}, \dots, \mathbf{c}_{Mj}]$ , with  $\mathbf{c}_{ij} \in \mathbb{C}^{K \times 1}$ ). Let  $\mathbf{b}_j$  be the vector containing all the possible input data symbols for the encoder of layer  $j$ , i.e.,  $\mathbf{b}_j = [\mathbf{b}_{1j}, \dots, \mathbf{b}_{Mj}]$  with  $\mathbf{b}_{ij} \in \mathbb{B}^{m \times 1}$ , and  $m = \log_2(M)$ . Also, we define  $\mathbf{B} = \mathbf{b}_1 \times \dots \times \mathbf{b}_J$  as the set of all possible  $J$ -tuple  $(\mathbf{b}_{i1}, \dots, \mathbf{b}_{iJ})$ . Each  $\mathbf{b}_{ij}$  is mapped into a unique  $K$ -dimensional codeword  $\mathbf{c}_{ij}$ , with only  $S < K$  non-zero elements. The mapping function, related to layer  $j$ , is bijective and denoted as  $f_j$ , i.e.,  $\mathbf{c}_{ij} = f_j(\mathbf{b}_{ij})$  and  $\mathbf{b}_{ij} = f_j^{-1}(\mathbf{c}_{ij})$ . Let  $\mathbf{c}_{ij}^{[k]}$  denote the  $k$ th dimension of codeword  $\mathbf{c}_{ij}$ . Typically, the  $S$  non-zero entries of each layer are statically configured and represented by a factor graph  $\mathbf{F} \in \mathbb{B}^{K \times J}$ . Each element  $F_{kj} \in \mathbf{F}$  is equal to 1 if and only if  $\mathbf{c}_{ij}^{[k]} \neq 0, \forall i = 1, \dots, M$ . The number of overlapped layers in a single RE is fixed and equal to  $d_f$ , i.e.,

$$\sum_{j=1}^J F_{kj} = d_f, \forall k \in \{1, \dots, K\}. \quad (3.104)$$

The overall encoding procedure and data transmission is described in the following. Each layer selects an input symbol  $\mathbf{x}_j \in \mathbf{b}_j$  and transmits the SCMA codeword  $\mathbf{e}_j = f_j(\mathbf{x}_j)$ . From the encoded signals  $\{\mathbf{e}_j\}_{j=1}^J$  transmitted by all layers, the summed signal  $\mathbf{s} \in \mathbb{C}^{K \times 1}$  can be expressed as

$$\mathbf{s} = \sum_{j=1}^J \text{diag}(\mathbf{h}_j) \mathbf{e}_j = \sum_{j=1}^J \text{diag}(\mathbf{h}_j) f_j(\mathbf{x}_j), \quad (3.105)$$

where  $\mathbf{h}_j = [h_{1j}, \dots, h_{Kj}]^T$  is the channel gain vector for the  $K$  REs of layer  $j$ , and  $\text{diag}(\cdot)$  denotes the diagonalized matrix. In addition to the summed signal  $\mathbf{s}$ , we consider the AWGN  $\mathbf{n} = [n_1, \dots, n_K]^T$ , with  $n_k = \mathcal{CN}(0, \sigma_N^2), \forall k = 1, \dots, K$ . So, the received signal  $\mathbf{r}$  can be written as  $\mathbf{r} = \mathbf{s} + \mathbf{n}$ .

As regards the receiver side, after the reception of the signal  $\mathbf{r}$ , the SCMA decoder aims to reconstruct the original symbols transmitted by the  $J$  layers, i.e.,  $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_J]^T$ , given  $\{\mathbf{h}_j\}_{j=1}^J$ , and the codebook set  $\{\mathbf{c}_j\}_{j=1}^J$ . We denote with  $\mathbf{y}$  the output vector of the decoder. Obviously, the ideal decoder gives as output  $\mathbf{y} = \mathbf{x}$  for any combination of symbols sent by the layers, and any channel condition.

The optimum multi-user detection problem can be solved by finding the MAP pmf of all layers' transmitted symbols [49]. Given  $\mathbf{r}$  and assuming that both  $\{\mathbf{h}_j\}_{j=1}^J$  and  $\sigma_N^2$  are available at the receiver, the decoder will estimate the optimal  $J$ -tuple  $(\mathbf{c}'_{i1}, \dots, \mathbf{c}'_{iJ})$  as follows

$$\begin{aligned} (\mathbf{c}'_{i1}, \dots, \mathbf{c}'_{iJ}) = \arg \max_{(\mathbf{c}_{i1}, \dots, \mathbf{c}_{iJ}) \in \mathbf{C}} \prod_{k=1}^K \frac{1}{\sqrt{2\pi\sigma_N}} \cdot \\ \exp \left( -\frac{1}{2\sigma_N^2} \left\| r_k - [s_k \mid (\mathbf{c}_{i1}, \dots, \mathbf{c}_{iJ}), \{\mathbf{h}_j\}_{j=1}^J] \right\|^2 \right), \end{aligned} \quad (3.106)$$

where  $\mathbf{C} = \mathbf{c}_1 \times \dots \times \mathbf{c}_J$  is the set of all possible combinations of code-words,  $[s_k \mid (\mathbf{c}_{i1}, \dots, \mathbf{c}_{iJ}), \{\mathbf{h}_j\}_{j=1}^J]$  represents the  $k$ th dimension of the summed signal calculated in (3.105) given  $\mathbf{e}_j = \mathbf{c}_{ij}, \forall j = 1, \dots, J$  and the channel gain vectors  $\{\mathbf{h}_j\}_{j=1}^J$ . Then, the output of the decoder is  $\mathbf{y} = [f_1^{-1}(\mathbf{c}'_{i1}), \dots, f_J^{-1}(\mathbf{c}'_{iJ})]^T$ . However, the complexity of this detector is very high and increases exponentially with  $J$  and polynomially with  $M$ .

In order to decrease the complexity of MAP detector, the Log-MPA detector has been introduced as a near optimal iterative detection algorithm. In the same way as the MAP, it assumes that both  $\{\mathbf{h}_j\}_{j=1}^J$  and  $\sigma_N^2$  are available at the receiver. Unlike the MAP solutions (3.106), this

### 3.6. A Wasserstein GAN autoencoder for SCMA networks

algorithm works in the logarithm domain in order to slightly reduce the system complexity. The Log-MPA detector is based on an iterative exchange of information between two types of nodes belonging to a bipartite factor graph. The first one represents the  $K$  REs, and the second type the  $J$  layers. The edges of the graph are constructed in such a way that the  $j$ th node is connected to the node  $k$  if and only if layer  $j$  transmits in RE  $k$ . The overall procedure is described in [51]. However, despite these improvements, the complexity of this detector remains very high.

#### 3.6.2 Our solution

In this section, we first describe our high-level approach to deal with the joint encoding and decoding problem. Then, we provide the implementation of the proposed strategy.

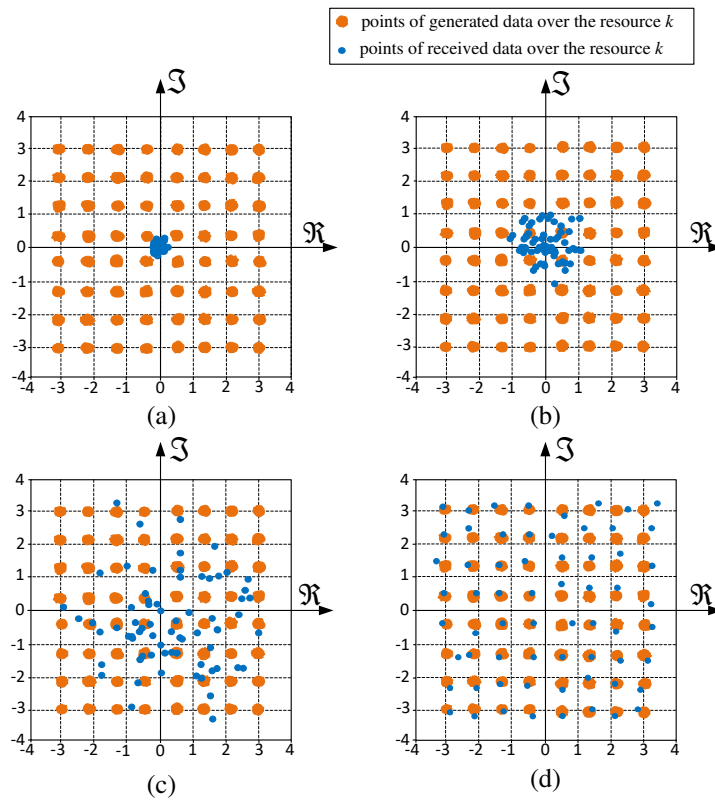
##### The proposed approach

At the receiver side, the task can be considered as a classification problem that identifies to which  $J$ -tuple of  $\mathbf{B}$  the observation  $\mathbf{r}$  belongs to. Clearly, the correct functioning of the decoding activity does not depend only from the strength of the decoder, but also from the goodness of the summed signal  $\mathbf{s}$ , thus, from the codebook set  $\{\mathbf{c}_j\}_{j=1}^J$  adopted. This means that, for each  $J$ -tuple in  $\mathbf{B}$  as input, the received signal  $\mathbf{r} = \mathbf{s} + \mathbf{n}$  should ideally correspond to a unique output  $\mathbf{y}$ , regardless of the noise  $\mathbf{n}$ .

At the aim of reaching this goal, we introduce the following innovative idea. To explain it, we analyze the received signal in a generic dimension  $k$ :

$$r_k = s_k + n_k = \sum_{j=1}^J F_{kj} h_{kj} e_{kj} + n_k. \quad (3.107)$$

Considering (3.104), it follows that the summed signal  $s_k$  is composed by the superposition of  $d_f$  contributions. So, given  $\{\mathbf{h}_j\}_{j=1}^J$  and by adopting a candidate codebook set, the transmission of each possible combination of input data symbols in  $\mathbf{B}$  corresponds to a constellation of  $M^{d_f}$  complex points. On the basis of the codebook set, some of the points may overlap or nearly overlap, invalidating the correct decoding, even in the absence of noise or in the presence of high  $E_b/N_0$ .



**Figure 3.39:** Desired evolution of the WGAN-SCMA training procedure from the beginning (a) to the objective (d), with  $M = 4$  and  $d_f = 3$ . Each sub-figure (a)-(d) represents the complex plane related to dimension  $k$ .



### 3.6. A Wasserstein GAN autoencoder for SCMA networks

Aiming to increase the robustness of  $s_k$  to the noise, we consider one noise-robust  $M^{d_f}$ -ary modulation available in literature (e.g.,  $M^{d_f}$ -QAM, or star  $M^{d_f}$ -QAM, etc.), and we generate  $M^{d_f}$  non-overlapping regions on the complex plane, starting from the signal constellation points. Then, the optimal solution is achieved if each out of the  $M^{d_f}$  points of  $s_k$  belong uniquely to one of the non-overlapping regions, regardless of the channel noise.

Specifically, we generate a prior distribution  $\mathbf{p}^{[k]}$  by adding to the constellation points several realizations of the AWGN with a very small variance, denoted as  $\sigma_p^2$ . We report an example in Fig. 3.39, where we show the complex plane related to the dimension  $k$ , with  $M = 4$  and  $d_f = 3$ . The orange points represent a prior distribution  $\mathbf{p}^{[k]}$  generated from the 64-QAM modulation, and the blue points represent the  $M^{d_f}$  possible values of  $s_k$ , obtained by adopting a given codebook set, plus one realization of the channel noise  $n_k$ , with  $\sigma_N^2$  depending on the considered  $E_b/N_0$  value. As shown in Fig. 3.39a, the blue points overlap and, consequently, the decoder cannot decode properly the received signal. This means that the codebook set needs to be improved, and the optimal solution is achieved when the condition shown in Fig. 3.39d occurs.

To achieve this, we introduce a proper WGAN model with the goal of generating the optimal codebook set so that the distribution of  $\mathbf{r}^{[k]}$ , under different realizations of the channel noise, fits  $\mathbf{p}^{[k]}$ , in each dimension  $k$ . The lower  $E_b/N_0$ , the harder the goal becomes. For low  $E_b/N_0$  values, the task is not fully reachable because, as the level of noise increases, the constellation points are randomly moved away from the desired position. This means that the functioning of the system cannot be perfect, but good performance can be achieved with both a proper encoding technique and an optimal decoder tailored for reconstruct these constellation points at the receiver side.

#### Implementation

We implement the approach reported in the previous subsection by introducing a new model composed of a WGAN and an autoencoder. Since we aim to fit the distribution  $\mathbf{r}^{[k]}$  with  $\mathbf{p}^{[k]}$ ,  $\forall k = 1, \dots, K$ , we introduce  $K$  Critics, each one having a critic function denoted as  $f_{c_k}(\cdot)$ . The critic  $k$  is trained aiming to minimize the following loss function:

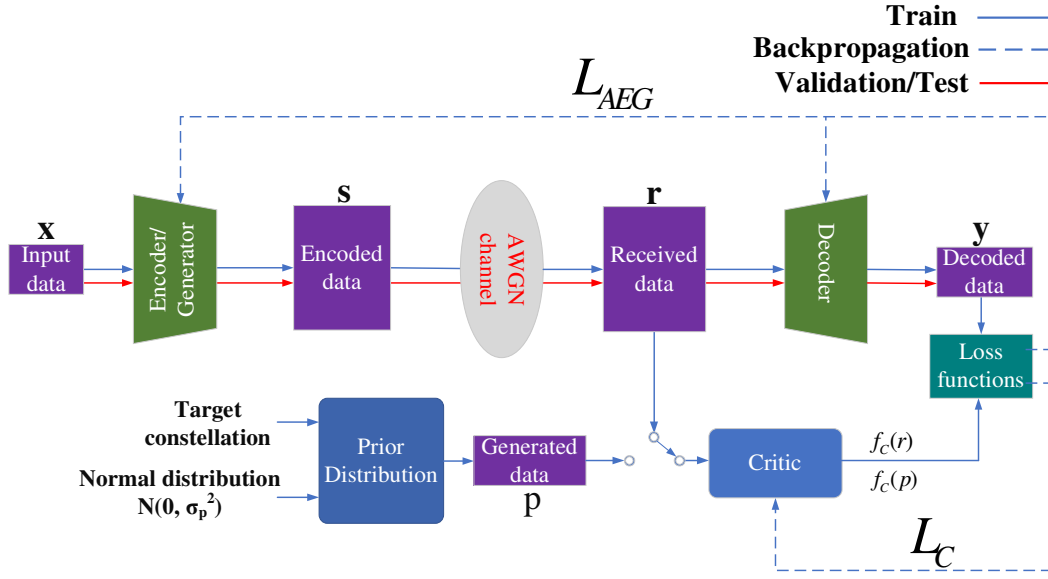


Figure 3.40: Overview of the proposed WGA-SCMA architecture.

$$L_C^{[k]} = \frac{1}{T} \sum_{t=1}^T f_{c_k} \left( p_t^{[k]} \right) - \frac{1}{T} \sum_{t=1}^T f_{c_k} \left( r_t^{[k]} \right), \quad (3.108)$$

where  $T$  is the size of the test set,  $p_t^{[k]}$  is the  $t$ th point of the distribution  $\mathbf{p}^{[k]}$ , and  $r_t^{[k]}$  is the  $t$ th point of the distribution  $\mathbf{r}^{[k]}$ . As regards the Generator side, for each dimension  $k$ , the following generator loss function is minimized during the training:

$$L_G^{[k]} = -\frac{1}{T} \sum_{t=1}^T f_{c_k} \left( r_t^{[k]} \right). \quad (3.109)$$

Starting from Fig. 3.39a, during the training procedure, the weights of the Critic and Generator are updated so that, epoch-by-epoch, the points of  $\mathbf{r}^{[k]}$  spread in the complex plane (see Figs. 3.39b-c) and, after several epochs, fit the distribution  $\mathbf{p}^{[k]}$ , as shown in Fig. 3.39d. Since it is up to the Generator to create the optimal codebook, this network act as Encoder.

As regards the task of providing an optimal decoder at the receiver side, we consider the autoencoder structure proposed in [105]. The Encoder is composed of  $K \cdot d_f$  Deep Neural Networks (DNNs), denoted as  $E_{kj}$ , each one related to layer  $j$  and resource  $k$ . The Encoder  $E_{kj}$  receives  $\mathbf{x}_j$  in input and gives  $e_{kj}$  as output. The input node  $\mathbf{x}_j$  and the

### 3.6. A Wasserstein GAN autoencoder for SCMA networks

Encoder  $E_{kj}$  are interconnected each other by following the factor graph  $\mathbf{F}$ , i.e., the node  $\mathbf{x}_j$  is connected to  $E_{kj}$  if  $F_{kj} = 1$ . Then, for each dimension  $k$ , the  $e_{kj}$  values are added together to obtain  $s_k$ . Finally, the noise  $n_k$  is added to  $s_k$ , obtaining  $r_k$ . The Decoder part receives  $\mathbf{r}$  in input aiming to give in output  $\mathbf{y} = \mathbf{x}$ . Since the data streams are multiplexed over  $K$  resources, a unique Decoder that combines all information spread over all the resources is adopted. During the training phase, the authors of [105] adopt the Mean Squared Error (MSE) as loss function, i.e., the Decoder is trained such that the difference between  $\mathbf{x}$  and  $\mathbf{y}$  is minimized. However, this approach is not optimized for the decoding task, since the goal of the SCMA decoder is to identify which data symbol  $\mathbf{x}_j$  out of  $M$  has been transmitted by the layer  $j$ . In other words, it is a mutually exclusive classification problem. Keeping this in mind, in our implementation, unlike [105], we adopt the one-hot encoding representation for the input vector  $\mathbf{x}_j, \forall j = 1, \dots, J$ , i.e., we represent it as an  $M$ -dimensional vector, where all of the elements are 0, except for one, which has 1 and represents the transmitted data symbol. In addition, we choose the softmax as the activation function for the output layer, and the Binary Cross-Entropy (BCE) function as autoencoder loss function. For the dimension  $k$ , this loss function is calculated as

$$L_{AE}^{[k]} = -\frac{1}{T} \sum_{t=1}^T \left[ \mathbf{x}_t^{[k]} \log \left( \mathbf{y}_t^{[k]} \right) + \left( 1 - \mathbf{x}_t^{[k]} \right) \log \left( 1 - \mathbf{y}_t^{[k]} \right) \right]. \quad (3.110)$$

The overall proposed architecture is called WGA-SCMA and is reported in Fig. 3.40. Since the Generator coincides with the Encoder, it is denoted as Encoder/Generator. The top row represents the autoencoder, while the bottom row the Critic part of the WGAN.

#### Training, validation, and testing procedures

As regards the training procedure, in Fig. 3.40 the blue solid line shows the dataflow at the train phase, and the blue dotted line the related back-propagation flow. The Encoder/Generator and the Critic have different training processes, and may be affected by the divergence phenomenon. To overcome this issue, we train our model by following different phases and adopting a model validation procedure for implementing the regu-

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios

---

larization by early stopping the training procedure when the error on the validation set increases or keeps constant for several epochs.

First, we train the Encoder/Generator and the Decoder for one batch. At this aim, we introduce the following autoencoder-generator loss function.

$$L_{AEG} = \frac{1}{K} \sum_{k=1}^K [\beta L_G^{[k]} + (1 - \beta) L_{AE}^{[k]}], \quad (3.111)$$

where  $\beta \in [0, 1]$  weighs the importance given to the generator function rather than to the autoencoder. The weights and the biases of both the Encoder/Generator and the Decoder entity are updated using the Gradient Descent (GD) method. We keep the Critic constant during the autoencoder-generator training phase.

Second, we train the  $K$  Critics by using the same batch and adopting the loss functions (3.108). Similarly, we keep both the Encoder/Generator and the Decoder constant during this phase, and the weights and the biases of the Critic are updated using the GD method. Phase 1 and 2 are repeated until one epoch is completed.

Third, we perform the model validation procedure. As shown in Fig. 3.40, this phase follows only the top row, i.e., the autoencoder. We give as input of the Encoder/Generator the validation set, and evaluate the  $L_{AE}^{[k]}$  values, considering several  $E_b/N_0$  values. Phase 1, 2 and 3 are repeated until the training procedure should be stopped.

Finally, the test procedure is similar to the validation one, but it adopts the test set.

#### 3.6.3 Performance evaluation

In this section, we compare the performance of the proposed WGA-SCMA with the DL-SCMA. Moreover, a comparison with the MAP and the Log-MPA decoder with 3 iterations is carried out, using the largely used codebook set given in [106].

##### Simulation setup

In order to train, validate and test our WGAN-SCMA model, we artificially created by simulations the datasets  $\mathbf{x}$  and  $\mathbf{n}$ , with  $J = 6$ ,  $K = 4$ ,

### 3.6. A Wasserstein GAN autoencoder for SCMA networks

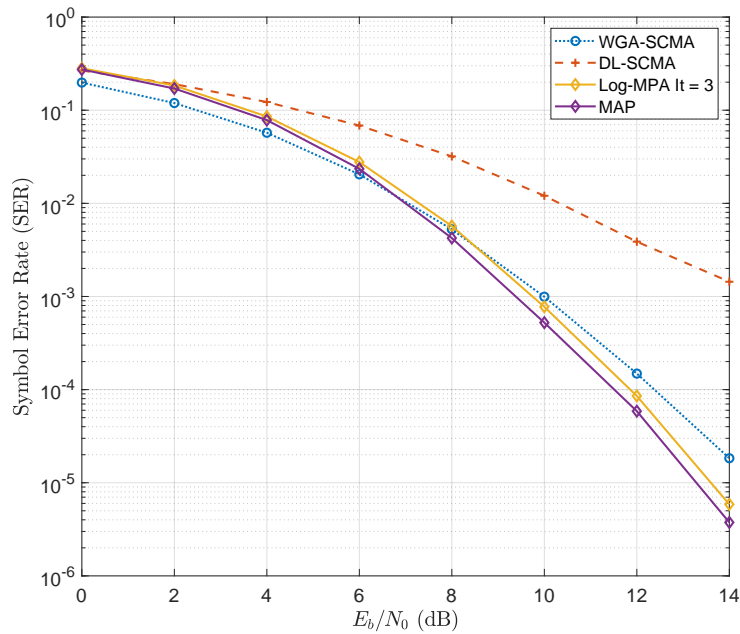
$S = 2$ ,  $M = 4$ , and  $d_f = 3$ . The software environment includes Python 3.9.7, the PyTorch 1.9.0 tensor library for deep learning, and the NVIDIA CUDA Toolkit version 11.3.0 for GPU-accelerated training and testing. The input dataset  $\mathbf{x}$  contains all the  $J$ -tuple in  $\mathbf{B}$  repeated several times so that its cardinality is equal to 60 millions. The dataset  $\mathbf{n}$  is composed of 12 groups, each one related to  $E_b/N_0 = q$  dB, where  $q \in \{0, 2, 4, \dots, 20\} \cup \{5\}$ . Finally, the 80% of the generated points compose the training set, while the remaining points are halved between test and validation sets. Through a large number of experiments, the best performances were obtained by adopting as training set the noise vector  $\mathbf{n}$  related to  $E_b/N_0 = 5$  dB. As regards the validation set, we adopted all the noise vectors  $\mathbf{n}$ , with the exception of that related to  $E_b/N_0 = 5$  dB. As regards the prior distribution, we adopted as target constellation the 64-QAM and as normal distribution the dataset  $\mathbf{n}$ , with  $E_b/N_0 = 20$  dB.

#### Implementation of the benchmark schemes and comparison

We implemented the DL-SCMA model presented in [105], trained and tested it by using the same training and test sets adopted for our model. As regards the mathematical decoders, we adopted the codebook set presented in [106], and the performance are evaluated by adopting our test sets. Let us note that, in order to work well, these mathematical methods need to be implemented considering the perfect knowledge of the noise power  $\sigma_N^2$ . Instead, both DL-SCMA and the proposed model do not need this ideal knowledge.

In Fig. 3.41, we compare the Symbol Error Rate (SER) below different values of  $E_b/N_0$ . As shown, WGAN-SCMA exhibits a notable improvement over DL-SCMA especially for low values of  $E_b/N_0$ . The improvements obtained depend both on the introduction of the innovative WGAN approach and on the proposed multi-phase training. We emphasize that these improvements are obtained without any increment of the latency. In fact, during the runtime part, i.e., the inference phase, only the Decoder is exploited. Since the Decoder structure is very similar to that of the DL-SCMA, it has the same computational complexity, that consequently results much lower than that of mathematical methods [105]. In this regard we also observe in Fig. 3.40 that the WGAN-SCMA shows performance similar to that showed by conventional decoding schemes

### Chapter 3. Traffic Prediction, Resource allocation, and En/Decoding strategies for SCMA-based mMTC scenarios



**Figure 3.41:** *SER comparison.*

although the computational complexity is much lower. These results are promising. In fact, we believe that our approach has good possibility to be improved even with respect to conventional MAP and Log-MPA decoders, since these systems have been tested assuming that the  $\sigma_N^2$  value is perfectly known.

---

## Learning Generalized Wireless MAC Communication Protocols via Abstraction

---

### 4.1 Overview

---

Envisioning to support heterogeneous services and applications of future Beyond 5G (B5G) and 6G networks, in this Chapter, we consider new types of communication protocols tailored to specific applications. In this context, Machine Learning (ML) can be used to design new protocols with reduced time, effort, and cost compared to conventional methods [107]. In particular, Multi-Agent Reinforcement Learning (MARL) [108] methods enable agents to learn an optimal policy by interacting with non-stationary environments. Recent advances in deep Reinforcement Learning (RL) and learning-to-communicate techniques (e.g., Differentiable Inter-Agent Learning (DIAL), and Reinforced Inter-Agent Learning (RIAL) [109]) have led to the emergence of protocol learning for the Physical (PHY) and Medium Access Control (MAC) layers [110–113]. Among them, to the best of our knowledge, the only works that assess the problem of MAC protocol learning in both control and data planes are [110, 111]. Therein, User Equipments (UEs) are cast as agents that learn from their partial observation of the global state how

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction

---

to deliver MAC Protocol Data Units (PDUs) to the Base Station (BS) throughout the radio channel. To generate optimal policies, the Centralized Training and Decentralized Execution (CTDE) method is adopted, where agents are trained offline using centralized information but execute in a decentralized manner. Specifically, in [110], both the BS and the UEs are cast as RL agents, and the multi-agent deep deterministic policy gradient (MADDPG) algorithm is adopted, that is, a commonly used CTDE-based actor-critic method. In [111], the BS is modeled as an expert agent adopting a predefined protocol, while the UEs are RL agents trained to learn a shared channel-access policy following the target signaling policy set by the BS. The policy is learned by exploiting a tabular Q-learning algorithm that follows the CTDE method. However, despite the good performances showed in the training environments, the learned protocols (i.e., policies) fail to generalize outside of their training distribution, as showed in [111]. We note that this drawback stems from the fact that agents learn their policies in the observational space that is specified for the environment instead of learning observation representations that are invariant over multiple environments, which enables better generalization and robustness. The notion of abstraction is based on learning task structure and invariance across tasks, while filtering out irrelevant information [114]. Leveraging abstraction when learning a communication protocol requires tackling two questions: i) how many abstracted observations agents need for optimal decision-making? and ii) how do agents choose their expert policy (or a set thereof) to extract essential information?

The main contribution of this Chapter is to leverage abstraction to learn new wireless MAC communication protocols with good generalization capabilities compared to state-of-the-art solutions. Towards this, we consider the same communication scenario presented in [111] and introduce the concept of observation abstraction. Specifically, we first present a new Autoencoder (AE) architecture to calculate the optimal abstraction space. Then, we solve the cooperative MARL problem by adopting the Multi-Agent Proximal Policy Optimization (MAPPO) algorithm [115] in the obtained abstracted observation space. We adopt MAPPO since it is one of the most promising algorithms following the CTDE approach for addressing cooperative MARL tasks [115]. Finally,



the performances of the proposed solution are compared with the same MARL problem solved by adopting the MAPPO without observation abstraction and with [111]. Simulation results show that the proposed approach yields policies that perform well not only in the training environment (as in the benchmark solutions), but also, and more importantly, in new and more complex environments that change in terms of number of PDUs, number of UEs, and channel conditions.

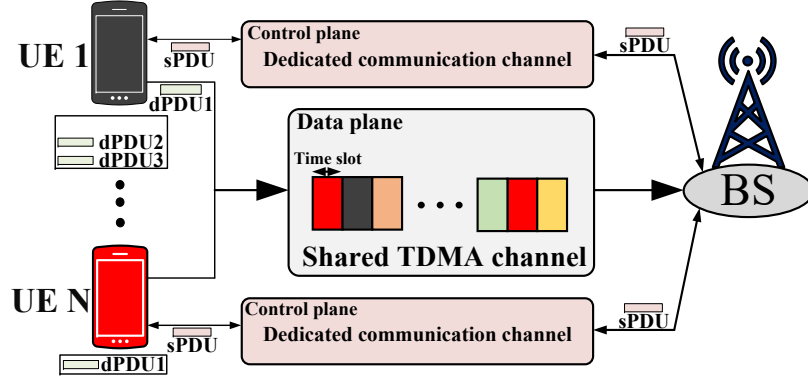
## 4.2 System Model

---

We consider an uplink radio network composed by a set  $\mathcal{N}$  of  $N$  homogeneous UEs and one BS, as shown in Fig. 4.1. We consider both data plane, where the UEs transmit the Uplink (UL) MAC Protocol Data Unit (PDU) to the BS, and control plane, where UEs exchange with the BS signaling MAC PDUs. In the following, we denote a data PDU transmitted in the data plane with "dPDU", a signaling PDU transmitted in the control plane with "sPDU", and a dPDU successfully transmitted to the BS and correctly deleted from the buffer with "dPDU successfully delivered". The task of each UE  $i \in \mathcal{N}$  is to successfully deliver  $P$  dPDUs. For sake of simplicity and to compare with [111], we adopt a Time Division Multiple Access (TDMA) channel access scheme. In the absence of transmission errors introduced by the radio channel, a dPDU is successfully received by the BS only if a single UE out of  $N$  has transmitted its dPDU. If multiple UEs simultaneously transmit their dPDUs, a collision occurs and the BS cannot correctly decode the received dPDUs.

To successfully transmit their dPDUs in the same contended radio channel, UEs (MAC learning agents) should learn an efficient MAC protocol. In the control plane, the learned MAC protocol should provide the optimal exchange of sPDUs between the UE and the BS to send the dPDU. Let  $\mathcal{M}_{\text{UE}}$  be the set of possible uplink control plane messages,  $a_{i,s} \in \mathcal{M}_{\text{UE}}$  be the signaling message sent by the  $i$ th UE, and  $\mathcal{M}_{\text{BS}}$  be the set of downlink (DL) control messages. Moreover, we consider that the data plane transmissions are modeled as a packet erasure channel, i.e., the dPDU is correctly received with a fixed probability equal to the Transport Block Error Rate (TBLER). Conversely, we assume that the control channels are error-free and dedicated to each UE.

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction



**Figure 4.1:** High-level depiction of the system model.

At each time slot  $t$ , each UE can send one sPDU to the BS in the dedicated control plane and one dPDU in the shared data plane. Each UE  $i$  has a dPDUs storage capability, modeled as a buffer with First-In First-Out (FIFO) policy, which can contain at most  $P \leq Q$  dPDUs. We denote with  $b_i^t \in \mathcal{B} = \{0, 1, \dots, Q\}$  the buffer status at time  $t$ , and we assume  $b_i^0 = P$  for all  $i \in \mathcal{N}$ . At each time slot  $t$ , a UE can only either transmit the first dPDU in the buffer or delete it. This means that the second dPDU in the buffer can only be considered after the first dPDU has been deleted.

Furthermore, we assume that the BS is a MAC expert agent, i.e., it adopts a protocol that is not learned and it operates only in the DL control plane. Specifically, for each time slot  $t$ , the BS sends a control message  $m_i^t \in \mathcal{M}_{\text{BS}} = \{0, 1, 2\}$  to each UE  $i$ . Here,  $m_i^t = 2$  represents an ACK message that confirms a dPDU sent from UE  $i$  has been correctly received at the BS in the previous time slot  $t - 1$ ,  $m_i^t = 1$  refers to a scheduling grant message to the UE  $i$ , and  $m_i^t = 0$  to indicate that no access is granted for the UE  $i$ . Clearly, the  $m_i^t = 2$  message can be sent to one UE at most, since the dPDU can be successfully received only if a collision was not occurred. We note that, if the UE  $i$  had successfully transmitted the dPDU to the BS concurrently with the access request, then we set  $m_i^t = 2$ . As regards  $m_i^t = 1$ , since only one UE can be scheduled each time slot, the BS sends this message to one UE randomly chosen from the ones having transmitted the access request in  $t - 1$  and with  $m_i^t \neq 2$ .

### 4.3 UEs-BS Interaction as an MARL Problem

In the considered system, the UE cannot receive information about other UEs, but makes decisions only on the basis of its observation of the global state. Moreover, UEs collaborate with one another to avoid collisions in the shared uplink radio channel, and thus they share the same global reward. As a consequence, the protocol learning problem is cast as a cooperative and Multi-Agent Partially Observable Markov Decision Process (MPOMDP) defined by  $\langle \mathcal{N}, \mathcal{A}, \mathcal{S}, \mathcal{O}, \pi_i, R, \gamma \rangle$ .

- $\mathcal{N}$ : set of all agents.
- $\mathcal{A}$ : shared action space. In this work, each agent  $i \in \mathcal{N}$  shares the same action space. Agent  $i$  performs an action  $a_i = (a_{i,u}, a_{i,s}) \in \mathcal{A}$ , that involves both data and control plane, where the data plane action  $a_{i,u} \in \{0, 1, 2\}$  and  $a_{i,s} \in \mathcal{M}_{\text{UE}} = \{0, 1\}$ . Specifically,  $a_{i,u} = 1$  means that the agent transmits the first dPDU in its buffer (if any),  $a_{i,u} = 2$  means it deletes the first dPDU in the buffer, and  $a_{i,u} = 0$  to do nothing. For the control plane,  $a_{i,s} = 1$  means sending an access request in order to reserve one time slot for its own transmission in the next time slot, while  $a_{i,s} = 0$  means do not transmit any signaling message.
- $\mathcal{S}$ : state space of the environment. At time  $t$ , the state  $s^t \in \mathcal{S}$  describes the environment by  $s^t = (\mathbf{b}^t, \mathbf{b}^{t-1}, \mathbf{a}^{t-1}, \mathbf{m}^{t-1}, \dots, \mathbf{b}^{t-M}, \mathbf{a}^{t-M}, \mathbf{m}^{t-M})$ , where  $\mathbf{b}^t = [b_1^t, b_2^t, \dots, b_N^t]$  is the vector containing all the buffer states,  $\mathbf{a}^t = [a_1^t, a_2^t, \dots, a_N^t]$  is the joint action vector,  $\mathbf{m}^t = [m_1^t, m_2^t, \dots, m_N^t]$  is the vector containing the DL control messages  $m_i^t$  received by each agent  $i$  from the BS, and  $M$  is the memory length.
- $\mathcal{O}$ : set of possible observations for each agent  $i$ . Each agent shares the same observation space. At time  $t$ , each agent has a partial observation of the global state  $s^t \in \mathcal{S}$ , defined as  $o_i^t \in \mathcal{O}$ . Its observation is the tuple  $o_i^t = (b_i^t, b_i^{t-1}, a_i^{t-1}, m_i^{t-1}, \dots, b_i^{t-M}, a_i^{t-M}, m_i^{t-M})$ .
- $\pi_i$ : policy of agent  $i$ , that is the probability of choosing a given action  $a_i$  given its partial observation  $o_i$ :

$$\pi_i: \mathcal{O} \rightarrow \Delta(\mathcal{A}), \quad (4.1)$$

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction

---

where  $\Delta$  defines a probabilistic space. Specifically, we denote with  $\pi_i(o_i, a_i)$  the probability to take  $a_i$  when observing  $o_i$ , and with  $\pi_i(o_i)$  the probability distribution among all possible actions in  $\mathcal{A}$  given observation  $o_i$ .

- $R \in \{-1, -\rho, +\rho\}$  is the global reward which quantifies the benefit of the joint actions performed by the agents. In this regard, the agents are penalized in the following case:
  1. if there exists an agent that deletes the dPDU without having previously transmitted it with success.

Instead, the agents are positively rewarded under these conditions:

1. if there exists an agent that deletes its dPDU having previously transmitted it with success.
2. If there exists an agent that has transmitted with success the dPDU for the first time.

Given these conditions, at the end of each time step  $t$ , each agent  $i \in \mathcal{N}$  receives the same global reward  $R^t$  as follows:

$$R^t = \begin{cases} -\rho & \text{if 1) is **True**,} \\ +\rho & \text{if 1) is **False** \wedge} \\ & \text{[ 2) is **True** \vee 3) is **True**],} \\ -1 & \text{otherwise.} \end{cases} \quad (4.2)$$

We underline that we set  $R^t = -1$  when no condition is true to minimize the number of time slots. The values assigned to the reward  $R^t$  follow from [110]. However, unlike that work, where the agents are positively rewarded if condition 3) is true, we also give a positive reward if condition 2) is true, since it allows the agent to transmit, in the subsequent time steps the next packet in the buffer.

- $\gamma \in [0, 1]$ : discount factor, which determines the impact of future rewards on the current decision. Therefore, we define the discounted accumulated reward  $G^t$  at time  $t$  as:

$$G^t = R^t + \gamma R^{t+1} + \gamma^2 R^{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R^{t+k}. \quad (4.3)$$

Due to the homogeneous nature of UEs, i.e., they share the same action space  $\mathcal{A}$ , the observation space  $\mathcal{O}$ , and the global reward  $R$ , instead of finding an optimal policy  $\pi_i^*$  for each UE  $i$ , we learn a shared optimal policy  $\pi^*$  via the parameter sharing technique [116].

To learn it, we adopt the CTDE paradigm. In general, several MARL techniques can be used, ranging from off-policy learning frameworks, such as MADDPG [110], value-based approaches (e.g., tabular Q-learning [111]), to on-policy algorithms such as MAPPO [115]. Among them, in this work we adopt MAPPO, since its on-policy nature is well-suited to the task of learning new MAC protocols.

## 4.4 Policy learning via abstraction

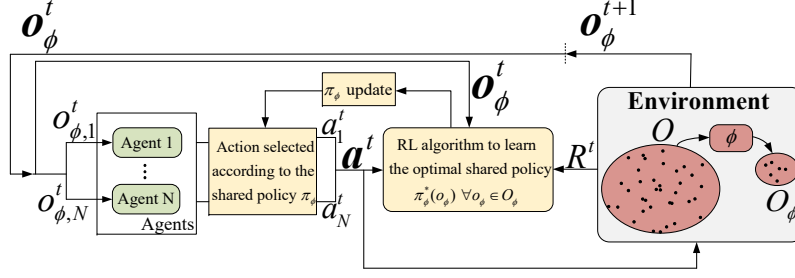
---

Typically, learning by abstraction is instrumental in reducing the size of the observation set  $\mathcal{O}$  that can be large and contains redundant information. This can be done by clustering and aggregating similar observations to form Abstracted Observations (AOs). In the context of RL, abstraction can overcome the fact that the policy is overfitted to a set of redundant and noisy observations, which hurts the ability to generalize. Concretely, if during the evaluation phase a new observation that was never encountered in the reduced training phase arises (e.g., a larger buffer dimension) the learned optimal policy will perform poorly. In this new environment, a new policy should be re-learned from scratch considering a large number of observations. In contrast, better generalization can be achieved by learning a policy in the abstracted observation space by finding the optimal solution in the presence, during the evaluation phase, of one or many never-seen observations that are mapped into the abstracted observation space.

### 4.4.1 Abstract Formulation

We define the abstracted MPOMDP by four components: the original MPOMDP presented in the previous Section, the observation abstraction (OA) function  $\phi$ , an abstraction of  $\mathcal{O}$  denoted as  $\mathcal{O}_\phi$ , and a shared abstracted policy operating on  $\mathcal{O}_\phi$  denoted as  $\pi_\phi$ . Specifically,  $\phi: \mathcal{O} \rightarrow \mathcal{O}_\phi$  maps each observation  $o_i \in \mathcal{O}$  into an abstracted observation  $o_{\phi,i} \in \mathcal{O}_\phi$ . The function is injective and each agent  $i$  makes use of the  $\phi$  function

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction



**Figure 4.2:** Proposed training procedure in the abstracted observation space.

as depicted in Fig. 4.2. At each time  $t$ , the environment provides each agent  $i$  the related partial observation of state  $s^t$ , denoted as  $o_i^t \in \mathcal{O}$ . This original observation is first passed through  $\phi$  yielding the abstracted state  $\phi(o_i^t) = o_{\phi,i}^t$ . Then, agent  $i$  uses  $o_{\phi,i}^t$  to take the action  $a_i^t$  according to  $\pi_\phi$ . After all agents take their actions, the environment provides one global reward  $R^t$ . This value, together with the vector of partial abstracted observations  $\mathbf{o}_\phi^t = [o_{\phi,1}^t, o_{\phi,2}^t, \dots, o_{\phi,N}^t]$  and the joint action vector  $\mathbf{a}^t$ , are used by the RL algorithm to update  $\pi_\phi$ . As a consequence, all agents learn the abstracted shared policy

$$\pi_\phi: \mathcal{O}_\phi \rightarrow \Delta(\mathcal{A}). \quad (4.4)$$

We resort to the concept of apprenticeship learning [117] to find a new representation  $\mathcal{O}_\phi$  with  $|\mathcal{O}_\phi| \ll |\mathcal{O}|$  that contains the most useful information yielding efficient decision-making, i.e.,  $\pi_\phi$ . Therein, learning  $\pi_\phi$  is carried out by observing an expert demonstrator following the policy  $\pi_E$  in the original observation domain. Hence, the goal of the observation abstraction is tantamount to compressing  $\mathcal{O}$  into  $\mathcal{O}_\phi$ , so that,  $\mathcal{O}_\phi$  provides agents with an effective understanding of the environment to allow them to follow the expert policy in the abstracted space. This gives rise to an interesting trade-off between *observation compression* and the ability of agents to follow the expert policy, expressed as a *divergence* between the expert policy  $\pi_E$  and the abstracted policy  $\pi_\phi$  in the compressed space  $\mathcal{O}_\phi$ . To quantify this divergence, we adopt the average Kullback-Leibler (KL) divergence:

$$d\{\pi_E, \pi_\phi\} = \mathbb{E}_{o \in \mathcal{O}} \{D_{\text{KL}}(\pi_E(o) \parallel \pi_\phi(\phi(o)))\}, \quad (4.5)$$

where

$$D_{\text{KL}}(\pi_{\text{E}}(o) \parallel \pi_{\phi}(\phi(o))) = \sum_{a \in \mathcal{A}} \pi_{\text{E}}(a, o) \log \left( \frac{\pi_{\text{E}}(a, o)}{\pi_{\phi}(a, \phi(o))} \right). \quad (4.6)$$

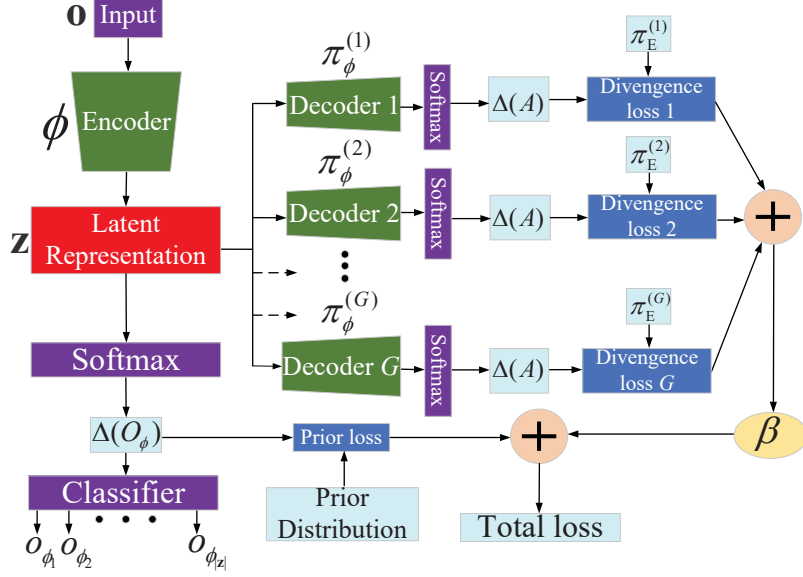
Departing from apprenticeship learning for the single-agent MDPs that rely on a unique optimal expert policy, for the multi-agent scenario of this interest, we allow agents to adopt and exploit the information gathered from different expert policies towards improving the robustness. In this view, we introduce a set  $\mathcal{P}_{\text{E}} = \{\pi_{\text{E}}^{(1)}, \pi_{\text{E}}^{(2)}, \dots, \pi_{\text{E}}^{(G)}\}$  of  $G$  expert policies defined in the original observation space  $\mathcal{O}$  and a set  $\mathcal{P}_{\phi} = \{\pi_{\phi}^{(1)}, \pi_{\phi}^{(2)}, \dots, \pi_{\phi}^{(G)}\}$  of corresponding abstracted policies defined in the abstracted observation space  $\mathcal{O}_{\phi}$ . The objective is re-defined as finding the optimal OA function with  $|\mathcal{O}_{\phi}| \ll |\mathcal{O}|$ , which minimizes the following divergence loss function:

$$L_{\text{div}} = \sum_{g=1}^G d \left\{ \pi_{\text{E}}^{(g)}, \pi_{\phi}^{(g)} \right\}. \quad (4.7)$$

To solve (4.7), we use AE architecture, which is composed of two Deep Neural Networks (DNNs), namely encoder and decoder. The encoder maps the high dimensional input into a low dimensional latent representation  $\mathbf{z}$  of size  $|\mathbf{z}|$  that contains only the important information needed to represent the original input. The decoder reproduces the original data from  $\mathbf{z}$  so that the output is a representation as close as possible to the original input. Both encoder and decoder are trained jointly to minimize the mean square error between the input and output.

Starting from the conventional AE architecture, we propose a new architecture represented in Fig. 4.3. Therein, the proposed AE receives the observations  $o \in \mathcal{O}$  as the input. The encoder reduces the cardinality of the input to ensure  $|\mathcal{O}_{\phi}| \ll |\mathcal{O}|$ , rather than reducing the input dimension, as in conventional AE. This is realized by enforcing the proposed encoder model to act as a multi-class classifier, in which, each sample  $o \in \mathcal{O}$  is assigned to one and only one abstracted observation  $o_{\phi_k} \in \mathcal{O}_{\phi}$  with  $k \in \{1, 2, \dots, |\mathbf{z}|\}$  and  $|\mathbf{z}| \ll |\mathcal{O}|$ . Note that  $|\mathcal{O}_{\phi}| \in \{1, 2, \dots, |\mathbf{z}|\}$  is held in general, since some abstracted observations  $o_{\phi_k}$  may not be assigned to any input  $o$ . Therefore, the encoder represents the observation abstraction function  $\phi$ . The decoder serves as an abstract policy network

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction



**Figure 4.3:** The proposed AE-based abstraction framework trading-off compression with value.

that maps each abstracted observation  $o_{\phi_k}$  to a distribution over action space  $\mathcal{A}$  instead of reconstructing the inputs as in conventional AE. Since we consider a set  $\mathcal{P}_E$  of  $G$  expert policies, we adopt  $G$  decoders, where each decoder is trained to produce the  $g$ th abstracted policy  $\pi_\phi^{(g)}$ . The aim of each  $g$ th network is to minimize the KL divergence with respect to  $\pi_E^{(g)}$  as per (4.7). Similar to the conventional AE, both encoder and decoders are jointly trained. Finally, the proposed loss function is composed of the sum of two parts. The first one, named divergence loss, aims to achieve the goal (4.7), while the second one, named prior loss, acts as a regularization term on the latent representation to make the distributions returned by the encoder close to a prior distribution  $\mathbf{p}$ . We propose to regularize the training with a prior distribution to avoid overfitting in the latent representation of the data so that the decoder networks can provide proper abstracted policies. For this, the regularization term is expressed as the KL divergence between the distribution at the output of the encoder and the prior  $\mathbf{p}$  as a uniform distribution among all the possible labels:

$$L_{\text{prior}} = \mathbb{E}_{o \in \mathcal{O}} \{D_{\text{KL}}(\Delta(O_\phi), \mathbf{p})\}. \quad (4.8)$$

The trade-off between the divergence loss and the regularization term is expressed by means of the hyper-parameter  $\beta \in \mathbb{R}_{\geq 0}$ . The total loss is



expressed as:

$$L_{\text{tot}} = L_{\text{prior}} + \beta L_{\text{div}}. \quad (4.9)$$

As  $\beta \rightarrow 0$ , the prior becomes more important, whereas as  $\beta \rightarrow \infty$ , minimizing divergence is prioritized. During training, the weights and biases of both encoder and decoder models are randomly initialized and updated via (4.9) by using the gradient descent (GD) method for  $N_{\text{abs}}$  episodes with the Adam optimizer and a learning rate  $l_{\text{abs}}$ . During evaluation, only the encoder part is adopted to provide, for each  $o_i \in \mathcal{O}$ , the proper label  $o_{\phi_k} \in \mathcal{O}_{\phi}$  at the output of the classifier.

---

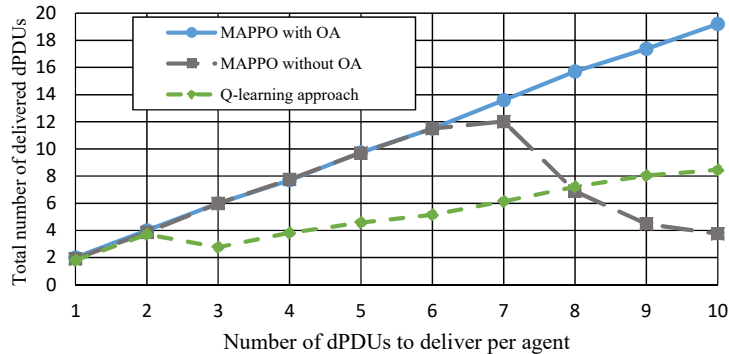
## 4.5 Performance Evaluation

In this section, we examine the performance of the proposed protocol learning approach leveraging abstraction, in terms of generalization to the number of dPDUs, to the TBLER, and to the number of UEs.

### 4.5.1 Setting

The encoder is a DNN composed of 3 hidden layers, each one with 512 neurons and the Rectified Linear Unit (ReLU) as activation function. We adopt 2 decoders (i.e.,  $G = 2$  expert policies), in which each decoder is a DNN with one hidden layer of 100 neurons and the ReLU as activation function. Moreover, we set  $l_{\text{abs}} = 25 \cdot 10^{-4}$  and  $N_{\text{abs}} = 10\,000$ ,  $\beta = 1000$ . The input set  $\mathcal{O}$  contains all possible arrangements between the elements in  $\mathcal{B}$  and the  $M$ -arrangements with repetition (i.e.,  $M$ -permutations with repetition) of the elements in  $\mathcal{B}$ ,  $\mathcal{A}$ , and  $\mathcal{M}_{\text{BS}}$ , with  $M = 1$  and  $P = 10$ . As a consequence,  $|\mathcal{O}| = 2178$ . As expert policies, we first adopt the conventional grant-based transmission, where the UE only transmits the dPDU following the reception of a scheduling grant, and deletes a dPDU following the reception of the ACK. A second expert policy is based on a grant-free transmission, where the UE transmits the dPDU immediately after it is available in the buffer, and deletes it after the transmission without waiting for the ACK message. The abstraction performance is evaluated as follows. Starting from  $|\mathbf{z}| = 1$  (i.e.,  $|\mathcal{O}_{\phi}| = 1$ ), we increased the size of  $|\mathbf{z}|$  with 1 unit until the loss (4.7) in the evaluation phase reached a plateau. The optimal cardinality of  $|\mathcal{O}_{\phi}|$  resulted equal to 8. Therefore, we adopt it for the subsequent simulations.

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction



**Figure 4.4:** Average total number of successfully delivered dPDUs by the  $N = 2$  agents.

For training, all the network parameters are: the number of UEs  $N = 2$ , the number of dPDUs to transmit  $P = 2$ , and the TBLER  $= 10^{-4}$ . Moreover, the training parameters together with the hyper-parameters are reported in Table 4.1. The training procedure for MAPPO follows the approach reported in Fig. 4.2, both adopting as observation space  $\mathcal{O}$  and the AO space  $\mathcal{O}_\phi$ . Thus, training generates two different solutions, named  $M_{\mathcal{O}}$  and  $M_{\mathcal{O}_\phi}$ , respectively. The performance evaluation takes into account the generation of different tasks. Each evaluation task is different in terms of  $P$ , TBLER, and  $N$ . All the performance results are obtained by averaging over 50 independent simulations per configuration, and carried out in Python environment.

### 4.5.2 Benchmarks

We compare the proposed solution, i.e.,  $M_{\mathcal{O}_\phi}$ , with the  $M_{\mathcal{O}}$  (no abstraction), and the approach proposed in [111], named  $Q_{\mathcal{O}}$ . In particular,  $Q_{\mathcal{O}}$  is trained by using the same hyper-parameters as in the original paper [111], and the same network parameters and reward structure of  $M_{\mathcal{O}_\phi}$  and  $M_{\mathcal{O}}$ . Moreover, due to the value-based nature of the  $Q_{\mathcal{O}}$  algorithm, it acts by choosing a random action when during evaluation it encounters an observation never seen during training.

### 4.5.3 Results and Discussion

We compare the solutions in terms of generalization to the number of dPDUs, to the TBLER, and to the number of UEs.

**Generalization to number of dPDUs.** Fig. 4.4 shows the performance in terms of total average number of successfully delivered packets when the evaluation is carried out keeping the same training parameters but  $P \in [1, 2, \dots, 10]$ . The results show that  $Q_{\mathcal{O}}$  performs well only for the value of  $P$  it was trained on, whereas its performance degrades for higher value of  $Q_{\mathcal{O}}$ . Conversely,  $M_{\mathcal{O}}$  shows an intrinsic generalization capability, since the on-policy training induces a probabilistic behavior (trajectory) that induces a good behavior within a certain range of variation. In this case, the performances are almost perfect for  $P < 7$ , while in the other cases a lower performance is incurred, achieving even lower performances than the  $Q_{\mathcal{O}}$  approach. Finally,  $M_{\mathcal{O}_{\phi}}$  exploits the intrinsic generalization capabilities of the on-policy algorithm and jointly reduces the uncertainties related to the different observation spaces through OA, achieving almost perfect performance for all considered ranges of  $P$ , achieving in the most difficult configuration, i.e.,  $P = 10$ , an increment of performance of 226.95% and 512% with respect to  $Q_{\mathcal{O}}$  and  $M_{\mathcal{O}}$ .

**Table 4.1:** Training algorithm Parameters

Common Parameter	Symbol	Value
Discount factor	$\gamma$	0.99
Epsilon value	$\epsilon$	0.1
Max. duration of episode (TTIs)	$t_{\max}$	300
Reward function parameter	$\rho$	3
$M_{\mathcal{O}}$ and $M_{\mathcal{O}_{\phi}}$ Parameter	Symbol	Value
Num. of neurons per hidden layer, evaluator		64
Num. of neurons per hidden layer, actor		64
Memory length	$M$	1
Learning rate	$lr_M$	$10^{-3}$
Number of training episodes	$N_{\text{tr}}$	20k
Act. function per layer, evaluator		{t, t, i} <sup>1</sup>
Act. function per layer, actor		{t, t, s}
Clipping value	$\psi$	0.2

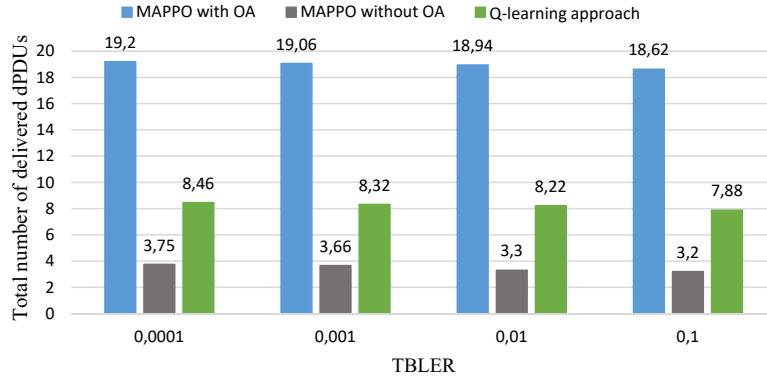
**Generalization to TBLEP.** To study the performances deviation related to a different value of TBLEP, in Fig. 4.5 we report the performance in terms of total average number of dPDUs successfully delivered

<sup>1</sup>t = tanh function, i = identity function

<sup>2</sup>t = tanh function, s = softmax function

## Chapter 4. Learning Generalized Wireless MAC Communication Protocols via Abstraction

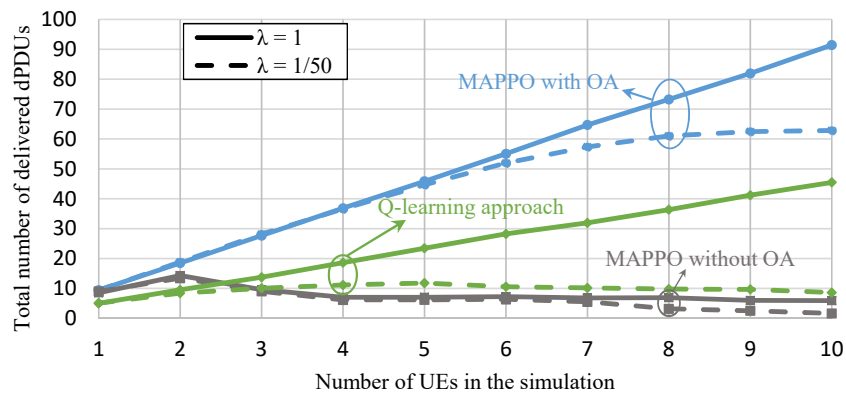
when the evaluation is carried out with  $P = 10$ ,  $N = 2$ , and  $\text{TBLER} \in [10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}]$ . We notice that the degradation of performance is contained for each method, and for each value of TBLER used in the evaluation. This result is obtained thanks to the technique of parameter sharing, that is used for each solution.



**Figure 4.5:** Average total number of delivered dPDUs under different solutions. Training procedure with  $\text{TBLER} = 10^{-4}$  but the performance is evaluated with different values of TBLER.

**Generalization to number of UEs.** Finally, in Fig. 4.6 we report generalization in terms of number of simultaneous active UEs while setting  $P = 10$  and  $\text{TBLER} = 10^{-4}$ . The UE arrivals are described by a Poisson distribution with a mean arrival rate of  $\lambda$ , and the simulations are carried out varying both the total number of UEs  $N$  in the simulation time and  $\lambda$ . The  $M_{\mathcal{O}_\phi}$  method significantly outperforms the other solutions in any condition in terms of average total number of dPDUs successfully delivered. However, when the number of simultaneous active UEs is high (i.e.,  $\lambda = 1$  and  $N \geq 7$ ), the performance starts to saturate.

## 4.5. Performance Evaluation



**Figure 4.6:** Average total number of delivered dPDUs under different number of agents and Poisson arrival rate ( $\lambda$ ).  $P = 10$  for each agent.



---

# CHAPTER 5

---

## Conclusions and Perspectives

---

### 5.1 Conclusions

---

In this thesis, we addressed the problem of studying and optimizing, by means of innovative techniques, several aspects the 5G NR interface that is facing services with very heterogeneous requirements. Specifically, we presented three research activities.

In the first one, we focused on the coexistence among a wide variety of services inside the same OFDM grid. We proposed a two-levels RRM framework that adaptively subdivides the band spectrum among the numerologies and properly allocates the PRBs to each numerology. Also, we implemented a new simulation environment constituted of the defined framework and a physical level simulator. We conducted a simulation study under different traffic types and channel conditions.

In the second research activity, we proposed a joint control of the dynamic resource allocation between access and transmission resources, and a new random access procedure based on an adaptive ACB. Then, we presented the idea of exploiting the unused PUSCH resources to serve an additional part of MTC devices. Moreover, we proposed an accurate current access attempts estimation method, based on DNN, which accepts

## Chapter 5. Conclusions and Perspectives

---

as input only the information really available at the gNB. Finally, to improve the transmission performance of SCMA in practical networks we designed an end-to-end SCMA en/decoding structure robust to the channel noise. Simulation results for each one of these enhancements shown improvements with respect to the state of the art solutions.

In the last research activity, we studied the problem of learning generalized MAC protocols that consider both user and control planes. To do so, we proposed a novel wireless MAC protocol learning framework for an uplink TDMA transmission scenario, based on abstraction. The simulation results showed that the proposed solution learns generalized MAC protocols that efficiently perform the transmission task, generalizing in terms of number of dPDUs to transmit, TBLER, and number of UEs.

### 5.2 Works in progress and Future works

---

The work done in this Dissertation in the different areas of 5G cellular networks can be considered as a step towards Beyond 5G (B5G) and 6G cellular networks. In fact, although 5G wireless systems are not fully deployed yet, B5G and 6G wireless systems are gaining more importance. These system are expected to require artificial intelligence as an essential component of their technology.

In fact, starting from the trend of 5G, the vision of B5G or 6G is to use largely machine learning techniques as tools to further optimize the performance of the wireless communication, to optimize communications building blocks or network function blocks. Inspired by the great success of the typical AI technologies, especially ML and DL in areas like computer vision, automatic speech recognition, and natural language processing, many researchers are attempting to introduce AI into mobile network systems with the capability to optimize a variety of wireless network problems [118]. In this context, during my PhD course, I spent abroad a period of 4 months and half at the Centre for Wireless Communications (CWC), Oulu, Finland. In this period, I deepened the knowledge about RL and DRL. Moreover, I attended, *inter alia*, the "Huawei Workshop on Intelligent IoT for 6G", and followed the "Machine Learning" course from the Department of Mathematics and Computer Sciences



## 5.2. Works in progress and Future works

---

of the University of Catania.

As regards works in progress, inspired by the work [19], we are developing a flexible encoding/decoding procedures for 6G SCMA Wireless Networks via adversarial machine learning techniques. Specifically, we defined a procedure for jointly training and validating the model. This procedure is used to avoid the overfitting problem, and also to provide a flexible en/decoding procedure so that the trained model is robust to different channel conditions. In addition, we are now proposing an efficient Next Generation Multiple Access (NGMA) scheme where several innovative solutions combine and integrate. The first goal is to maximize the spectral efficiency of the PUSCH by adopting the most appropriate NOMA introducing an efficient contention-based approach for data transmission in the unused PUSCH resources. Second, to set the uplink radio resources dimensioning by searching the optimal tradeoff between lowering the collision probability and providing enough radio resources in the PUSCH to permit the data transmission to all the succeeded access requests. Third, to make feasible the traffic load estimation, even when the PRACH is in overload condition. Fourth, to exploit Artificial Intelligence technologies to further optimize the performances obtained. On the other hand, in the 6G vision, the traditional communication protocols should not be designed from the scratch and handcrafted for a generic purpose, but tailored for the specific service and on-the-fly designed. For this reason, we are now exploring new techniques for enhancing the emergence of generalized wireless MAC protocols.

As future works, in the context of 6G SCMA Wireless Networks, we are also investigating an innovative procedure that detects a change in the channel realization and permits to iteratively retrain the model with the aim of further improving the decoding performances. As regards the context of MAC protocol learning, we are considering the introduction of meta-tuned RL algorithms to yield a faster training convergence in unseen environments with a considerably low computation complexity and energy saving.

The 6G network is seen as a revolutionary step to enable devices and networks to be more autonomous, robust, resilient and sustainable. They would be able to continuously adapt and generalize across different tasks, environments, and types of communication. To do so, we should move

## Chapter 5. Conclusions and Perspectives

---

away from machines that learn statistical models with no ability to understand what is happening, and move towards the ability to understand and reason about data, being able to infer what is missing and fixing what is incorrect.

---

---

## Bibliography

---

- [1] L. Miuccio, D. Panno, P. Pisacane, and S. Riolo, “Channel-aware and QoS-aware downlink resource allocation for multi-numerology based 5G NR systems,” in *2021 19th Mediterranean Communication and Computer Networking Conference (MedComNet)*, pp. 1–8, 2021.
- [2] ITU-T, “IMT Vision - Framework and overall objectives of the future development of IMT for 2020 and beyond,” Recommendation M.2083, International Telecommunication Union, Geneva, Sep. 2015.
- [3] 3GPP, “Physical Channels and Modulation,” TS 36.211, 3rd Generation Partnership Project (3GPP), 9 2019. Version 15.7.0.
- [4] 3GPP, “Study on New Radio Access Technology,” TR 38.802, 3rd Generation Partnership Project (3GPP), 09 2017. Version 14.2.0.
- [5] K. Pedersen, G. Pocovi, J. Steiner, and A. Maeder, “Agile 5G scheduler for improved E2E performance and flexibility for different network implementations,” *IEEE Communications Magazine*, vol. 56, no. 3, pp. 210–217, 2018.
- [6] G. Pocovi, B. Soret, K. I. Pedersen, and P. Mogensen, “MAC layer enhancements for ultra-reliable low-latency communications in cellular networks,” in *2017 IEEE International Conference on*

## Bibliography

---

- Communications Workshops (ICC Workshops)*, pp. 1005–1010, 2017.
- [7] A. Anand, G. de Veciana, and S. Shakkottai, “Joint scheduling of URLLC and eMBB traffic in 5G wireless networks,” *IEEE/ACM Transactions on Networking*, vol. 28, no. 2, pp. 477–490, 2020.
- [8] G. Wunder, P. Jung, M. Kasparick, T. Wild, F. Schaich, Y. Chen, S. T. Brink, I. Gaspar, N. Michailow, A. Festag, L. Mendes, N. Cassiau, D. Ktenas, M. Dryjanski, S. Pietrzyk, B. Eged, P. Vago, and F. Wiedmann, “5GNOW: non-orthogonal, asynchronous waveforms for future mobile applications,” *IEEE Communications Magazine*, vol. 52, no. 2, pp. 97–105, 2014.
- [9] Y. Endo, Y. Kishiyama, and K. Higuchi, “Uplink non-orthogonal access with MMSE-SIC in the presence of inter-cell interference,” in *2012 International Symposium on Wireless Communication Systems (ISWCS)*, pp. 261–265, 2012.
- [10] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, “Non-orthogonal multiple access (NOMA) for cellular future radio access,” in *2013 IEEE 77th Vehicular Technology Conference (VTC Spring)*, pp. 1–5, 2013.
- [11] Q. Qi and X. Chen, “Wireless powered massive access for cellular Internet of Things with imperfect SIC and nonlinear EH,” *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 3110–3120, 2019.
- [12] X. Zhang, G. Han, D. Zhang, D. Zhang, and L. Yang, “An efficient SCMA codebook design based on lattice theory for information-centric IoT,” *IEEE Access*, vol. 7, pp. 133865–133875, 2019.
- [13] S. R. Pokhrel, J. Ding, J. Park, O.-S. Park, and J. Choi, “Towards enabling critical mMTC: A review of URLLC within mMTC,” *IEEE Access*, vol. 8, pp. 131796–131813, 2020.
- [14] L. Miuccio, D. Panno, P. Pollara, and S. Riolo, “Implementation of a simulation environment for multi-numerology scenarios in 5G Vienna link level simulator,” in *2021 IEEE International Conference on Computing (ICOCO)*, pp. 134–139, 2021.

- 
- [15] L. Miuccio, D. Panno, P. Pisacane, and S. Riolo, “A QoS-aware and channel-aware Radio Resource Management framework for multi-numerology systems,” *Computer Communications*, vol. 191, pp. 299–314, 2022.
- [16] L. Miuccio, D. Panno, and S. Riolo, “Joint control of random access and dynamic uplink resource dimensioning for massive MTC in 5G NR based on SCMA,” *IEEE Internet of Things Journal*, vol. 7, no. 6, pp. 5042–5063, 2020.
- [17] L. Miuccio, D. Panno, and S. Riolo, “A new contention-based PUSCH resource allocation in 5G NR for mMTC scenarios,” *IEEE Communications Letters*, pp. 1–1, 2020.
- [18] L. Miuccio, D. Panno, and S. Riolo, “A DNN-based estimate of the PRACH traffic load for massive IoT scenarios in 5G networks and beyond,” *Computer Networks*, vol. 201, p. 108608, 2021.
- [19] L. Miuccio, D. Panno, and S. Riolo, “A Wasserstein GAN autoencoder for SCMA networks,” *IEEE Wireless Communications Letters*, vol. 11, no. 6, pp. 1298–1302, 2022.
- [20] L. Miuccio, S. Riolo, S. Samarakoon, D. Panno, and M. Benis, “Learning generalized wireless MAC communication protocols via abstraction,” *ArXiv*, vol. abs/2206.06331, 2022.
- [21] F. D’Urso, C. Santoro, and F. F. Santoro, “Wale: A solution to share libraries in docker containers,” *Future Generation Computer Systems*, vol. 100, pp. 513–522, 2019.
- [22] A. Yazar, B. Pekoz, and H. Arslan, “Fundamentals of multi-numerology 5G New Radio,” 2019.
- [23] 3GPP, “NR; NR and NG-RAN Overall description; Stage-2,” TS 38.300, 3rd Generation Partnership Project (3GPP), 01 2021. Version 16.4.0.
- [24] M. Kawser, H. Farid, A. Hasin, A. Sadik, and I. Razu, “Performance comparison between Round Robin and Proportional Fair scheduling methods for LTE,” *International Journal of Information and Electronics Engineering*, vol. 2, no. 5, pp. 678–681, 2012.

## Bibliography

---

- [25] D. Panno and S. Riolo, “A new joint scheduling scheme for GBR and non-GBR services in 5G RAN,” in *2018 IEEE 23rd International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, pp. 1–6, Sep. 2018.
- [26] X. Zhang, M. Jia, L. Chen, J. Ma, and J. Qiu, “Filtered-OFDM- enabler for flexible waveform in the 5th generation cellular networks,” in *2015 IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6, 2015.
- [27] V. Vakilian, T. Wild, F. Schaich, S. ten Brink, and J.-F. Frigon, “Universal-filtered multi-carrier technique for wireless systems beyond LTE,” in *2013 IEEE Globecom Workshops (GC Wkshps)*, pp. 223–228, 2013.
- [28] X. Zhang, L. Zhang, P. Xiao, D. Ma, J. Wei, and Y. Xin, “Mixed numerologies interference analysis and inter-numerology interference cancellation for windowed OFDM systems,” *IEEE Transactions on Vehicular Technology*, vol. 67, no. 8, pp. 7047–7061, 2018.
- [29] S. Pratschner, B. Tahir, L. Marijanovic, M. Mussbah, K. Kirev, R. Nissel, S. Schwarz, and M. Rupp, “Versatile mobile communications simulation: the Vienna 5G Link Level Simulator,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2018, p. 226, Sept. 2018.
- [30] A. Akhtar and H. Arslan, “Downlink resource allocation and packet scheduling in multi-numerology wireless systems,” in *2018 IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*, pp. 362–367, 2018.
- [31] E. Memisoglu, A. B. Kihero, E. Basar, and H. Arslan, “Guard band reduction for 5G and beyond multiple numerologies,” *IEEE Communications Letters*, vol. 24, no. 3, pp. 644–647, 2020.
- [32] A. B. Kihero, M. S. J. Solaija, and H. Arslan, “Inter-numerology interference for beyond 5G,” *IEEE Access*, vol. 7, pp. 146512–146523, 2019.

- [33] 3GPP, “NR; Physical layer procedures for data,” TS 38.214, 3rd Generation Partnership Project (3GPP), 01 2021. Version 16.4.0.
- [34] 3GPP, “E-UTRA; User Equipment (UE) radio transmission and reception,” TS 36.101, 3GPP, 01 2021. Version 16.8.0.
- [35] 3GPP, “NR;Physical channels and modulation,” TS 38.211, 3rd Generation Partnership Project (3GPP), 04 2020. Version 16.1.0.
- [36] 3GPP, “System architecture for the 5G System (5GS),” TS 23.501, 3GPP, 09 2021. Version 17.2.0.
- [37] R. Jain, D. M. Chiu, and W. R. Hawe, “A quantitative measure of fairness and discrimination for resource allocation in shared computer systems,” *ArXiv Computer Science e-prints*, 1998.
- [38] H. Li, Q. Guo, L. Fang, and D. Huang, “Fairness and capacity analysis of opportunistic feedback protocol with proportional fair or maximum throughput scheduling,” in *2012 International Conference on Wireless Communications and Signal Processing (WCSP)*, pp. 1–5, Oct 2012.
- [39] A. Yazar and H. Arslan, “A flexibility metric and optimization methods for mixed numerologies in 5G and beyond,” *IEEE Access*, vol. 6, pp. 3755–3764, 2018.
- [40] E. Memisoglu, A. E. Duranay, and H. Arslan, “Numerology scheduling for PAPR reduction in mixed numerologies,” *IEEE Wireless Communications Letters*, vol. 10, no. 6, pp. 1197–1201, 2021.
- [41] A. Yazar and H. Arslan, “Reliability enhancement in multi-numerology-based 5G new radio using INI-aware scheduling,” *EURASIP Journal on Wireless Communications and Networking*, 2019.
- [42] J. Abdoli, M. Jia, and J. Ma, “Filtered OFDM: A new waveform for future wireless systems,” in *2015 IEEE 16th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, pp. 66–70, 2015.

## Bibliography

---

- [43] X. Zhang, L. Chen, J. Qiu, and J. Abdoli, "On the waveform for 5G," *IEEE Communications Magazine*, vol. 54, no. 11, pp. 74–80, 2016.
- [44] S. Pratschner, M. K. Muller, F. Ademaj, A. Nabavi, B. Tahir, S. Schwarz, and M. Rupp, "Verification of the Vienna 5G Link and System Level Simulators and their interaction," in *2019 16th IEEE Consumer Comms. Network. Conference (CCNC)*, pp. 1–8, 2019.
- [45] L. Miuccio, D. Panno, P. Pisacane, and S. Riolo, "Channel-aware and QoS-aware downlink resource allocation for multi-numerology based 5G NR systems," in *2021 19th Mediterranean Communication and Computer Networking Conference (Med-ComNet)*, pp. 1–8, 2021.
- [46] A. Laya, L. Alonso, and J. Alonso-Zarate, "Is the random access channel of LTE and LTE-A suitable for M2M communications? A survey of alternatives," *IEEE Commun. Surveys Tuts.*, vol. 16, pp. 4–16, First 2014.
- [47] L. Ferdouse, A. Anpalagan, and S. Misra, "Congestion and overload control techniques in massive M2M systems: a survey," *Trans. on Emerging Telecomm. Technol.*, vol. 28, no. 2, p. e2936, 2017. e2936 ett.2936.
- [48] M. Taherzadeh, H. Nikopour, A. Bayesteh, and H. Baligh, "SCMA codebook design," in *2014 IEEE 80th Vehicular Technology Conference (VTC2014-Fall)*, pp. 1–5, Sep. 2014.
- [49] R. Hoshyar, F. P. Wathan, and R. Tafazolli, "Novel low-density signature for synchronous CDMA systems over AWGN channel," *IEEE Transactions on Signal Processing*, vol. 56, no. 4, pp. 1616–1626, 2008.
- [50] K. Xiao, B. Xiao, S. Zhang, Z. Chen, and B. Xia, "Simplified multiuser detection for SCMA with sum-product algorithm," in *2015 International Conference on Wireless Communications Signal Processing (WCSP)*, pp. 1–5, 2015.



- [51] J. Liu, G. Wu, S. Li, and O. Tirkkonen, "On fixed-point implementation of Log-MPA for SCMA signals," *IEEE Wireless Communications Letters*, vol. 5, no. 3, pp. 324–327, 2016.
- [52] J. Vidal, L. Tello-Oquendo, V. Pla, and L. Guijarro, "Performance study and enhancement of Access Barring for massive Machine-Type Communications," *IEEE Access*, vol. 7, pp. 63745–63759, 2019.
- [53] S. Duan, V. Shah-Mansouri, Z. Wang, and V. W. S. Wong, "D-ACB: Adaptive congestion control algorithm for bursty M2M traffic in LTE networks," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 12, pp. 9847–9861, 2016.
- [54] X. Tang, L. Li, L. Hao, and M. Peng, "Joint random access control scheme based on PRACH channel quality and access class barring," in *2020 IEEE 92nd Vehicular Technology Conference (VTC2020-Fall)*, pp. 1–5, 2020.
- [55] H. D. Althumali, M. Othman, N. K. Noordin, and Z. M. Hanapi, "Dynamic backoff collision resolution for massive M2M random access in cellular IoT networks," *IEEE Access*, vol. 8, pp. 201345–201359, 2020.
- [56] G. Lin, S. Chang, and H. Wei, "Estimation and adaptation for bursty LTE random access," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 4, pp. 2560–2577, 2016.
- [57] F. Vázquez-Gallego, C. Kalalas, L. Alonso, and J. Alonso-Zarate, "Contention tree-based access for wireless machine-to-machine networks with energy harvesting," *IEEE Transactions on Green Communications and Networking*, vol. 1, no. 2, pp. 223–234, 2017.
- [58] N. Jiang, Y. Deng, A. Nallanathan, and J. Yuan, "A decoupled learning strategy for massive access optimization in cellular IoT networks," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 3, pp. 668–685, 2021.

## Bibliography

---

- [59] 3GPP, “Service accessibility,” TS 22.011, 3rd Generation Partnership Project (3GPP), 06 2010. Version 9.4.0.
- [60] L. Tello-Oquendo, J. Vidal, V. Pla, and L. Guijarro, “Dynamic access class barring parameter tuning in LTE-A networks with massive M2M traffic,” in *2018 17th Med-Hoc-Net Workshop*, pp. 1–8, June 2018.
- [61] L. Tello-Oquendo, D. Pacheco-Paramo, V. Pla, and J. Martinez-Bauset, “Reinforcement learning-based ACB in LTE-A networks for handling massive M2M and H2H communications,” in *2018 IEEE International Conference on Communications (ICC)*, pp. 1–7, May 2018.
- [62] N. Zhang, G. Kang, J. Wang, Y. Guo, and F. Labeau, “Resource allocation in a new random access for M2M communications,” *IEEE Communications Letters*, vol. 19, pp. 843–846, May 2015.
- [63] L. Miuccio, D. Panno, and S. Riolo, “Dynamic uplink resource dimensioning for massive MTC in 5G networks based on SCMA,” in *European Wireless 2019 (EW 2019)*, (Aarhus, Denmark), May 2019.
- [64] T. Xue, L. Qiu, and X. Li, “Resource allocation for massive M2M communications in SCMA network,” in *2016 IEEE 84th Vehicular Technology Conference (VTC-Fall)*, pp. 1–5, Sep. 2016.
- [65] M. Polignano, D. Vinella, D. Laselva, J. Wigard, and T. B. Sorensens, “Power savings and QoS impact for VoIP application with DRX/DTX feature in LTE,” in *2011 IEEE 73rd Vehicular Technology Conference (VTC Spring)*, pp. 1–5, May 2011.
- [66] L. Miuccio, D. Panno, and S. Riolo, “Joint congestion control and resource allocation for massive MTC in 5G networks based on SCMA,” in *15th International Conference on Telecommunications (ConTEL 2019)*, (Graz, Austria), July 2019.
- [67] 3GPP, “Study on RAN Improvements for Machine-type Communications,” TS 37.868, 3rd Generation Partnership Project (3GPP), 10 2011. Version 11.1.0.

- [68] L. Miuccio, D. Panno, and S. Riolo, “Dynamic uplink resource dimensioning for massive MTC in 5G networks based on SCMA,” in *European Wireless 2019; 25th European Wireless Conference*, pp. 1–6, 2019.
- [69] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein generative adversarial networks,” in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, pp. 214–223, Aug. 2017.
- [70] S. Moon, H. Lee, and J. Lee, “SARA: Sparse code multiple access-applied random access for IoT devices,” *IEEE Internet of Things Journal*, vol. 5, pp. 3160–3174, Aug 2018.
- [71] S. Duan, V. Shah-Mansouri, Z. Wang, and V. W. S. Wong, “D-ACB: Adaptive congestion control algorithm for bursty M2M traffic in LTE networks,” *IEEE Transactions on Vehicular Technology*, vol. 65, pp. 9847–9861, Dec 2016.
- [72] 3GPP, “NR; Physical layer procedures for control,” Technical Specification (TS) 38.213, 3rd Generation Partnership Project (3GPP), 09 2019. Version 15.7.0.
- [73] B. Fekade, T. Maksymyuk, M. Kyryk, and M. Jo, “Probabilistic recovery of incomplete sensed data in IoT,” *IEEE Internet of Things Journal*, vol. 5, pp. 2282–2292, Aug 2018.
- [74] J. Lee and J. Lee, “Prediction-based energy saving mechanism in 3GPP NB-IoT networks,” in *Sensors*, 2017.
- [75] 3GPP, “Study on Non-Orthogonal Multiple Access (NOMA) for NR (Release 16),” TR 38.812, 3rd Generation Partnership Project (3GPP), 12 2018. Version 16.0.0.
- [76] GSMA, “5G Implementation Guideline,” July 2019. Version 2.0.
- [77] H. Nikopour and H. Baligh, “Sparse code multiple access,” in *2013 IEEE 24th Annual Intl. Sym. on PIMRC*, pp. 332–336, Sep. 2013.

## Bibliography

---

- [78] M. Moltafet, N. Mokari, M. R. Javan, H. Saeedi, and H. Pishro-Nik, "A new multiple access technique for 5G: Power Domain Sparse Code Multiple Access (PSMA)," *IEEE Access*, vol. 6, pp. 747–759, 2018.
- [79] M. Moltafet, N. M. Yamchi, M. R. Javan, and P. Azmi, "Comparison study between PD-NOMA and SCMA," *IEEE Transactions on Vehicular Technology*, vol. 67, pp. 1830–1834, Feb 2018.
- [80] M. Vameghestahbanati, I. D. Marsland, R. H. Gohary, and H. Yanikomeroğlu, "Multidimensional constellations for uplink SCMA systems - A comparative study," *CoRR*, vol. abs/1804.05814, 2018.
- [81] M. Tavares, D. Samardzija, H. Viswanathan, H. Huang, and C. Kahn, "A 5G lightweight connectionless protocol for massive Cellular Internet of Things," in *2017 IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*, pp. 1–6, March 2017.
- [82] H. S. Dhillon, H. Huang, and H. Viswanathan, "Wide-area wireless communication challenges for the Internet of Things," *IEEE Communications Magazine*, vol. 55, pp. 168–174, February 2017.
- [83] C. Kahn and H. Viswanathan, "Connectionless access for mobile cellular networks," *IEEE Communications Magazine*, vol. 53, pp. 26–31, Sep. 2015.
- [84] H. S. Jang, S. M. Kim, H. Park, and D. K. Sung, "An early preamble collision detection scheme based on tagged preambles for cellular M2M random access," *IEEE Transactions on Vehicular Technology*, vol. 66, pp. 5974–5984, July 2017.
- [85] H. Shariatmadari, R. Ratasuk, S. Iraji, A. Laya, T. Taleb, R. Järntti, and A. Ghosh, "Machine-type communications: current status and future perspectives toward 5G systems," *IEEE Communications Magazine*, vol. 53, pp. 10–17, Sep. 2015.
- [86] M. Centenaro, L. Vangelista, S. Saur, A. Weber, and V. Braun, "Comparison of collision-free and contention-based radio access

- protocols for the Internet of Things,” *IEEE Transactions on Communications*, vol. 65, pp. 3832–3846, Sep. 2017.
- [87] H. He, Q. Du, H. Song, W. Li, Y. Wang, and P. Ren, “Traffic-aware ACB scheme for massive access in machine-to-machine networks,” in *2015 IEEE Intl. Conf. on Commun.*, pp. 617–622, June 2015.
- [88] L. Zhe, A. Meizhen, B. Linhou, and Z. Enyong, “Correlation analysis on telemetry data of manned spacecraft,” in *2018 Chinese Control And Decision Conference (CCDC)*, pp. 377–380, June 2018.
- [89] 3GPP, “Radio Resource Control (RRC) protocol specification (Release 15),” TS 38.331, 3rd Generation Partnership Project (3GPP), 9 2019. Version 15.7.0.
- [90] W.-j. Liu, X.-l. Hou, and L. Chen, “Enhanced uplink non-orthogonal multiple access for 5G and beyond systems,” *Frontiers of Information Technology & Electronic Engineering*, vol. 19, pp. 340–356, Mar 2018.
- [91] K. Mikhaylov, V. Petrov, R. Gupta, M. A. Lema, O. Galina, S. Andreev, Y. Koucheryavy, M. Valkama, A. Pouttu, and M. Dohler, “Energy efficiency of Multi-Radio Massive Machine-Type Communication (MR-MMTC): Applications, challenges, and solutions,” *IEEE Communications Magazine*, vol. 57, pp. 100–106, June 2019.
- [92] Nokia, “DRX parameters in LTE,” Tech. Rep. R2-071285, March 2007. Version 13.1.0.
- [93] M. S. Mushtaq, S. Fowler, and A. Mellouk, “Power saving model for mobile device and virtual base station in the 5G era,” in *2017 IEEE International Conference on Communications (ICC)*, pp. 1–6, May 2017.
- [94] S. M. Birajdar and A. K. Tamboli, “LTE-A network through DRX configuration and development of power saving device,” in *Inter-*

## Bibliography

---

- national Journal of Engineering Trends and Technology (IJETT)*, vol. 36, pp. 102–105, June 2016.
- [95] A. B. Kihero, M. S. J. Solaija, and H. Arslan, “Inter-numerology interference for beyond 5G,” *IEEE Access*, vol. 7, pp. 146512–146523, 2019.
- [96] Y. Wu, C. Wang, Y. Chen, and A. Bayesteh, “Sparse code multiple access for 5G radio transmission,” in *2017 IEEE 86th Vehicular Technology Conference (VTC-Fall)*, pp. 1–6, Sep. 2017.
- [97] C. Wei, G. Bianchi, and R. Cheng, “Modeling and analysis of random access channels with bursty arrivals in OFDMA wireless networks,” *IEEE Transactions on Wireless Communications*, vol. 14, no. 4, pp. 1940–1953, 2015.
- [98] L. Miuccio, D. Panno, and S. Riolo, “Joint congestion control and resource allocation for massive MTC in 5G networks based on SCMA,” in *2019 15th Intern. Conf. on Telecomm. (ConTEL)*, pp. 1–8, 2019.
- [99] N. Jiang, Y. Deng, O. Simeone, and A. Nallanathan, “On-line supervised learning for traffic load prediction in framed-ALOHA networks,” *IEEE Communications Letters*, vol. 23, no. 10, pp. 1778–1782, 2019.
- [100] 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, 2020.
- [101] S. Vural, N. Wang, P. Bucknell, G. Foster, R. Tafazolli, and J. Muller, “Dynamic preamble subset allocation for RAN slicing in 5G networks,” *IEEE Access*, vol. 6, pp. 13015–13032, 2018.
- [102] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” *Journal of Machine Learning Research - Proceedings Track*, vol. 9, pp. 249–256, 01 2010.
- [103] J. Wang, S. Li, Z. An, X. Jiang, W. Qian, and S. Ji, “Batch-normalized deep neural networks for achieving fast intelligent fault diagnosis of machines,” *Neurocomputing*, vol. 329, 10 2018.

- 
- [104] A. Sivanathan, H. H. Gharakheili, F. Loi, A. Radford, C. Wijenayake, A. Vishwanath, and V. Sivaraman, “Classifying IoT devices in smart environments using network traffic characteristics,” *IEEE Transactions on Mobile Computing*, vol. 18, no. 8, pp. 1745–1759, 2019.
- [105] M. Kim, N. Kim, W. Lee, and D. Cho, “Deep learning-aided SCMA,” *IEEE Communications Letters*, vol. 22, no. 4, pp. 720–723, 2018.
- [106] Altera Innovate Asia website, “1st 5G algorithm innovation competition-env1.0-SCMA,” 2015.
- [107] S. Han, T. Xie, C.-L. I, L. Chai, Z. Liu, Y. Yuan, and C. Cui, “Artificial-intelligence-enabled air interface for 6G: Solutions, challenges, and standardization impacts,” *IEEE Communications Magazine*, vol. 58, no. 10, pp. 73–79, 2020.
- [108] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, “Deep reinforcement learning: A brief survey,” *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [109] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, “Learning to communicate with deep multi-agent reinforcement learning,” *CoRR*, vol. abs/1605.06676, 2016.
- [110] M. P. Mota, A. Valcarce, J. M. Gorce, and J. Hoydis, “The emergence of wireless MAC protocols with multi-agent reinforcement learning,” in *2021 IEEE Globecom Workshops (GC Wkshps)*, pp. 1–6, 2021.
- [111] A. Valcarce and J. Hoydis, “Toward joint learning of optimal MAC signaling and wireless channel access,” *IEEE Transactions on Cognitive Communications and Networking*, vol. 7, no. 4, pp. 1233–1243, 2021.
- [112] H. B. Pasandi and T. Nadeem, “Towards a learning-based framework for self-driving design of networking protocols,” *IEEE Access*, vol. 9, pp. 34829–34844, 2021.

## Bibliography

---

- [113] F. Al-Tam, N. Correia, and J. Rodriguez, “Learn to schedule (LEASCH): A deep reinforcement learning approach for radio resource scheduling in the 5G MAC layer,” *IEEE Access*, vol. 8, pp. 108088–108101, 2020.
- [114] D. Abel, D. Arumugam, L. Lehnert, and M. Littman, “State abstractions for lifelong reinforcement learning,” in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, pp. 10–19, PMLR, 10–15 Jul 2018.
- [115] C. Yu, A. Velu, E. Vinitzky, Y. Wang, A. M. Bayen, and Y. Wu, “The surprising effectiveness of MAPPO in cooperative, multi-agent games,” *CoRR*, vol. abs/2103.01955, 2021.
- [116] J. K. Terry, N. Grammel, S. Son, and B. Black, “Parameter sharing for heterogeneous agents in multi-agent reinforcement learning,” *CoRR*, vol. abs/2005.13625, 2020.
- [117] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, “A survey of robot learning from demonstration,” *Robotics and Autonomous Systems*, vol. 57, no. 5, pp. 469–483, 2009.
- [118] L. U. Khan, I. Yaqoob, M. Imran, Z. Han, and C. S. Hong, “6G wireless systems: A vision, architectural elements, and future directions,” *IEEE Access*, vol. 8, pp. 147029–147044, 2020.