UNIVERSITAT POLITÈCNICA DE CATALUNYA

PH.D. DISSERTATION

# AI/ML for multi-technology RAN automation with full and limited infrastructure support

*Author:*
Zoraze ALI

*PhD Advisors:*
Dr. Lorenza Giupponi
and
Dr. Josep Mangues-Bafalluy

*Thesis Tutor:*
Prof. Dr. Miguel Soriano

# Abstract

This thesis studies and proposes solutions to some of the most relevant challenges in the Radio Access Network (RAN) management arising from its evolution beyond 5G and towards 6G. The tackled problems are selected due to the increasing inherent complexity with which these technologies come along in general 5G/6G scenarios, which justifies the need for Artificial Intelligence and Machine Learning (AI/ML) techniques. In particular, with the aim to generalize the applicability of our approach to 5G/6G scenarios, we model a general 5G/6G scenario as a function of the complexity it can present in terms of infrastructure and supported technologies. Based on this modelling, we identify two axes of complexity. On the x-axis, lies the "infrastructure support complexity axis", where the complexity varies based on what support the infrastructure provides, i.e., networks that operate with and without the support of Base Stations (BS) (yet, in the latter case, limited support from a roadside-unit may be offered). On the other hand, on the y-axis, lies the "Technology complexity axis", which captures the complexity variation based on the number of technologies to be operated in a coordinated way. It includes single-technology RAN, which is comprised of only one access technology, e.g., Long Term Evolution (LTE) or New Radio (NR), as well as multi-access technology RAN, which also contains other technologies, such as WiFi. Then, based on these complexity axes, we define three RAN scenarios: 1) infrastructure-based single-technology scenarios, 2) infrastructure-based multi-technology scenarios, and 3) limited infrastructure-based single-technology scenarios. This thesis follows a systematic three-step approach to study these scenarios in depth and identify a set of representative use cases along these axes, which can be addressed with AI/ML solutions to automate RAN management and, at the same time, improve the overall network performance.

In the first step, we focus our study on the infrastructure-based single-technology RAN scenarios. In this category, we identify two use cases, 1) handover management and 2) initial Modulation and Coding Scheme (MCS) selection. Enhancing handover efficiency and optimizing MCS selection are persistent challenges in mobile networks, which further get complicated by the introduction of newer technologies. Focusing on the traditional handover schemes, they present the drawback of considering only the quality of signals from the serving and the target cell to make a handover decision. Also, the initial MCS at the start of the session is usually handled conservatively, i.e., the lowest MCS is assigned to a mobile device that connects to a BS when it first switches on or connects to a new cell after handover. To address these drawbacks, we propose AI/ML solutions that 1) consider the Quality of Experience (QoE) resulting from past handover decisions as the driving principle to select a new target BS to handover and 2) use the experience extracted from the network data to make smarter initial MCS allocations. Specifically, the thesis first presents the AI/ML models designed to address these use cases separately, labeled as single-task solutions. However, we observe that such a technique can present challenges in terms of training cost that increases linearly with respect to the number of use cases to be addressed. Therefore, we propose a generalized AI/ML framework, called a multi-task solution, capable of handling multiple use cases that can operate concurrently at the same or different layers of a mobile protocol stack, e.g., the handover management at

layer 3 and the initial MCS selection at layer 2. To train the proposed models (i.e., single-task and multi-task) and evaluate their performance, we implement a realistic simulation scenario in the *ns-3* simulator that builds an extensive database. Results show that the proposed AI/ML models outperform the 3GPP standardized handover and initial MCS selection approaches by improving the QoE of users resulting from a handover and the throughput obtained upon establishing a new connection with the network.

In the second step, we deal with the complexity of infrastructure-based multi-access technology scenarios to handle the challenges arising from the simultaneous operation of multiple technologies in RAN that share common resources. In particular, we focus on the use case of coexistence in the unlicensed spectrum by studying the channel access technologies known as License Assisted Access (LAA) and LTE- Unlicensed (LTE-U) and their coexistence with WiFi. The main challenge for LAA and LTE-U in these scenarios is that they cannot function without guaranteeing their fair coexistence with WiFi. Therefore, to avoid collisions, the LAA channel access mechanism uses feedback(s) from mobile devices to compute the time it must wait before transmitting to the channel. However, this feedback has an inherent delay, which causes LAA to monopolize the channel. In the thesis, we show that such constraints of the LAA channel access procedure could degrade the performance of neighboring Wi-Fi networks. To solve this problem, we propose an AI/ML-based scheme that learns from experience and infers this feedback under certain channel conditions without delay. Performance evaluation through simulations shows that the proposed scheme provides the best trade-off between the fairness to WiFi and the LAA performance in terms of throughput and latency when compared to the benchmark approaches. We extend our study by proposing a statistical framework to evaluate the fairness offered by LAA and LTE-U when they coexist with WiFi. The comparative analysis confirms that LAA provides better fairness, and LTE-U introduces more collisions.

Finally, in the third step, we direct our attention to the limited infrastructure-based single-technology RAN scenarios. Since, in this case, the RAN does not include BS, operations, such as radio resource selection and scheduling, are uncoordinated and uncontrolled. Mobile devices can communicate directly by selecting the resources autonomously, leading to another level of complexity that needs to be addressed. In this context, this thesis targets the use case of vehicle-to-vehicle communication with limited infrastructure, i.e., besides the vehicles in the scenario, we also consider that there is a roadside unit capable of broadcasting a basic set of information using the 3GPP NR-V2X technology. Nevertheless, to avoid interference among vehicles, the resource selection by a vehicle in NR-V2X is performed by continuously sensing the channel. However, it comes at the cost of higher energy consumption. In contrast, if sensing is not employed with the aim of saving energy, it can result in increased interference. Thus, an energy-performance trade-off arises. To exploit this trade-off, we propose a partial sensing mechanism based on AI/ML to dynamically achieve a balance between performance and energy consumption of V2X users that cannot be obtained by manually configuring parameters of standard sensing procedure.

As for the evaluation methodologies, this thesis also makes an important contribution to the research community by actively contributing and supporting an open-source network simulator. The aim is to foster the reproducibility of our research findings and collaboration in our research community. To this end, all the simulation and data gener-

ation activities are based on the *ns-3* simulator. It offers full-stack, standard-compliant models for major 3GPP cellular technologies (LTE/LAA/LTE-U/NR), which are mainly developed and maintained at CTTC, along with support for WiFi. This guarantees the ability to evaluate all the scenarios mentioned above realistically. In particular, the *ns-3* LTE and LAA models have been extended, and the *5G-LENA* NR-V2X model has been developed in the context of this thesis in collaboration with NIST (part of the U.S. Department of Commerce). We integrate our proposed AI/ML frameworks with these simulation models to conduct end-to-end performance evaluations.

In summary, through extensive evaluations conducted in large-scale representative scenarios that capture the evolving complexity of networks along the identified axes, this thesis successfully demonstrates the potential of AI/ML techniques in addressing the most significant challenges arising in current and future generation networks.

# Resumen

Esta tesis propone soluciones para algunos de los desafíos más relevantes en la gestión de la Red de Acceso por Radio (RAN) que surgen de su evolución más allá del 5G y hacia el 6G. Los problemas abordados han sido seleccionados en función de la complejidad inherente que estas tecnologías conllevan en escenarios genéricos 5G/6G, lo cual justifica la necesidad de técnicas de Inteligencia Artificial y Aprendizaje Automático (AI/ML). Con el objetivo de generalizar la aplicabilidad de nuestro enfoque a escenarios de 5G/6G, modelamos un escenario genérico 5G/6G en función de la complejidad que puede presentar en cuanto a infraestructura y tecnologías soportadas. Basándonos en este modelo, identificamos dos ejes de complejidad. En el eje x se encuentra el "eje de soporte de infraestructura", en el que la complejidad varía según el tipo de soporte ofrecido por la infraestructura desplegada, es decir, redes que operan con y sin el soporte de Estaciones Base (BS) (aunque en este último caso, se puede ofrecer un soporte limitado por parte de una unidad al borde de la calle, o road-side unit, en inglés). Por otro lado, el eje y es el "eje de variedad de tecnologías", que captura la variación de complejidad según el número de tecnologías que operan de manera coordinada. En este sentido, "single-technology" consta de una única tecnología de acceso, como Long Term Evolution (LTE) o New Radio (NR), y "multi-technology" RAN incluye otras tecnologías, como por ejemplo WiFi. A continución, basándonos en estos ejes de complejidad, definimos tres escenarios de RAN: 1) infrastructure-based single-technology, 2) infrastructure-based multi-access technology, and 3) limited infrastructure-based single-technology. Esta tesis sigue un enfoque sistemático de tres pasos para estudiar estos escenarios en profundidad e identificar un conjunto de casos de uso representativos definidos por estos ejes, que pueden abordarse con soluciones de AI/ML para automatizar la gestión de RAN y, al mismo tiempo, mejorar el rendimiento general de la red.

En el primer paso, enfocamos nuestro estudio en los escenarios de RAN infrastructure-based single-technology. En esta categoría, identificamos dos casos de uso: 1) gestión de handover y 2) selección inicial de Modulación y Esquema de Codificación (MCS). Mejorar la eficiencia del handover y optimizar la selección de MCS son desafíos persistentes en las redes móviles, que se vuelven más complicados con la introducción de nuevas tecnologías. Al centrarnos en los esquemas tradicionales de handover, estos presentan la desventaja de considerar solo la calidad de las señales de la celda de servicio y la celda objetivo para tomar una decisión de handover. Además, el MCS inicial en el establecimiento de la sesión se maneja generalmente de manera conservadora, es decir, se asigna el MCS más bajo a un dispositivo móvil que se conecta a una estación base cuando se pone en marcha por primera vez o se conecta a una nueva celda después de un handover. Para abordar estas desventajas, proponemos soluciones de AI/ML que 1) consideran la Calidad de Experiencia (QoE) de decisiones anteriores para seleccionar una nueva estación base y 2) utilizan la experiencia extraída de los datos de la red para realizar asignaciones más inteligentes de MCS iniciales. Específicamente, la tesis presenta primero los modelos de AI/ML diseñados para abordar estos casos de uso por separado, etiquetados como soluciones "single-task". Sin embargo, observamos que esta técnica puede presentar desventajas en términos de coste de entrenamiento, que aumenta linealmente con respecto

al número de casos de uso que se deben abordar. Por lo tanto, proponemos un marco generalizado de AI/ML, llamado solución "multi-task", capaz de gestionar múltiples casos de uso que pueden funcionar simultáneamente en la misma o en diferentes capas de la pila de protocolos móvil, por ejemplo, la gestión de handover en la capa 3 y la selección inicial de MCS en la capa 2. Para entrenar los modelos propuestos (es decir, de una sola tarea y de múltiples tareas) y evaluar su rendimiento, implementamos un escenario de simulación realista en el simulador *ns-3* que crea una base de datos extensa. Los resultados muestran que los modelos propuestos de AI/ML superan las soluciones/técnicas estandarizadas de handover y selección inicial de MCS del 3GPP al mejorar la QoE resultante de los usuarios de un handover y el rendimiento obtenido al establecer una nueva conexión con la red.

En el segundo paso, abordamos la complejidad de los escenarios "infrastructure-based multi-access technology" para tratar los desafíos que surgen de la operación simultánea de múltiples tecnologías RAN que comparten recursos. En particular, nos centramos en el caso de uso de coexistencia en el espectro no licenciado mediante el estudio de las tecnologías de acceso conocidas como License Assisted Access (LAA) y LTE-Unlicensed (LTE-U) y su coexistencia con WiFi. El principal desafío para LAA y LTE-U en estos escenarios es garantizar su coexistencia adecuada con WiFi. Por lo tanto, para evitar colisiones, el mecanismo de acceso al canal de LAA utiliza feedback de los dispositivos móviles para calcular el tiempo que debe esperar antes de transmitir al canal. Sin embargo, estos feedback(s) tienen un retraso inherente, lo que provoca que LAA monopolice el canal. En la tesis, mostramos que estas limitaciones del procedimiento de acceso al canal de LAA podrían degradar el rendimiento de las redes WiFi próximas. Para resolver este problema, proponemos un esquema basado en AI/ML que aprende de la experiencia e infiere estos feedback(s) en ciertas condiciones del canal sin retraso. La evaluación del rendimiento a través de simulaciones muestra que el esquema propuesto proporciona el mejor equilibrio entre la equidad para WiFi y el rendimiento de LAA en términos de throughput y latencia en comparación con los enfoques del estado de arte. Extendemos nuestro estudio proponiendo un marco estadístico para evaluar la equidad ofrecida por LAA y LTE-U cuando coexisten con WiFi. El análisis comparativo confirma que LAA proporciona una mejor equidad y que LTE-U introduce más colisiones.

Finalmente, en el tercer paso, nos centramos en los escenarios de RAN "limited infrastructure-based single-technology". En este caso, la RAN no incluye estaciones bases y las operaciones, p.ej., la selección y la repartición de recursos de radio, no son coordinadas ni controladas. Los dispositivos móviles pueden comunicarse directamente seleccionando los recursos de manera autónoma, lo que conlleva otro nivel de complejidad que debe abordarse. En este contexto, esta tesis se centra en el caso de uso de comunicación de vehículo a vehículo con soporte limitado por parte de la infraestructura. Es decir, además de los vehículos en el escenario, también consideramos que hay una unidad en el borde de la carretera (road-side unit) capaz de transmitir un conjunto básico de información utilizando la tecnología NR-V2X del 3GPP. Sin embargo, para evitar interferencias entre vehículos, la selección de recursos por parte de un vehículo en NR-V2X se realiza mediante la detección (sensing) continua del canal. Sin embargo, esto conlleva un mayor consumo de energía. En contraste, si no se emplea la detección con el objetivo de ahorrar energía, puede haber un aumento de la interferencia. Por lo tanto, surge un compromiso entre energía y rendimiento. Para aprovechar este compromiso, proponemos un mecanismo de detección parcial basado en AI/ML para lograr dinámicamente un

equilibrio entre el rendimiento y el consumo de energía de los usuarios V2X que no se puede obtener mediante la configuración manual de los parámetros del procedimiento de detección estándar.

En cuanto a las metodologías de evaluación, esta tesis también hace una importante contribución a la comunidad de investigación aportando un simulador de red de código abierto. El objetivo es fomentar la reproducibilidad de nuestros hallazgos de investigación y la colaboración en nuestra comunidad de investigación. Con este fin, todas las actividades de simulación y generación de datos se basan en el simulador *ns-3*. Este simulador ofrece modelos "full-stack" y compatibles con estándares para las principales tecnologías celulares 3GPP (LTE/LAA/LTE-U/NR) así como el soporte para WiFi, que se desarrollan y se mantienen principalmente en el CTTC. Esto garantiza la capacidad de evaluar todos los escenarios mencionados anteriormente de manera realista. En particular, los modelos de LTE y LAA de *ns-3* se han ampliado y se ha desarrollado el modelo *5G-LENA* NR-V2X en el contexto de esta tesis, en colaboración con NIST (Departamento de Comercio de los Estados Unidos). Integramos nuestros entornos propuestos de AI/ML con estos modelos de simulación para realizar evaluaciones de rendimiento extremo a extremo.

En resumen, a través de evaluaciones exhaustivas realizadas en escenarios representativos a gran escala que capturan la complejidad de la evolución de las redes a lo largo de los ejes de infraestructura y tecnología, esta tesis demuestra el potencial de las soluciones de AI/ML para automatizar y mejorar la gestión de RAN más allá del 5G.

x

# Acknowledgements

It feels like I dived into the ocean named Ph.D., and while holding my breath, I collected some of the most precious, valuable, and rewarding moments of my life. And, of course, I couldn't hold that breath without the help of the people in my life who acted as oxygen for me during this long dive.

First and foremost, I am immensely grateful to my supervisors, Dr.Lorenza Giupponi and Dr.Josep Mangues Bafalluy, for their patience in guiding me at every step and for believing in me that I can swim in this Ph.D. ocean. Their expertise and constant support have shaped this research and pushed me to achieve my best.

To Nicola Baldo, thank you so much for all your guidance and support in the initial days of my Ph.D., which laid the foundation for my research journey.

Biljana Bojovic, thanking you would be an understatement. Your support as a friend has been invaluable throughout my journey.

To my amazing CTTC team members, who are like family to me and were always there for me. Thank you, Lorenza, Biljana, Sandra, Katerina, and Natale.

I am and will always be grateful to Katerina and Karrar Rizvi for their support and being the source of strength in the final days of my Ph.D.

I am incredibly thankful to Dr.Jessica Moysen and Dr.Istiak Hussain for reviewing the thesis despite their busy schedule.

Indeed, mere words cannot adequately convey the depth of my appreciation to God for the gift of the family He has blessed me. I would like to thank my wife, Sobia, for being the lighthouse of my life, being the best wife, and putting her trust in me. I want to thank my grandpa, grandma, siblings, and uncles for their support in getting me where I am today.

Finally, I dedicate this thesis to my parents, Nusrat Parveen and Meherban Ali, for being the most amazing, supportive, and loving parents. It is you who taught me to work hard and aim for excellence. I will forever be grateful for your sacrifices and the values you instilled in me. This thesis is a testament to your unstoppable dedication and your profound impact on my life. Love you and miss you, papa!

<div align="right">Zoraze Ali</div>

# Contents

# List of Figures

# List of Tables

# Acronyms

$P_{\text{rsvp}}$     Resource Reservation Period
3GPP     3rd Generation Partnership Project
5G     Fifth-generation wireless
5GAA     5G Automotive Association
6G     Sixth-generation wireless
AE     Auto Encoder
AI     Artificial Intelligence
AP     Access Point
ARPU     Average Revenue per User
BS     Base Station
BSR     Buffer Status Report
BWP     Bandwidth Part
C-V2X     Cellular V2X
CAPEX     Capital Expenditures
CCA     Clear Channel Assessment
CDF     Cumulative Distribution Function
CP     Cyclic Prefix
CQI     Channel Quality Indicator
CSAT     Carrier Sense Adaptive Transmission
CSI     Channel State Information
CW     Contention Window
D2D     Device-to-Device
DB     DataBase
DC     Duty Cycle
DL     Downlink
DMRS     Demodulation Reference Signals
DNN     Deep Neural Network
DSRC     Dedicated Short Range Communications
ECDF     Empirical Cumulative Distribution Function
eICIC     enhanced Inter-Cell Interference Coordination
eNB     Evolved Node B
EPC     Evolved Packet Core
EPS     Evolved Packet System
FDD     Frequency Division Duplex
FDMA     Frequency-Division Multiple Access
FFNN     Feed-forward Neural Network

| | |
|---|---|
| FTP | File Transfer Protocol |
| gNB | next-Generation Node B |
| GPU | Graphical Processing Unit |
| HARQ | Hybrid Automatic Repeat Request |
| HO | HandOver |
| HTTP | Hypertext Transfer Protocol |
| ICIC | Inter-Cell Interference Coordination |
| IE | Information Element |
| IEEE | Institute of Electrical and Electronics Engineers |
| IP | Internet Protocol |
| ITS | Intelligent Transport System |
| KPI | Key Performance Indicator |
| KQI | Key Quality Indicator |
| KS-test | Kolmogorov-Smirnov test |
| LAA | Licensed-Assisted Access |
| LBT | Listen-Before-Talk |
| LC | Logical Channel |
| LCID | Logical Channel Identifier |
| LDPC | Low Density Parity Check |
| LR-WPAN | Low-Rate Wireless Personal Area Network |
| LSTM | Long Short Term Memory |
| LTE | Long Term Evolution |
| LTE-U | LTE Unlicensed |
| MAC | Medium Access Control |
| MCPTT | Mission Critical Push To Talk |
| MCS | Modulation and Coding Scheme |
| MDP | Markov Decision Process |
| MDT | Minimization of Drive Test |
| MIMO | Multiple Input Multiple Output |
| ML | Machine Learning |
| MLP | Multi-Layer Perceptron |
| mmWave | Millimeter Wave |
| MNO | Mobile Network Operator |
| MOS | Mean Opinion Score |
| MSE | Mean Square Error |
| MTL | Multi-Task Learning |
| NACK | Negative Acknowledgement |
| NAS | Non-Access Stratum |
| NDI | New Data Indicator |
| NG-RAN | Next Generation RAN |
| NG-SON | Next Generation Self-Organizing Network |
| NN | Neural Network |
| NOMA | Non-Orthogonal Multiple Access |
| NR | New Radio |
| NR-U | New Radio-based access to Unlicensed spectrum |

| | |
|---|---|
| NRMSE | Normalize Root Mean Square Error |
| NSA | Non-Standalone |
| NTN | Non-Terrestrial Networks |
| O-RAN | Open-RAN |
| OFDM | Orthogonal Frequency Division Modulation |
| OPEX | Operational Expenditures |
| PCAP | Packet CAPture |
| PDCP | Packet Data Convergence Protocol |
| PDSCH | Physical Downlink Shared Channel |
| PDU | Packet Data Unit |
| PGW | Packet data network GateWay |
| PHY | Physical Layer |
| PIR | Packet Inter-reception Delay |
| POMDP | Partially Observable Markov Decision Process |
| ProSe | Proximity Services |
| PRR | Packet Reception Ratio |
| PSBCH | Physical Sidelink Broadcast Channel |
| PSC | Public Safety Communications |
| PSCCH | Physical Sidelink Control Channel |
| PSFCH | Physical Sidelink Feedback Channel |
| PSSCH | Physical Sidelink Shared Channel |
| PUSCH | Physical Uplink Shared Channel |
| QoE | Quality of Experience |
| QoS | Quality of Service |
| RAN | Radio Access Network |
| RAT | Radio Access Technology |
| RB | Resource Block |
| RE | Resource Element |
| REM | Radio Environment Map |
| RIC | Radio Intelligent Controller |
| RLC | LTE Radio Link Control |
| RLF | Radio Link Failure |
| RNN | Recurrent Neural Network |
| RRC | Radio Resource Control |
| RSRP | Reference Signal Receive Power |
| RSRQ | Reference Signal Received Quality |
| RSU | Road Side Unit |
| RV | Redundancy Version |
| SCI | Sidelink Control Information |
| SCS | sub-carrier spacing |
| SDL | Supplemental DownLink |
| SGW | Service GateWay |
| SI | Study Item |
| SINR | Signal to Interference plus Noise Ratio |
| SLA | Service Level Agreement |
| SLRRC | Sidelink Resource Reselection Counter |

| | |
|---|---|
| SNR | Signal to Noise Ratio |
| SON | Self-Organizing Network |
| SOTA | State-of-the-Art |
| SPS | Semi-Persistent Scheduling |
| SQL | Structured Query Language |
| SRS | Sounding Reference Signal |
| STA | Station |
| SUMO | Simulation of Urban MObility |
| SVM | Support Vector Machine |
| TBS | Transport Block Size |
| TCP | Transmission Control Protocol |
| TDD | Time Division Duplex |
| TDMA | Time-Division Multiple Access |
| TFT | Traffic Flow Template |
| TR | Technical Report |
| TS | Technical Specification |
| TTI | Transmission Time Interval |
| TxOP | Transmission Opportunity |
| UAV | Unmanned Aerial Vehicle |
| UE | User Equipment |
| UL | Uplink |
| UM | Unacknowledged Mode |
| URLLC | Ultra-Reliable and Low-Latency Communications |
| UTRAN | UMTS Terrestrial Radio Access |
| V2I | Vehicle-to-Infrastructure |
| V2N | Vehicle-to-Network |
| V2P | Vehicle-to-Pedestrian |
| V2V | Vehicle-to-Vehicle |
| V2X | Vehicular-to-everything |
| WAVE | Wireless Access in Vehicular Environments |
| WCDMA | Wideband Code Division Multiple Access |
| WI | Work Item |
| WiGig | Wireless Gigabit |
| WiMAX | Worldwide Interoperability for Microwave Access |

# Chapter 1

# Introduction

## 1.1 Motivation

Mobile communications have experienced during the last decades an incredible evolution. Since its inception, mobile device connections have surpassed the number of people worldwide, making it the fastest growing technology ever [1]. Despite that, it is well known that the revenue generated by Mobile Network Operators (MNOs) per user (Average Revenue per User (ARPU)) has been steadily decreasing for the last decade. The Capital Expenditures (CAPEX) of Fifth-generation wireless (5G) networks are still not completely clear and include: 1) more spectrum, with expensive auction fees, 2) deployment of new antennas and equipment upgrade, 3) large-scale small-cell deployments, to pursue the Millimeter Wave (mmWave) vision. Therefore, the reduction of Operational Expenditures (OPEX) is fundamental to the evolution of Beyond 5G mobile communication systems. Another interesting data is that 70 % of the total cost involving deployment, optimization, and operation of a network comes from the Radio Access Network (RAN) segment [2]. It is the reason that there is a significant interest in improving the efficiency of the RAN management, which consequently reduces its OPEX. Currently, there are two main trends to achieve these objectives.

On the one hand, automation and self-organization of the RAN have become two fundamental ingredients for optimal resource utilization and management. It has been almost a decade since when Self-Organizing Network (SON) was defined and introduced as a feature of Long Term Evolution (LTE), in 3rd Generation Partnership Project (3GPP) Release 8 [3]. Since then, it has been evolving through the releases and into

the concept of Next Generation Self-Organizing Network (NG-SON) for 5G and Sixth-generation wireless (6G) networks [4]. These future cellular networks are characterized by highly complex, dense, and heterogeneous deployments to increase network coverage and capacity. Besides traditional sub-6 GHz and licensed bands, the access can span a wide bandwidth range, including mmWave and unlicensed spectrum. The high diversity of mobile devices and new applications further complicates the network architecture and its management. In this context, mobile networks generate a massive amount of measurements, control, and management information during their normal operation [5] [6]. This huge amount of information could be efficiently utilized to address future mobile network management challenges. With mobile networks being rolled out nowadays, it is already necessary to feature SON solutions (e.g., for Inter-Cell Interference Coordination (among others)) to operate properly. While narrowly focused SON solutions can still be designed with traditional engineering approaches (design based on analytical modeling of the problem, followed by successive empirical refinements), this strategy is not feasible anymore as more parameters, protocols, and optimization objectives are taken into the picture. To address this issue, several techniques recently developed in the fields of ML appear promising to leverage the huge amount of information to optimize the network [7] [8]. The key point is that this massive amount of data is overwhelming for traditional engineering practices; hence only ML approaches are expected to be able to exploit them successfully for network management and optimization purposes. In this line, the recent evolution in computational capabilities, e.g., the availability of the Graphical Processing Units (GPUs) to train neural networks at a relatively faster pace, has allowed us to take advantage of ML and novel deep learning solutions to tackle multiple problems in different disciplines. In 5G and its evolution, the possibilities now available for ML and novel deep learning implementations are infinite and pave the way to an evolved vision of NG-SON to be able to address end-to-end solutions.

On the other hand, the evolution towards 6G networks calls for further architectural transformations required to support service heterogeneity, coordination of multi-connectivity, on-demand service deployment, and network automation. Therefore, the telecommunication industry, through different consortiums, e.g., 3GPP, and the Open-RAN (O-RAN) alliance, have acknowledged Artificial Intelligence (AI)/ML as one of the important components of future mobile networks [9–11]. In 3GPP, the standardization laying the ground for AI/ML studies in RAN started in 5G Phase 1, in Release 15, which continued in 3GPP releases 16, 17 and 18, targeting 6G and beyond [12]. The studies in the 3GPP RAN groups, i.e., RAN1[1] and RAN3[2], aim to standardize the use of AI/ML in specific use cases [13,14]. These studies cover network energy saving, load balancing, mobility optimization, Channel State Information (CSI) feedback enhancement, beam management, and position accuracy enhancement, as initial use cases. While 3GPP standards are being adopted to deploy early 5G commercial networks [15], in Feb. 2018, a group of mobile network operators founded the O-RAN Alliance to further enhance RAN performance through virtualized network elements, openness, and intelligence. Openness aims to eliminate proprietary hardware and software implementations by establishing open

---

[1]The RAN1 group is responsible for the standardization of the physical layer procedures of the radio Interfaces for User Equipment (UE), Evolved UMTS Terrestrial Radio Access (UTRAN), Next Generation RAN (NG-RAN), and beyond.

[2]The RAN3 group is responsible for the standardization of the UTRAN/Evolved UTRAN/NG-RAN architecture and the protocols needed for network interfaces.

standard interfaces, reducing operating costs. Intelligence is necessary for deploying, optimizing, and operating Beyond 5G networks. O-RAN introduces new Radio Intelligent Controller (RIC) modules and enables them with AI/ML features to enhance traditional network functions with intelligence. In one of its defining white papers [11], different use cases are proposed, like traffic steering, Quality of Experience (QoE) optimization, Quality of Service (QoS) based resource optimization, RAN Slice Service Level Agreement (SLA) assurance, and context-based dynamic HO management for Vehicular-to-everything (V2X) [16] [17].

Both these trends and visions converge to the already widely agreed need of AI/ML as fundamental ingredients of Beyond 5G and 6G networks. Therefore, solutions based on AI/ML have been lately intensively investigated in the literature of mobile communications to solve a wide range of problems in various domains, including the RAN [18]. For example, research has focused on augmenting the SON functionality with AI/ML targeting various use cases such as, resource optimization [19, 20], mobility management [21–25], and load balancing [26, 27]. In the context of O-RAN, studies have aimed on improving network's energy efficiency [28], QoS [29], QoE [30], traffic steering capability [31], and radio resource scheduling [32], among others, through AI/ML based frameworks.

In this line, this thesis targets to study and exploit the possibilities of AI/ML for improving the RAN operation. More details related to the identified RAN challenges and applicability of AI/ML in RAN are given in the following section.

## 1.2 Problem statement

RAN operation presents many challenges due to its high complexity. For instance, it requires the ability to continuously adapt to the environment's ever-changing conditions regarding propagation, diverse users' needs, system load, high mobility, among others. Taking one step further, the number of tasks that a Beyond 5G RAN has to execute is vast and includes all traditional, as well as future envisioned SON use cases [9]. Specifically, the heterogeneity of the RAN is not anymore only limited to the dense deployment of different types of base stations (e.g., high, mid, and low-powered cells), which in its own presents many issues, as mentioned before [33–35]. Future mobile RAN ecosystem is evolving in many dimensions. Some examples include the support of co-existence in the unlicensed spectrum, D2D, V2X, and now recently, Non-Terrestrial Networks (NTN) communications. Moreover, the 3GPP standard TS 22.261 already contemplates the idea of 5G and beyond cellular systems being capable of simultaneously supporting multiple access technologies, e.g., LTE, New Radio (NR), Licensed-Assisted Access (LAA) for one or more services active on a mobile station [36]. On the one hand, this evolution of RAN is the enabler for mobile networks to accelerate towards 6G [37]. On the other hand, it makes the management of RAN very complex. To handle such complexity, the academia and the industry have already foreseen the benefits of applying AI/ML techniques to improve the RAN performance, as presented in the previous section. In fact, AI/ML can find its applicability in various use cases, while its efficiency can further be improved.

**Figure 1.1:** A visual representation of scenarios based on the two identified complexity axes, where I = Infrastructure, LI = Limited Infrastructure, ST = Single access Technology, MT = Multi access Technology

With this in mind, this thesis leverages AI/ML to improve the RAN management and network performance while reducing its OPEX. Our initial goal is to study different RAN architectures comprised of multiple access technologies and to find unexplored use cases where AI/ML can be applied to improve network performance. In particular, this thesis identifies two complexity axes to target three representative RAN scenarios, as shown in Fig. 1.1. Collectively, these three scenarios encompass various challenges that may arise in a RAN. Among these challenges, we carefully choose a specific issue and demonstrate that AI/ML can effectively tackle it. Furthermore, it would be valuable to explore solutions enabling simultaneous learning and training of new use cases, as this would address the feasibility of implementing the AI/ML vision. By utilizing these RAN scenario categories, which are explained in detail in Sections 1.2.1, 1.2.2, and 1.2.3, this thesis aims to answer the following High-level Question (HQ):

**HQ: How can AI/ML be used to automate the increasing RAN management complexity along two axes: 1. infrastructure- and limited infrastructure-based RAN scenarios and 2. single- and multi-access technology RAN scenarios?**

In order to answer this high-level question, we need to study each of the above three identified RAN scenarios separately and target specific use cases within them. What follows is a brief description of these use cases and the more specific research questions that are embedded in the main question described above that this thesis aims to answer.

## 1.2.1  Infrastructure-based single-technology RAN

This scenario corresponds to the infrastructure-based traditional RAN that is comprised of base stations and devices operating in licensed spectrum, using a single access technology, e.g., LTE or NR. More than a decade ago, when the LTE system requirements were finalized, the "seamless mobility" was among one of the key features that mobile networks has to offer [38]. The current deployments of the mobile networks are carried out with the joint installation of densely deployed heterogenous cells. This by itself already poses many challenges to the mobility management. The future mobile networks, such as 5G and 6G, target to serve even higher number of users by further densifying the network. In these next-generation networks, when there is high mobility and the coverage gets impacted, the UEs perform frequent handovers to maintain the connection. The user´s QoE during mobility gets affected and might get highly sensitive to subtle movements in the coverage area. The routine for handover is a challenge itself as it is a trial-and-error based mechanism with chances of being impacted by unpredictable behaviors caused by packet drop, longer delay because of speedy movements, and inconsistency and randomness of the transmission environment, causing glitches in the feedback procedure and thus leading to failures. Thus, the move towards beyond 5G networks would not only increase the complexity of the RAN's topology but also complicates its management, imposing the need for more advanced techniques to handle mobility. Moreover, regardless of the advancements introduced by new mobile generations, the issue of efficiently managing mobility/handover persists, and in some cases, even gets worse. On this matter, this thesis focuses on extensively studying the existing handover schemes, specifically the target cell selection mechanism to choose the next base station to connect. The standard approach for such functionality is that it selects the next cell solely based on the strongest signal strength before the handover. Undoubtedly, this solution for target cell selection is simple and easy to implement. However, the problem that arises is that this solution does not consider the QoE of the user after the handover, which can get impacted due to complex radio environments that may exist in current and future mobile networks. In this line, by selecting mobility management as the first use case under the infrastructure-based single-technology RAN scenario, we focus on answering the following Research Question (RQ):

**RQ1: How to use AI/ML in mobility management to achieve better QoE?**

Following the introduction of the HO use case, which is the layer-3 problem and serves as a highly representative element for existing and future mobile networks, our attention shifts to different yet related issues. In the existing networks, when a UE is first turned on in a specific cell or after a handover, the new base station has no information about the radio conditions of that UE. Therefore, the base station takes a conservative approach of selecting the lowest possible initial Modulation and Coding Scheme (MCS) for transmission, which limits the data rate for the UE, impacting its throughput until the base station receives the channel status report, which is used to update the MCS. Given the importance of MCS selection on spectral efficiency and user experience, it becomes crucial to enhance this functionality by leveraging the network's experience to determine an appropriate initial MCS for a specific network location. In this respect, the thesis aims to answer the following research question:

**RQ2: How to use AI/ML to select the initial MCS for newly connected mobile devices to achieve better throughput?**

After formulating the above two research questions that tackle the layer-3 HO management and the layer-2 initial MCS use cases separately, we turn towards exploring the possibilities of finding a generalized AI/ML framework that can be used to handle multiple RAN use cases. The reason is that in AI/ML, we typically care about optimizing for a particular metric. To do this, we generally train a model to perform a desired task, or as in our case, the use case. We then fine-tune these models until the desired accuracy is obtained. By doing so independently for all the different RAN use cases, we may be ignoring information that might help us do even better on the metric of interest. Specifically, this useful information may come from different and related tasks/RAN use cases. By sharing relevant information between these tasks, for example, using a wide feature space, which is not reduced only to a specific problem to solve, we can enable our model to generalize better on our original task. Moreover, it can also help to reduce the complexity in terms of the training cost of the AI/ML models, which can increase linearly with the number of use cases a RAN has to handle. In this context, this thesis aims to answer the following question:

**RQ3: How to generalize an AI/ML solution to address diverse RAN use cases?**

## 1.2.2 Infrastructure-based multi-access technology RAN

Besides the approach of densifying the RAN with many cells to increase network capacity, the mobile networks have evolved to also operate in unlicensed frequency bands, overcoming the scarcity of expensive licensed spectrum. In this second scenario, we focus on the use case of coexistence in unlicensed spectrum by highlighting the complexities arising from the multi-access technology scenarios (i.e., the second complexity axis). In these scenarios, base stations and mobile devices operate in both licensed and unlicensed spectrum using the access technologies such as LAA or LTE Unlicensed (LTE-U) (presented in Chapter 2). The complexity in the management of this type of mobile network originates due to the fact that it has to coexist with other wireless technologies in the unlicensed bands, e.g., widely used WiFi networks, which is not the case when using licensed spectrum. Therefore, it is of utmost importance for LAA or LTE-U devices to access the channel fairly so that they do not hamper the performance of other WiFi devices in terms of throughput and latency. Between LAA and LTE-U, it is believed that the LAA channel access mechanism could provide better fairness to coexisting WiFi networks. However, its level of fairness depends on the configuration of the parameters related to its channel access mechanism. To avoid collisions, the LAA channel access procedure uses HARQ feedback(s) from mobile devices to compute the time it must wait before transmitting to the channel. However, in LAA, the HARQ feedback suffers a delay due to the inherent latencies in the mobile protocol stack. Additionally, the feedback from multiple non-co-located users in a subframe is combined to increase the probability of decoding a packet, which is different than WiFi. This thesis shows that such characteristics of the LAA channel access procedure can lead to an unfair coexistence

with WiFi. In this line, this thesis aims to address these drawbacks by answering the following research question:

**RQ4: How to use AI/ML to guarantee fair coexistence of LAA in the unlicensed spectrum?**

## 1.2.3 Limited infrastructure-based single-technology RAN

With the above two scenarios, this thesis targets the use cases originating from a RAN network comprised of single or multiple access technologies. However, in such scenarios, the communication between the devices and the network is always coordinated and controlled by a base station. Based on this, we categorized them as infrastructure/full-infrastructure-based scenarios. Nevertheless, as envisioned today, the future 5G and 6G networks will operate when needed or instructed by the network, without a base station. This functionality is particularly useful when providing network support is not possible, e.g., in case of natural disasters or where the latency incurred by the network can be decreased by enabling direct communication between mobile devices. The 3GPP technologies that enable such features include the D2D, Cellular V2X (C-V2X) (both based on LTE), and recently the NR-V2X. This thesis focuses on the NR-V2X technology to investigate the complexities in a RAN that has limited infrastructure support. Let us notice that, the term "limited" implies that the scenario under study does not involve the use of base stations; however, it does include a Road Side Unit (RSU)[3] capable of broadcasting a minimum set of information, such as, the total number of vehicles in a specific area, that can help vehicles to have a broader view of the network. Nevertheless, in such a scenario, radio resources for transmission are autonomously selected by the vehicles, which can result in an increased level of interference among the vehicles within range of each other. In turn, this non-coordinated transmissions lead to an elevated level of complexity in the network that does not exist in the previously presented infrastructure-based scenarios. To minimize interference, the 3GPP standard defines the sensing procedure in which the vehicles need to continuously sense the surrounding transmissions. Despite the fact that this type of resource selection reduces the interference among vehicles, it consumes more energy. Moreover, depending on the scenario that could vary in terms of number of vehicles or congestion level in the network, how much sensing a vehicle must perform is an open discussion since there is a trade-off between the energy consumption and the performance of a V2X UE. As such, a static configuration of the sensing parameters can be suboptimal. With this in mind, this thesis aims to answer the following research question to explore the possibilities of leveraging AI/ML in such scenarios.

**RQ5: How to achieve a balance between the energy consumption and the performance of a NR V2X UE in limited infrastructure-based scenarios using AI/ML?**

---

[3]An RSU is a special node that is located at the roadside and is capable of wirelessly communicating safety warnings and traffic information to passing vehicles.

# 1.3 Research approach

Using AI/ML techniques across different disciplines, including wireless communication, has yielded promising results, as highlighted at the beginning of this chapter. However, its performance mainly relies upon training datasets that should be sufficiently large and comprehensive [39]. Thus, obtaining such data is of paramount importance. While the fields, such as computer vision, can leverage big data datasets, e.g., MNIST handwriting dataset [40] or ImageNet [41], in mobile communications it is not the case, even though there are some exceptions [42]. The reasons include 1) policies concerning user privacy and security that make it difficult for MNO to release user´s data, 2) the openly available data is aggregated, lacking the information about, e.g., the traffic type, the type of technology, the protocol layer, etc. To circumvent these limitations, the research community has turned towards using testbeds to generate datasets in close-to-real network setups [43,44]. These setups use commercial off-the-shelf equipment and devices to create a network closely mimicking real network implementation. However, they also have a few drawbacks: 1) It isn't easy to test complex RAN scenarios such as the one presented in chapters 3 and 4 of this thesis, 2) the lack of support for the new technologies in a timely manner, which is very important in research, e.g., in the context of research work presented in chapter 5 to 8 that focuses on LAA, LTE-U (chapter 5 and 6) and NR-V2X (chapter 7 and 8) technologies. Given these constraints, in this thesis, we decided to opt for a simulation framework that can primarily fulfill our high-level requirements of 1) generating synthetic datasets for the training and testing of the proposed machine learning models, 2) implementing innovative and complex RAN scenarios, and 3) validating the performance of the proposed solutions.

In academia and the industry, AI/ML-based data-driven research on wireless communication mainly uses two types of simulation tools, 1) Link-level and 2) System-level. The question is: which one should be used to achieve the research objective? Answering this question based on the above high-level requirements is not very straightforward. Therefore, one needs to narrow down the requirements further based on the overall research objective. On the one hand, the link-level simulators are developed to model detailed physical layer functionalities (e.g., bit scrambling, precoder, IFFT/FFT, channel equalizer, etc.) to emulate the physical layer of real wireless networks. Because of such detailed implementation, these simulators may be computationally demanding. Therefore, it might not be suitable for multi-cell and multi-RAT scenarios, such as the ones studied in this thesis. Moreover, as mentioned in the Subsection. 1.2, the RAN use cases we plan to focus on encompass different protocol layers, i.e., layers 1, 2, and 3 of mobile networks. In this context, a link-level simulator cannot fulfil our research objective.

On the other hand, the system-level simulators may model all the protocol layers, multi-cell, multi-RAT scenarios, and other parts of the network, e.g., the core network. These characteristics seem promising to achieve the goal of this thesis. Therefore, we have selected a system-level simulation framework as our primary tool. However, the completeness offered by a system-level simulator comes at the cost of a certain level of abstraction [45]. For example, the physical layer abstraction is achieved using Link-to-System mapping[4]. A high-level abstraction may cause the results to deviate too much

---

[4]A Link-to-System mapping is a technique to run simulations in a timely manner by accurately predicting the performance of a link in a computationally efficient way [46].

from the experimental results [47]. Therefore, extra attention is to be paid to the level of abstraction used to model such a simulator. Other essential aspects which increase the authenticity and acceptability of the simulator in the research community are its 1) standard-compliant and validated simulation models, 2) simplicity in terms of extending and implementing models, 3) extensive documentation, and most importantly, 4) open-source availability to facilitate the reproducibility of the results, which is difficult to achieve using proprietary simulators.

Currently, a handful of open-source system-level simulators are used in research, e.g., *ns-2*, *ns-3*, *SimuLTE* based on *OMNeT++*, *JiST*, and *SimPy*. Among these simulators, to the best of the author's knowledge, *ns-3* is the most cited, hence more trusted, system-level simulator that contains all the essential aspects listed above [45]. More importantly, it provides all the means to achieve the research objective of this thesis. Therefore, we consider it the primary simulation tool in the thesis to implement case-specific scenarios, extend existing simulation models, implement new models, generate synthetic data to train and validate proposed AI/ML models and perform end-to-end full protocol stack performance evaluations. What follows next is a brief overview of the *ns-3* simulator, highlighting its key features, and the approach used in this thesis to link *ns-3* simulator with the AI/ML framework.

## 1.3.1 *ns-3* overview

*ns-3* is a discrete event-based system-level simulator that has been openly available under the GNU General Purpose License, version 2 (GPLv2) since 2006. The term "discrete event" implies that the state of the simulation can only change upon an occurrence of an "event", at a particular time [47]. For example, an event can be a start/stop of an application that transmits packets, an update of a node's position that moves with a certain speed, etc.

From a software organization standpoint, the *ns-3* simulator is divided into C++ libraries called "modules" [48]. These modules are built to simulate high-fidelity models of a complete network protocol stack, i.e., application, transport, network, and specific MAC, and PHY layers implementation of different communication technologies, such as Ethernet, Wi-Fi, Low-Rate Wireless Personal Area Network (LR-WPAN), Wireless Access in Vehicular Environments (WAVE)[5], Worldwide Interoperability for Microwave Access (WiMAX), LTE, etc. Additionally, there are modules that can be used to model other important aspects of a simulation, e.g., the mobility of a node, wireless channel propagation model, antenna model, placement of buildings (to simulate indoor/outdoor scenarios), and the extraction of useful information during and after the simulation using commonly used formats, e.g., Packet CAPture (PCAP)[6], text, and Structured Query Language (SQL) database. At the time of writing this thesis, the *ns-3* simulator has 44 modules that provide comprehensive documentation on their modeling and usage. Details of these modules are out of the scope of this thesis. Therefore, we refer the interested

---

[5]WAVE is the technology for wireless access in vehicular environments, which is also known as IEEE 802.11/p [49].

[6]PCAP is a standardized format used by network packet analysis tools such as tcpdump and WireShark [47].

reader to *ns-3* manual, tutorial, and model library in [48], [50], and, [51]. However, in the following, we do provide a brief overview of the LTE and the NR modules, which are extensively used and extended for the studies conducted in this thesis. Lastly, in Table 1.1, we summarize 1) the key features supported by these modules and their extensions and 2) the contributions that the author of this thesis has made to these modules.

The LTE module of *ns-3*, also known as LENA [52], is a very commonly used LTE network simulation platform that allows the simulation of end-to-end LTE heterogeneous networks. This module is based on an industrial API (the small cell forum LTE MAC Scheduler interface specification). Because of this, the protocol stack is very similar to actual protocol implementations found in commercial products. This feature is important in a data-driven AI/ML research work, like the one presented throughout this thesis. As mentioned above, it is not easy for mobile network operators to openly provide network traces due to several reasons, including network security and users´ privacy issues [42]. In this line, the *ns-3* simulator, through its "Tracing System", allows extracting useful information from any protocol stack layer. One such example is presented in Chapter 4 of this thesis, in which we extract 84 measurements from the overall LTE protocol stack. The LTE module of *ns-3* was first released in 2013. It resulted from a collaboration between a small-cell vendor, Ubiquisys Ltd. (now part of Cisco), and CTTC. Over time, in the research, this module has been validated against a testbed and through calibration studies in [53] and [54]. In [53], the authors conclude that an emulation setup based on *ns-3* LTE module can achieve comparable performance in terms of the Mean-Opinion-Score (MOS) and latency for the voice application as an experimental testbed, consisting of real LTE equipment over a range of Signal to Noise Ratios (SNRs). Finally, the authors in [54] calibrated the LTE module using the parameters tested in 3GPP for urban and rural macro cell scenarios [55]. The results show that the LTE module obtains similar Signal to Interference plus Noise Ratio (SINR) distributions and users' perceived throughput as of the other 17 industrial simulators in 3GPP. Moreover, this module has been under continuous development, either by introducing new features in the module itself or extending its functionality by creating new independent modules that live outside the *ns-3* project code base, e.g., the NR module, which is also known as 5G-LENA (see Table 1 for more details).

Towards the end of 2017, the 3GPP, in release 15, standardized the NR technology to operate in Non-Standalone (NSA) mode[7]. The introduction of NR access technology started the 5G era, resulting in tremendous research interest from academia and industry. Driven by this, CTTC and Interdigital started working on extending *ns-3* to build a 5G simulator that is 1) open source, 2) able to support end-to-end full protocol stack simulations, 3) capable of working in FR1 and FR2 frequency ranges, as defined in 3GPP, 4) able to support the research focused on coexistence studies, and 5) as much as possible standard compliant. Building upon these objectives, in 2019, CTTC released its NR module that is easily pluggable to *ns-3*. It supports the NSA 5G architecture, which uses the LTE core network, i.e, Evolved Packet Core (EPC), implementation of the *ns-3* LTE module. Currently, the module relies on the LTE module for the layers above MAC. That is, it has a completely new MAC and PHY layers that support a flexible frame structure and variable sub-carrier spacings (i.e., numerologies), the Bandwidth

---

[7]In NSA architecture, 5G RAN is anchored to the LTE core.

Part (BWP) concept standardized in 5G, Low Density Parity Check (LDPC) coding for data channels, MCSs up to 256-QAM, Time Division Duplex (TDD), beamforming, among others (see Table 1). Similar to the LTE module, this module is also actively been extended to provide new interesting features, e.g., the New Radio-based access to Unlicensed spectrum (NR-U) to support the coexistence studies involving Wireless Gigabit (WiGig) in 60 GHz band, the support of dual-polarized Multiple Input Multiple Output (MIMO), and the NR V2X (contributed under this thesis framework). Moreover, recently, this module has been calibrated according to the 3GPP NR reference scenarios for outdoor deployments. Thus, validating the module's capability to achieve comparable performance to that of industrial proprietary simulators and real networks [56].

**Table 1.1:** *ns-3* LTE/NR models

| Modules | Characteristics | Thesis author´s contribution |
|---|---|---|
| LTE [52] | Full LTE protocol stack for Evolved Node B (eNB) and UE. | Added feature to simulate coverage holes (see Chapter 3 and [57]). |
| | LTE core network, i.e., EPC with single MME and multiple Service GateWay (SGW) and Packet data network GateWay (PGW) nodes. | Implemented deterministic handover algorithm (see Chapter 3 and [57]). |
| | | Ported RLF functionality from ELENA. |
| | Key features: | LTE module maintainer from 2018 to 2021. |
| | • Frequency Division Duplex (FDD) | |
| | • MIMO | |
| | • Hybrid Automatic Repeat Request (HARQ) | |
| | • Uplink power control | |
| | • RLC-UM | |
| | • RLC-AM | |
| | • Handover | |
| | • Radio Link Failure (RLF) | |
| | • Carrier Aggregation | |
| | • Fractional frequency reuse | |
| LAA/ LTE-U [58] | Extension of *ns-3* LTE module. | Helped to improve the code by performing extensive testing. |
| | Key features: | One of the co-authors of the most cited LAA and LTE-U paper based on *ns-3*. |
| | • LTE-U implementation based on Carrier Sense Adaptive Transmission (CSAT) | |
| | • 3GPP LAA implementation | |
| | *Continued on next page* | |

**Table 1.1 – continued from previous page**

| Modules | Characteristics | Thesis author´s contribution |
|---|---|---|
| ELENA [59] | Extension of *ns-3* LTE module.<br><br>Key features:<br><br>• Paging<br>• Time alignment<br>• RLF<br>• Handover failure<br>• Idle mode cell reselection<br>• Improved random access<br>• Improved Radio Resource Control (RRC) state machine | Reviewer of ELENA LTE module new features.<br><br>Ported and documented ELENA RLF implementation to the official *ns-3*. |
| psc-ns3 [60] | Extension of *ns-3* LTE module to support 3GPP LTE D2D communication.<br><br>Key features:<br><br>• Out-of-coverage synchronization<br>• In and out of coverage D2D Discovery<br>• In and out of coverage D2D Communication<br>• Various 3GPP aligned propagation models<br>• on-network and off-network Mission Critical Push To Talk (MCPTT) models<br>• Unmanned Aerial Vehicle (UAV) mobility energy model<br>• Hypertext Transfer Protocol (HTTP) application<br>• UE-to-Network Relay<br>• Video streaming model | Ported initial Public Safety Communications (PSC) module implementation to ns-3.29.<br><br>Author of additional examples and test scripts.<br><br>Author of initial PSC LTE model and user documentation. |
| ns-3 mmwave [61] | Based on *ns-3* LTE module.<br><br>Key features:<br><br>• Ray tracing<br>• NR-specific PHY and MAC layers<br>• Carrier Aggregation at the MAC layer<br>• Enhanced RLC layer<br>• Dual connectivity | — |

Continued on next page

**Table 1.1** – **continued from previous page**

| Modules | Characteristics | Thesis author´s contribution |
|---|---|---|
| nr [62] | Based on *ns-3* LTE module to provide 3GPP compliant NR technology implementation.<br><br>Key features:<br><br>• FDD<br>• TDD (configurable pattern)<br>• Multiple NR numerologies<br>• Time-Division Multiple Access (TDMA)<br>• Frequency-Division Multiple Access (FDMA)<br>• 3GPP compliant buffer status report<br>• 3GPP NR specific processing delays and control timings<br>• BWP support<br>• LDPC codes<br>• NR MCS table 1 and table 2<br>• Radio Environment Map (REM), Almost Blank Slots based<br>• Inter-Cell Interference Coordination (ICIC)<br>• Uplink power control<br>• Sounding Reference Signal (SRS) scheduling<br>• Realistic beamforming<br>• MIMO<br>• NR V2X | Acted as one of the maintainers of the module.<br><br>Lead developer for the NR V2X feature (see Chapter 7 and [62]). |

## 1.3.2   Linking *ns-3* with machine learning

In this thesis, we have used two well-known AI/ML techniques, i.e., Neural Network (NN) and Fuzzy Logic, to address the RAN use cases described in Sec 1.2. We have leveraged the supervised learning paradigm for the NN based models. Following this paradigm, the proposed models are trained offline using the synthetic dataset generated using *ns-3* simulation(s). In the scope of this thesis, the term "offline" implies that there is no close loop between the *ns-3* simulator and the AI/ML framework during the training, testing, and evaluation of the AI/ML model. Figure 1.2 presents the block diagram of the Offline-Training procedure adopted to train the proposed NN-based models. The measurements from a *ns-3* simulation environment (single or multiple executions) are preprocessed to form a labelled dataset. That is, besides containing the input parameters (also known

**Figure 1.2:** Offline-Training

as input features) to an ML model, each sample[8] in the dataset also contains an output. This output serves as a ground truth in the training phase to adjust the weights of the hidden layers of NN. This weight adjustment is needed to improve the inference capability of the model under training, which in turn decreases the error between the prediction and the expected output. Once the model is trained, its performance evaluation can be conducted using the unseen data[9] to assess the gain that can be achieved using the AI/ML solution(s) over the State-of-the-Art (SOTA) or benchmark schemes.

In this regard, we have used two performance evaluation methods 1) Offline-Evaluation, which is an approach most used in the simulation community for the AI/ML models based on supervised learning [45, 63, 64], and 2) Semi-Online-Evaluation method, which to the best of our knowledge is only used in [63], however, elaborated more comprehensively in this thesis. These two performance evaluation methods are illustrated in Fig. 1.3 and Fig. 1.4.

The Offline-Evaluation can be conducted using the following steps:

1. Conduct a simulation using the *ns-3* SOTA model, and save the resulted Key Performance Indicator (KPI) for the Evaluation stage.

2. Perform *ns-3* simulations using a deterministic algorithm (one such algorithm is explained in section 3.2.2) for the use case we intend to target using AI/ML solution. From these simulations, we construct a dataset that contains 1) input samples for the trained AI/ML model and 2) the output, i.e., the performance KPI, e.g., the QoE of the users, corresponding to each input sample that is to be used at the Evaluation stage.

3. Trigger the trained AI/ML model by providing one sample at a time from the dataset constructed in Step 2.

4. Store the inferred KPI by the AI/ML model for each input sample.

5. Select a configuration from the sample that achieved the best performance based on the inferred KPI.

6. Based on the selected configuration in Step 5, evaluate the performance by comparing the KPI from Step 2 with the KPI of the ns-3 SOTA model as acquired in Step 1.

---

[8]A sample is comprised of one or more input and output parameters.

[9]Unseen data comprises the measurements that were not part of the training dataset. In *ns-3*, these measurements are usually generated from the same simulation environment but using a different seed or Run number.

**Figure 1.3:** Offline-Evaluation

**Figure 1.4:** Semi-Online-Evaluation

In the Offline-Evaluation model, there is no interaction between the trained AI/ML model and *ns-3* simulation at the runtime. This is the key difference between the Offline-Evaluation and the Semi-Online-Evaluation, as shown in Fig. 1.4. The term "Semi-Online" refers to the fact that there is no retraining of the already trained AI/ML model during the Execution-Cycle.

The Semi-Online-Evaluation is comprised of the following steps:

1. Conduct a simulation using the *ns-3* SOTA model, and save the resulted KPI for the Evaluation stage.

2. In the Execution-Cycle, the process begins by starting a simulation that triggers the trained AI/ML model using a Bash script[10]. Thanks to the *ns-3* simulator, which uses C++ as its primary language, the BASH script is triggered using the C++ *"System"* function [65]. This function facilitates providing the inputs to the Bash script, which then are passed to the trained AI/ML model.

3. Upon receiving the input features from the previous step, the AI/ML model infers the desired KPI and writes it to a Text file.

4. The previous step's AI/ML model output is fed into the running *ns-3* simulation to (re-)configure the parameter(s) under observation.

5. The Execution-Cycle, i.e., Step 2 – Step 4, continues until the simulation ends.

6. Evaluate the performance by comparing the KPI computed during or at the end of the Execution-Cycle with the KPI of the *ns-3* SOTA model as acquired in Step 1.

The Offline-Evaluation model described above is used for the AI/ML models proposed to tackle the HO and initial MCS use-cases, whereas the Semi-Online-Evaluation is used for the models proposed for coexistence in unlicensed spectrum and V2X use-cases. Further details about the use-cases and the proposed AI/ML models are provided throughout Chapters 3-8 of this thesis.

## 1.4 Thesis organization and publications

Figure 1.5 shows the organization of the thesis. At the top of the figure is Chapter 2 that provides the fundamental knowledge of the RAN technologies and their functionalities that can leverage AI/ML-based solutions. For the reader who is not familiar with these RAN technologies, this chapter provides the knowledge to understand the details of each forthcoming chapters. In particular, it covers the following:

- The HO procedure that 3GPP has standardized. Specifically, the details on the different HO events based on which UE measurement reports get triggered and target cell selection criterion.

- The specifics about the LAA and LTE-U technologies and their respective channel access mechanisms.

- A comprehensive review of the history of V2X technology in 3GPP and its various releases.

We conclude the chapter by highlighting the gaps in the above RAN technologies targeted by the thesis. After Chapter 2, the Chapters 3-8 contain the technical contribution of the thesis that focuses on the four RAN use cases mentioned in the Subsection. 1.2. Each of these chapters first presents the "Related Work" section to review of the SOTA and then

---

[10]A Bash script is a set of UNIX commands written in a text file.

**Figure 1.5:** Thesis organization

discusses the work done beyond the SOTA. Due to the versatility of thesis that touches different RAN use cases, the choice of having dedicated "Related Work" section helps the reader to easily spot the differences between the already conducted research and the work done in this thesis. Hence, it improves the readability of the thesis.

The research work targeting the infrastructure-based single-technology RAN scenario is covered in Chapters 3 and 4. In **Chapter 3**, we present a simple single-task machine learning-based HO scheme to address the shortcoming of a standard HO algorithm. Specifically, we provide some interesting results that prove that a handover decision based only on the signal strength of the target Base Station (BS) is not enough in a challenging propagation scenario. The contributions of this chapter are published in the following two conference papers:

[C1] **Zoraze Ali**, Nicola Baldo, Josep Mangues-Bafalluy , Lorenza Giupponi, "Simulating LTE mobility management in presence of coverage holes with ns-3", *In Proceedings of the 8th International Conference on Simulation Tools and Techniques (SIMUTools)*, Athens, Greece, 24-26, August 2015

[C2] **Zoraze Ali**, Nicola Baldo, Josep Mangues-Bafalluy , Lorenza Giupponi, "Machine Learning Based Handover Management for Improved QoE in LTE", *In Proceedings of the IEEE/IFIP Network Operations and Management Symposium (NOMS)*, Istanbul, Turkey, 25-29, April 2016

**Chapter 4** extends the simulation scenario presented in Chapter 2 to depict a more realistic and complex deployment of multiple cells and more users. Using this scenario,

**17**

we deeply investigate the Multi-Task Learning (MTL) paradigms, i.e., parallel learning and incremental learning to address two different but related RAN use cases, i.e., HO management and the initial MCS selection when a UE establishes a new connection with the BS. The results presented in this chapter have been published in one conference paper and one journal paper.

[**C3**]  **Zoraze Ali**, Marco Miozzo, Lorenza Giupponi, Paolo Dini, Stojan Denic, Stavroula Vassaki, "Recurrent Neural Networks for Handover Management in Next-Generation Self-Organized Networks", *In Proceedings of the IEEE 31st Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, London, UK, 31 August 2020 - 03 September 2020

[**J1**]  **Zoraze Ali**, Lorenza Giupponi, Marco Miozzo, Paolo Dini "Multi-Task Learning for Efficient Management of Beyond 5G Radio Access Network Architectures", *IEEE Access*, Vol.9, pp. 158892-158907, 2021

In **Chapter 5** and **Chapter 6**, we focus on the infrastructure-based multi-technology RAN scenario. In particular, we study the coexistence performance of cellular technologies, i.e., LAA and LTE-U with WiFi in the unlicensed spectrum. In Chapter 5, we analyze the channel access mechanism of the LAA technology. The study reveals that LAA's contention window size procedure could hamper its fair coexistence with WiFi, which is essential for LAA to be deployed in an unlicensed spectrum. To address this issue, the chapter presents a machine learning solution that helps LAA to manage, i.e., increase or decrease, its contention window size to achieve fair coexistence with WiFi. To analyze the results, we have followed the 3GPP approach of comparing the throughput and latency Cumulative Distribution Function (CDF) plots. However, this approach is easy to employ when we have non-overlapping CDF curves, but it is hard to use with overlapping curves. Therefore, in Chapter 6, we propose a statistical frame work to evaluate better the fairness performance of the LAA and LTE-U technologies with WiFi. The results of these chapters have been published in the following two conference papers:

[**C4**]  **Zoraze Ali**, Lorenza Giupponi, Josep Mangues-Bafalluy, Biljana Bojovic, "Machine learning based scheme for contention window size adaptation in LTE-LAA", *In Proceedings of the IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, Montreal, QC, Canada, 08-13 October 2017

[**C5**]  **Zoraze Ali**, Lorenza Giupponi, Josep Mangues-Bafalluy, "On fairness evaluation: LTE-U vs. LAA", *in Proceedings of the 14th ACM International Symposium on Mobility Management and Wireless Access* (MOBIWAC), Malta 13-17 November 2016

Towards the end of the thesis, we aim to study the limited infrastructure-based single-technology RAN scenario, involving NR V2X technology. However, at that time, to the best of author's knowledge, there was no open-source standard compliant simulation tool to carry out our research. Therefore, the **Chapter 7** of this thesis targets the development of the first *ns-3* based NR V2X simulator. In particular, we extend the

*ns-3* LTE and NR (also known as 5G-LENA) modules to perform full-stack and end-to-end NR V2X simulations. The chapter covers in detail the 3GPP NR V2X technology, specifically Mode 2 for autonomous resource selection. In addition, we also provide the results of comprehensive simulation campaigns, studying the impact of key parameters on the standardized KPIs, i.e., Packet Inter-reception Delay (PIR) and Packet Reception Ratio (PRR). The contributions presented in this chapter have been published in one conference and one journal paper.

[**C6**] **Zoraze Ali**, Sandra Lagén, Lorenza Giupponi, "On the impact of numerology in NR V2X Mode 2 with sensing and random resource selection", *In the proceedings of the IEEE Vehicular Networking Conference (VNC)*, Ulm, Germany, 10-12 November 2021

[**J2**] **Zoraze Ali**, Sandra Lagén, Lorenza Giupponi, "3GPP NR V2X Mode 2: Overview, Models and System-Level Evaluation", *IEEE Access*, Vol.9, pp. 89554-89579, 2021

In **Chapter 8**, we use our NR V2X simulator to study the energy-performance trade-off when using sensing and random resource selection methods in NR V2X Mode 2. In particular, we argue that the trade-off mentioned above can be exploited better using a fuzzy inference system that dynamically adjusts the partial sensing duty cycle. The results are to be submitted to the Institute of Electrical and Electronics Engineers (IEEE) vehicular technology magazine.

[**J3**] **Zoraze Ali**, Sandra Lagén, Lorenza Giupponi, "NR V2X Mode 2 and the Energy-Performance Trade-off", *to be submitted to the IEEE vehicular technology magazine*

Lastly, **Chapter 9** concludes the thesis by drawing out the key aspects of our research and some future directions to help the research community advance in the areas tackled in this thesis.

## 1.4.1 Other journal/conference papers

In this section, we list the other conference and journal paper contributions that are related but are not included in this thesis.

[**J4**] Marco Miozzo, **Zoraze Ali**, Lorenza Giupponi, Paolo Dini "Distributed and Multi-Task Learning at the Edge for Energy Efficient Radio Access Networks", *IEEE Access*, Vol.9, pp. 12491-12505, 2019

[**J5**] Biljana Bojovic, Lorenza Giupponi, **Zoraze Ali**, Marco Miozzo "Evaluating Unlicensed LTE Technologies: LAA vs LTE-U", *IEEE Access*, Vol.7, pp. 89714-89751, 2019

[**J6**] Katerina Koutlia, Biljana Bojovic, **Zoraze Ali**, Sandra Lagén, "Calibration of the 5G-LENA system level simulator in 3GPP reference scenarios", *ELSEVIER*, Simulation Modelling Practice and Theory, Vol.119, pp. 102580, 2022

[**C7**] Tommaso Zugno, Matteo Drago, Sandra Lagén, **Zoraze Ali**, and Michele Zorzi, "Extending the ns-3 spatial channel model for vehicular scenarios", *In Proceedings of the Workshop on ns-3 (WNS3)*, New York, NY, USA, 21-25 June 2021

[**C8**] Biljana Bojovic, **Zoraze Ali**, Sandra Lagen, Katerina Koutlia, "ns-3 and 5G-LENA extensions to support dual-polarized MIMO", *In Proceedings of the Workshop on ns-3 (WNS3))*, Virtual workshop, 20-24 June 2022

# Chapter 2

# Fundamentals

This chapter provides the fundamental knowledge of the mechanisms and technologies needed to understand the thesis better. Moreover, as mentioned in Chapter 1, this thesis leverages *ns-3* and AI/ML to address the identified RAN use cases. Specifically, it proposes supervised learning-based solutions that use well-known Feed-forward Neural Network (FFNN), Recurrent Neural Network (RNN), and Fuzzy-logic techniques. Given the huge interest from academia in AI/ML, a comprehensive amount of literature is already available surveying the aforementioned AI/ML techniques. Therefore, for the sake of conciseness and focus, in this thesis, we do not cover such background and refer the reader to the vast available literature on these techniques (for instance, [66], [67], [45], [42], [63], [7], and [18]).

In this chapter, Section 2.1 presents the details of the HO procedure standardized by 3GPP. Section 2.2 gives a technical background of LTE LAA and LTE-U technologies. Then, Section 2.3 summarizes the evolution of 3GPP V2X technologies. Lastly, Section 2.4 concludes the chapter.

## 2.1   3GPP Handover

The handover procedure is an essential part of connected mode mobility management. It guarantees the continuity of the services provided by a mobile network when a UE in a connected state moves around. The 3GPP standards, since the inception of Wideband Code Division Multiple Access (WCDMA) technology, support the HO functionality, which is also part and parcel of the 4G, 5G, and, not yet fully defined, 6G technologies.

These technologies encompass critical use case scenarios, e.g., Ultra-Reliable and Low-Latency Communications (URLLC), which means high reliability and quick availability for smart manufacturing, connected vehicles, electrical power distribution, and more, such as drones controlled by the network [68]. This requires a fast, up to a few milliseconds, and reliable HO procedure to maintain an adequate level of QoE.

As per the 3GPP standards, a base station controls the mobility of its UEs by configuring one or more event-based measurements. As per this configuration, the UE conducts radio measurements of its surrounding environment. When the event-triggering conditions are met, the UE sends the measurement reports to the base station. The base station may trigger a handover based on these measurements. Up to now, the 3GPP standard supports the following two measurement quantities:

1. Reference Signal Receive Power (RSRP)

2. Reference Signal Received Quality (RSRQ)

There are five handover events, which get triggered based on either of the two measurement quantities [38,69]. These handover events are listed below.

- Event A1. The source cell becomes better than a threshold.

- Event A2. The source cell becomes worse than a threshold.

- Event A3. The neighbour cell becomes better than an offset relative to the source cell.

- Event A4. The neighbour cell becomes better than a threshold.

- Event A5. The source cell becomes worse than one threshold, and neighbour cell becomes better than another threshold.

Fig 2.1 illustrates these events using RSRP as a measurement quantity. Moreover, it shows different parameters, e.g., thresholds, offset, and hysteresis, that influence the "entry condition" of these events, after which the UE sends the measurement reports. A UE measurement report, among other useful information, includes the signal strength (RSRP or RSRQ) of the source and the neighbour cells that fulfil the configured event criterion. The source cell, upon receiving the measurement, chooses a neighbour that has the strongest signal strength as the target cell and initiates the handover process by sending it the "Handover Request" message. The target cell, if have enough radio resources to serve the UE, accepts the handover request. After that, the UE is instructed to terminate its radio link to the source cell and establish a new connection with the target cell.

**Figure 2.1:** Handover events

## 2.2 Unlicensed technologies background

The quest for increasing the network capacity and cost-effectively reducing the price per megabyte is the main challenge of network operators; various solutions are being considered, from the densification of small cells to the offload of traffic in the unlicensed band. The network densification through deploying different types of small cells, e.g., pico and femto (also known as small cells), has undoubtedly increased the network capacity [70]. However, it also introduces inter-cell-interference among these different types of small cells since the mobile networks usually operate in full-frequency reuse (i.e., reuse 1) manner. Therefore, the mechanisms such as ICIC and enhanced Inter-Cell Interference Coordination (eICIC) must be employed to reap the real benefits of such deployments [33]. However, the scarcity of the licensed spectrum, especially at low-frequency ranges, e.g., 410 MHz – 7125 MHz (also known as Frequency Range 1 (FR1)), is a bottleneck for such a solution. Thus, the complementary use of unlicensed spectrum has been part of the 3GPP standard since its Release 6. These technologies can be divided into two categories [45]:

- Technologies that provide access to the unlicensed band using inter-Radio Access Technology (RAT). For example, Interworking WLAN (I-WLAN) (Release 6 and 8 [71, 72]), RAN-controlled WLAN Interworking (Releases 12 and 13 [73, 74]), LTE-WLAN aggregation (LWA) (Releases 13 and 14 [74, 75]), and LTE-WLAN radio level integration with IPsec tunnel (LWIP) (Release 13 [74]).

- Technologies that provide access to the unlicensed band using the same RAT. For example, LTE-U [76], LAA (Release 13 and 14 [74, 75]), and NR-U (NR Release 16 [77]).

For the coexistence use case, this thesis focuses on the LAA and LTE-U technologies. The choice is motivated by the innovative nature of these two technologies, which bring a paradigm shift about how the network and devices based on the 3GPP technology coexist in unlicensed bands, which inherently are built to operate in licensed spectrum. This paradigm shift creates infinite opportunities for research to study the coexistence among these and other technologies, e.g., WiFi. Therefore, we provide a technical background of

**Table 2.1:** LAA channel access priority classes

| Channel access priority class | $m_p$ | $CW_{min}$ | $CW_{max}$ | $T_{mcot,p}$ | Allowed CW sizes |
|---|---|---|---|---|---|
| 1 | 1 | 3 | 7 | 2 ms | 3, 7 |
| 2 | 1 | 7 | 15 | 3 ms | 7, 15 |
| 3 | 3 | 15 | 63 | 8 or 10 ms | 15, 31, 63 |
| 4 | 7 | 15 | 1023 | 8 or 10 ms | 15, 31, 63, 127, 255, 511, 1023 |

these two technologies in the following subsections. To know about the other unlicensed technologies mentioned above, readers are referred to [45].

Through the principle of carrier aggregation, the LAA and LTE-U technologies aim toward the deployment of LTE networks in the unlicensed spectrum, i.e., the 5 GHz band. These technologies boost the performance of LTE networks by providing wider bandwidth in high-traffic areas. One of the principles for their design is the different regional regulatory regimes, which may or may not mandate the use of Listen-Before-Talk (LBT) procedures to access the wireless channel. For example, in Europe and Japan, access to the unlicensed spectrum is subject to LBT requirements. On the other hand, there is no such requirement in the USA, China, India, and Korea.

Nevertheless, these technologies cannot be used without ensuring a fair and friendly coexistence with other incumbent technologies in the unlicensed band [78]. Coexistence-fairness with existing technologies, especially Wi-Fi, is the fundamental requirement for deploying them in the unlicensed spectrum [79]. In what follows, we will explain in detail the working of both technologies.

## 2.2.1 LAA

Third Generation Partnership Project (3GPP) in LTE releases 13 and 14 introduced the LAA technology to meet ETSI's Clear Channel Assessment (CCA)/LBT requirements [79]. And because of this, it is considered to be a global solution framework that allows compliance with any regional regulatory requirements. In the beginning, 3GPP analyzed different categories of LBT for LAA, and finally, the most similar to the Carrier Sense Multiple Access with Collision Avoidance mechanism (CSMA/CA) of Wi-Fi was selected. This is referred to as Category 4 LBT [79]. In LTE release 13, LAA was standardized only to support the transmissions in the downlink channel, i.e., Physical Downlink Shared Channel (PDSCH). After that, in release 14, in the context of Enhanced LAA (eLAA), the support for the uplink channel, i.e., Physical Uplink Shared Channel (PUSCH) was added. However, many efforts were made by academia and industry to evaluate the LAA performance in the downlink [45, 79]. In this line, this thesis also targets the LAA operation in the downlink, and here, we will explain how an LAA eNB is enabled to access the unlicensed channel and how its Contention Window (CW) evolves upon collision as specified by the standard [80].

Fig. 2.2, shows the complete algorithm for category 4 LBT. An eNB, which intends to transmit the data in the downlink, first performs an *Initial CCA* during which it senses the channel for a defer duration of $T_d$. The defer duration $T_d$ is composed of duration

**Figure 2.2:** LBT 3GPP Category 4 algorithm

$T_f = 16$ $\mu$s and $m_p$ CCA slots, where each slot duration is $T_{sl} = 9$ $\mu$s. The value of the $m_p$ depends on the LBT priority class, as shown in Table 2.1 [80]. If the channel is idle during $T_d$, the eNB occupies the channel for maximum $T_{mcot,p}$ duration, which is also known as a Transmission Opportunity (TxOP). The duration of the TxOP depends on the LBT priority class, which categorizes the type of traffic scheduled in the unlicensed band [80]. An LAA eNB can occupy the channel up to 10 ms in case of Best Effort (BE) and Background (BK) traffic, i.e., priority classes 3 and 4, respectively. For other types of traffic requiring higher service quality, the length of the TxOP is shorter [80]. On the other hand, if the channel is busy during the $T_d$ period, the eNB performs an *Extended CCA*. Under it, similar to Wi-Fi, the eNB draws a random counter $N$ in the range of [0,$CW_p$], where $CW_p$ is the current CW size, which ranges between $CW_{min}$ and $CW_{max}$. If the eNB finds the channel idle for $N$ CCA slots, it occupies the channel for $T_{mcot,p}$ duration. The CW size is increased exponentially upon collisions, which in LAA are detected using HARQ feedbacks from a receiving node. In particular, the CW size

**Figure 2.3:** HARQ feedback timing diagram (FDD mode)

at the eNB is increased if 80% of the HARQ feedbacks belonging to the first subframe in the most recent TxOP are NACKs [80].

The rationale behind such a rule is twofold. On the one hand, the eNB may schedule more than one UE in a single subframe. Thus, it will receive multiple HARQ feedbacks that have to be translated into a single decision about if the collision has happened or not. On the other hand, the rule considers only the feedbacks from the first subframe of the TxOP to reduce the delay in updating the CW size. As shown in Fig. 2.3, the data transmitted by the eNB in subframe *n* is acknowledged by the UE in subframe *n+4*, i.e., after 4 ms from the data transmission [38]. Therefore, to update the CW based on the HARQ feedbacks from all the subframes in a TxOP, the eNB should wait till the HARQ feedback of the last subframe. This introduces a further delay in deciding whether a collision has occurred or not. It is worth mentioning that, while the LTE protocol stack induces these high delays, in Wi-Fi, a receiving Station (STA) upon the correct reception of a data frame transmits an acknowledgment (ACK) after Short Interframe Space (SIFS) of 16 $\mu$s.

## 2.2.2 LTE-U

The LTE-U coexistence paradigm is specified by the LTE-U Forum [81, 82]. It is an industry consortium formed in 2014 by Verizon, Ericsson, Alcatel-Lucent, Qualcomm Technologies, Inc., and Samsung. The LTE-U forum has the following main objectives:

- To develop a proprietary solution that would enable the coexistence of LTE and Wi-Fi networks in the 5 GHz unlicensed band.

- The solution would target the markets without the LBT requirements.

- The solution should be agile, i.e., it can be quickly deployed with minor changes to the LTE Release 10/11/12 carrier-aggregation protocol.

- It would be used in a Supplemental DownLink (SDL) carrier (i.e., the secondary carrier) in conjunction with a licensed carrier.

Based on the above objectives, the LTE-U technology does not use the LBT procedure like LAA to access the channel. Instead, it follows a CSAT procedure specified by the LTE-U Forum. Following the CSAT procedure, the LTE-U secondary cell duty cycles its transmissions, i.e., alternates ON and OFF periods, as shown in Fig. 2.4. The duty cycle length is the sum of $T_{ON}$ and $T_{OFF}$ durations. During the $T_{OFF}$ period, an LTE-U

**Figure 2.4:** LTE-U CSAT duty cycle period

secondary cell does not transmit, giving a chance to other coexisting technologies, e.g., Wi-Fi, to transmit. On the other hand, during the $T_{ON}$ period, LTE-U start transmitting without sensing the channel. During the $T_{ON}$ period, there are a few subframes that are periodically left blank, i.e., an LTE-U eNB does not schedule any transmission during these subframes. As per the LTE-U standard, this process is referred to as *puncturing*. It allows latency-sensitive applications that run over Wi-Fi to transmit. The quantity of such punctured subframes, usually 1 or 2, depends on the length of the $T_{ON}$ period and the LTE-U traffic load. The CSAT duty cycle parameters, i.e., the $T_{ON}$ and $T_{OFF}$ must be configured such that a fair coexistence can be achieved with other technologies operating in the same band [81]. Alternatively, these parameters can be configured adaptively, e.g., by estimating the most appropriate channel share it should occupy, depending on the other networks´ activity. In this line, the LTE-U Forum does not specify any algorithm; therefore, its implementation is vendor specific. The most representative algorithm of such nature is the Qualcomm CSAT/eCSAT algorithm [83]. The details of this algorithm are out of the scope of this thesis. Therefore, for a detailed overview of this algorithm, the interested reader is referred to [45].

## 2.3 Evolution of 3GPP V2X technologies

The automotive industry is currently transitioning towards automated driving and advanced driver assisted systems, where vehicles are able to react by themselves to changes in the driving environment. In this context, V2X is seen as a key technology to provide complete environmental awareness around the vehicle by exchanging messages with other vehicles, roadside units, and pedestrians with low latency and high reliability. V2X communications are expected to provide potentiality in different areas, like faster alerts and notifications, law enforcement, better service on roadways, reduced world-wide traffic load, reduced emissions, time savings, and increased automotive safety, thus contributing to prevent crashes/injuries and save lives [84]. Additionally, V2X-capable vehicles can assist in better traffic management also for non-safety applications. Several advanced V2X use cases have been already proposed within the 3GPP Release 15 such as vehicle platooning, extended sensors, advanced and remote driving, or cooperative collision avoidance [85]. Also, industrial associations like the 5G Automotive Association (5GAA) in Europe have been built to promote the vision of connected mobility, including autonomous driving and intelligent transportation [86].

As of today, the two key radio access technologies that enable vehicular communications are 1) Dedicated Short Range Communications (DSRC), standardized by IEEE in 802.11p [87] and the more recent 802.11bd [88], and 2) LTE C-V2X, based on 3GPP LTE Release 14 and Release 15 [89]. DSRC is designed to primarily operate in the 5.9 GHz band, while C-V2X is thought to operate in both 5.9 GHz and in cellular licensed carriers at sub 6 GHz carrier frequencies. Differently from DSRC that focused on Vehicle-to-Vehicle (V2V) and Vehicle-to-Infrastructure (V2I) communications, C-V2X encompasses V2V, Vehicle-to-Pedestrian (V2P), V2I and Vehicle-to-Network (V2N) [90]. C-V2X is designed to support basic safety message sharing among proximity users, such as collision warning, emergency stop warning, and adaptive cruise control. A comparison study in [91] shows that LTE C-V2X gets a superior reliability performance over DSRC, due to the more efficient Physical Layer (PHY) layer of LTE C-V2X.

V2X requirements can be met using LTE C-V2X, as long as the vehicular density is not too high [89]. However, as the quality of service requirements become more stringent, which is the case in many V2X applications, LTE C-V2X falls short, and 5G NR is called as a complementary solution [88]. Towards that goal, 3GPP Release 16 has included a Study Item (SI) to support new applications with more stringent requirements, which has resulted in Technical Report (TR) 38.885 [92]. Based on the study outcome captured in this TR, 3GPP Release 16 has completed a Work Item (WI) in July 2020 to standardize V2X on top of 5G NR standardized in Release 15 [93]. As a main design principle, NR is not designed to be backward compatible with LTE. Similarly, NR V2X is not backward compatible with LTE C-V2X. The NR V2X SI indicates that the design objective of NR V2X is not to replace LTE C-V2X, but to supplement C-V2X in supporting those use cases that cannot be supported by LTE C-V2X [92]. To ensure that NR V2X can provide a unified support for all V2X applications in the future, NR V2X must be capable of supporting not only advanced V2X applications but also basic safety applications that are supported today by LTE C-V2X.

Such a wide applications and use cases' support is possible in NR V2X because of the flexible framework inherited by the NR technology and the recent progresses envisioned in NR V2X standardization, which includes many enhancements over LTE C-V2X concepts. The NR radio access technology provides wide bandwidth support in various frequency ranges (including sub 6 GHz bands and mmWave bands), flexible frame structure with reduced transmission time intervals (by means of multiple numerologies and sub-carrier spacing (SCS) support), support for massive MIMO systems and high modulation orders, and advanced channel coding [94]. All these new features and functionalities intrinsically contribute to increase the data rate, reduce the latency, and improve the spectral efficiency of V2X communication systems. In addition, new enhancements and key procedures have been defined for NR V2X, specifically designed to improve the reliability of V2X communications systems, such as new communication types (unicast and groupcast), a new feedback channel, the support of feedback-based retransmissions, and new resource allocation and scheduling mechanisms [92].

In the following subsection, we review the history of sidelink communications and the NR V2X standard with particular emphasis on NR V2X Mode 2. Then, in Chapter 7, we present the proposed simulation tool and provide the implementation details of the developed NR V2X models, including the design choices and implementation changes that affect all the layers of the protocol stack. Based on these models, we present a

comprehensive set of simulation campaigns, in which we study the impact on the end-to-end network performance of different key parameters of the NR V2X system when using sensing-based resource selection in NR V2X Mode 2, as defined by 3GPP. In particular, we study the impact of the numerology, the number of retransmissions, the length of the resource selection window, the maximum number of resources per reservation, the probability of keeping the same resource over multiple reservation periods, and the MCS. Finally, we compare different resource selection procedures for NR V2X Mode 2 considered in 3GPP, including sensing-based and random resource selections. From these detailed end-to-end campaigns, we derive interesting insights on the technology, which are summarized at the end of the Chapter 7.

### 2.3.1 History of sidelink technology in 3GPP

The concept of sidelink was first introduced in Release 12, together with D2D communications extensions to the traditionally centralized paradigm of cellular communications promoted by 3GPP. LTE C-V2X first and NR V2X later are all significantly based on previous D2D efforts [95,96]. In this section we review the history of D2D and V2X technologies inside 3GPP, giving special emphasis to Releases 12, 14, and 16, which are those starting the definition of the new D2D, C-V2X, and NR V2X technologies, respectively.

Tables 2.2 and 2.3 summarize the evolution of sidelink communications in 3GPP, since its introduction with D2D (Release 12/13), through LTE C-V2X (Release 14/15) and up to date in NR V2X (Release 16/17), including various SIs and WIs related to sidelink communication studies and standardization. For each SI/WI, we specify the 3GPP release, the working/leading group in charge of such SI/WI, the objective of the SI/WI, the resulting TR for the case of SIs and the impacted Technical Specification (TS) for the case of WIs. The history of sidelink is split into two tables. Table 2.2 covers from Release 12 until Release 15 and Table 2.3 covers the SIs/WIs from Release 16 until Release 17.

### 2.3.2 D2D (Release 12)

D2D has been defined as a support for Proximity Services (ProSe). D2D enables the quick exchange of data over short distances via a direct link between nodes and introduces a new interface, the PC5, between nodes. This offers an efficient way to bypass the LTE base station (or eNB) and offload the eNB traffic. Besides content sharing, a D2D UE can act as a relay for another device with a poor connection to the eNB and, therefore, D2D can be used to extend cellular network coverage. Two modes have been defined for centralized and distributed scheduling of UE transmissions, namely Mode 1 and Mode 2. Centralized scheduling occurs at the eNB (in-coverage mode), whereas distributed scheduling is carried out by the D2D UEs themselves, with no need to be in the coverage area of an eNB (out-of-coverage mode). In Mode 1, the UEs are scheduled by the eNB over dedicated radio resources for data transmission. In Mode 2, a UE can autonomously select a radio resource from a resource pool, which is either configured by the network or pre-configured in the user device for its direct D2D communication over PC5 interface.

Both modes share the same resource allocation structure, in which the transmission of data is scheduled within the so-called sidelink control period. Within this period, a set of subframes are allocated for the PSCCH transmission and a different set of subframes are allocated for the PSSCH. The corresponding PSCCH for a given PSSCH is always sent before the PSSCH data. The PSCCH contains the Sidelink Control Information (SCI), also called scheduling assignment, which is used by the receiver to identify the occupation of the PSSCH radio resources. In both modes, the SCI is configured in format 0, and it is transmitted twice using two different subframes in which it occupies the same Resource Block (RB). The second transmission is needed to improve the reliability of the SCI message delivery at the receiver due to the lack of a feedback channel in sidelink communication. The receiver blindly detects the SCI by monitoring all possible PSCCH resources. The transport block is transmitted four times in four consecutive subframes within the resource pool. This allows the receiver UE to implement open loop HARQ by combining the four redundancy versions of the transport block.

The operational principle of Modes 1 and 2 is battery life improvement of mobile devices. Vehicular communications have, however, other constraints that cannot be accommodated with D2D ProSe. Specifically, the high latencies of D2D are not suitable for vehicular communications, where packet delays or packet losses can have severe and life-threatening consequences. In terms of requirements, the maximum allowed latency varies between 20 ms and 100 ms, depending on the application, with reliability from 80 % to 95 % [90].

**Table 2.2:** Evolution of sidelink in 3GPP (part 1: Release 12 till Release 15).

| Release | WI/SI | Group | Objective | TR/TS |
|---------|-------|-------|-----------|-------|
| Release 12 | SI: Study on LTE D2D Proximity Services - Radio Aspects | RAN1, RAN2, RAN3, RAN4 | To define the methodology to evaluate LTE D2D proximity services, identify PHY layer options and enhancements | TR 36.843 |
| Release 12 | WI: Proximity-based Services | SA1, SA2, SA3 | To specify service requirements for ProSe discovery and ProSe communication over E-UTRA | TS 21.905, 22.115, 22.278, 23.002, 23.122, 23.303, 23.401, 23.402, 24.301, 33.220, 33.303, 33.833, 36.413, 36.423 |
| Release 13 | WI: Enhanced LTE D2D Proximity Services | RAN2, RAN1, RAN3, RAN4 | To define enhancements to LTE D2D communications and discovery meeting requirements for public safety applications | TS 36.101, 36.104, 36.133, 36.141, 36.211, 36.213, 36.214, 36.300, 36.301, 36.304, 36.306, 36.321 36.331, 36.413, 36.423 |
| | | | | Continued on next page |

**Table 2.2 – continued from previous page**

| Release | WI/SI | Group | Objective | TR/TS |
|---|---|---|---|---|
| Release 13 | WI: Enhancements to Proximity-based Services | SA1 | To support stage 2/3 development during Release 13 and support end of Release 12 maintenance to review and ensure that Release 13 TS 22.278 and TS 22.115 contain all agreed ProSe Stage 1 requirements | TS 23.303, 33.303 |
| Release 14 | SI: Study on LTE support for V2X services | SA1 | To study service requirements for V2P, V2P, V2N/V2I | TR 22.885 |
| Release 14 | WI: LTE support for V2X services | SA1 | To specify service requirements for V2P, V2P, V2N/V2I | TS 22.185 |
| Release 14 | WI: Support for V2V services based on LTE sidelink | RAN1, RAN2, RAN3, RAN4 | To specify LTE sidelink enhancements for V2V services defined in TR 22.885 | TS 36.101, 36.104, 36.133, 36.141, 36.201, 36.211, 36.212, 36.213, 36.214, 36.300, 36.302, 36.304, 36.306, 36.307, 36.321, 36.323, 36.331, 36.413, 36.423 |
| Release 15 | WI: ProSe Support for Band 72 in LTE | RAN5 | To update the 3GPP RAN WG5 RF, RRM and Protocol conformance test specification with the support of ProSe for Band 72 | TS 36.508, 36.521 |
| Release 15 | WI: Remote UE access via relay UE | SA1 | To specify service requirements for a UE with UICC to connect with network via an Evolved ProSe UE-to-Network Relay | TS 22.011, 22.115, 22.278, |
| Release 15 | SI: Study on Enhancement of 3GPP support for V2X services | SA1 | To identify use cases and potential service requirements to enhance 3GPP support for V2X service in safety and non-safety V2X scenarios | TR 22.886 |
| Release 15 | WI: Enhancement of 3GPP support for V2X scenarios | SA1 | To specify service requirements to enhance 3GPP support for V2X scenarios valid for the 3GPP systems (i.e., 5G, EPS), including the transport layer support for safety and non-safety V2X scenarios | TS 22.186 |
| | | | | Continued on next page |

**Table 2.2 – continued from previous page**

| Release | WI/SI | Group | Objective | TR/TS |
|---------|-------|-------|-----------|-------|
| Release 15 | WI: V2X new band combinations for LTE | RAN4 | To specify RAN4 RF requirements for the concurrent operation of additional LTE Uu frequency bands and PC5 operation on Band 47 and for concurrent operation of LTE Uu Carrier Aggregation and PC5 operation on Band 47 | TS 36.101, 36.307 |
| Release 15 | WI: Enhancements on LTE-based V2X services | RAN1, RAN2, RAN3 | To define enhancements on LTE-based V2X Services | TS 23.285, 23.303, 24.334, 24.385, 24.386, 36.101, 36.133, 36.201, 36.211, 36.212, 36.213, 36.300, 36.302, 36.304, 36.306, 36.321, 36.323, 36.331 |
| Release 15 | SI: Study on security aspects for LTE support of V2X services | SA2 | To identify and evaluate potential architecture enhancements needed to operate LTE-based V2X (V2V, V2I/N, and V2P), based on vehicular services requirements defined in SA1 V2X LTE and determine which of the solutions can proceed to normative specification | TR 33.885 |
| Release 15 | SI: Study on evaluation methodology of new V2X use cases for LTE and NR | RAN1 | To establish the evaluation methodology to evaluate technical solutions supporting the full set of 5G V2X use cases as identified in TR 22.886 and the full set of 5G RAN requirements in TR 38.913 | TR 37.885 |
| Release 15 | SI: Study on further enhancements to LTE Device to Device (D2D), UE to network relays for IoT (Internet of Things) and wearables | RAN2, RAN1, RAN3, RAN4 | To study enhancements to Prose UE-to-network relaying and to the LTE D2D framework for commercial and public safety applications such as wearable devices | TR 36.746 |

### 2.3.3 LTE C-V2X (Release 14)

3GPP Release 14 extended the D2D ProSe functionality by adding two new modes, Modes 3 and 4, for LTE C-V2X connectivity. Basic safety messages and event-triggered messages are transmitted for collision avoidance. V2V mainly enables cooperative automated driving. V2P establishes the communications protocol between vehicles and pedestrians for pedestrian safety. V2I implies the communications with roadside units and allows

to make information about local road and traffic conditions readily available to vehicles. V2N enables commercial services by providing access to data stored in the Cloud.

Modes 3 and 4 have been designed to satisfy the latency requirements and accommodate high Doppler spreads and high density of vehicles for LTE C-V2X communications. Similarly to Mode 1, Mode 3 uses the centralized eNB scheduler. The vehicular UE and eNB use the Uu interface to communicate. This transmission mode is only available when the vehicles are under cellular coverage. UE context information in terms of traffic patterns, for example, can be reported to the eNBs in order to assist in the resource allocation procedure. Mode 4 employs distributed UE scheduling, as Mode 2. In contrast to Mode 3, Mode 4 can operate without cellular coverage. However, these modes share a completely different structure than Modes 1 and 2 described above, when it comes to the allocation of the PSCCH and PSSCH. First, PSCCH and PSSCH channels are not separated in the temporal domain, but in the frequency domain. The resource grid is divided into sub-bands or sub-channels in which the first RBs of these sub-channels form the PSCCH pool and, the other RBs, the PSSCH pool.

**Table 2.3:** Evolution of sidelink in 3GPP (part 2: Release 16 till Release 17).

| Release | WI/SI | Group | Objective | TR/TS |
|---------|-------|-------|-----------|-------|
| Release 16 | SI: Study on Improvement of V2X Service Handling | SA1 | To identify use cases and potential service requirements to enhance 3GPP support for V2X | TR 22.886 |
| Release 16 | WI: Improvement of V2X Service Handling | SA1 | To define use cases and potential service requirements to enhance 3GPP support for V2X, based on the studies in TR 22.886 | TS 22.186 |
| Release 16 | SI: Study on application layer support for V2X services | SA6 | To develop key issues, corresponding architecture requirements and solution recommendations to enable the application layer support for V2X services over 3GPP systems | TR 23.795 |
| Release 16 | WI: Application layer support for V2X services | SA6 | To define architecture requirements, functional architecture, procedure and information flows, based on solutions and conclusions reached in TR 23.795 | TS 23.286, 23.795, 24.486, 24.587, 27.007, 29.486 |
| Release 16 | SI: Study on architecture enhancements for the Evolved Packet System (EPS) and the 5G System (5GS) to support advanced V2X services | SA2 | To identify and evaluate potential architecture enhancements of EPS and 5G System design needed to support advanced V2X services identified in TR 22.886 | TR 23.786 |
| | | | | Continued on next page |

**Table 2.3 – continued from previous page**

| Release | WI/SI | Group | Objective | TR/TS |
|---|---|---|---|---|
| Release 16 | WI: Architecture enhancements for 3GPP support of advanced V2X services | SA2 | To specify architecture enhancements of 5G system to support advanced V2X services as per conclusions reached within TR 23.786 | TS 23.008, 23.122, 23.285, 23.287, 23.501, 23.502, 23.503, 24.007, 24.301, 24.385, 24.386, 24.501, 24.587, 24.588, 27.007, 29.122, 29.230, 29.272, 29.274, 29.388, 29.502, 29.503, 29.504, 29.505, 29.510, 29.512, 29.513, 29.514, 29.518, 29.519, 29.520, 29.522, 29.525, 29.571, 31.102, 33.185, 33.535, 33.536, 38.413, TS 38.423 |
| Release 16 | SI: Study on NR Vehicle-to-Everything (V2X) | RAN1, RAN2, RAN3 | To study sidelink design, Uu enhancements for advanced V2X use cases, Uu-based sidelink resource allocation/configuration, RAT/Interface selection for operation, QoS management, and coexistence | TR 38.885 |
| Release 16 | SI: Study on V2X Media Handling and Interaction | SA4 | To study use cases relevant to transmission of multimedia over 3GPP and detail the requirements and procedures for media capturing, compression, and transmission | TR 26.985 |
| Release 16 | SI: Study on Security Aspects of 3GPP support for Advanced V2X Services | SA3 | To provide security and privacy analysis of eV2X system architecture, derive potential security and privacy requirements, and evaluate security and privacy solutions for protection of it | TR 33.836 |
| | | | | Continued on next page |

**Table 2.3** – continued from previous page

| Release | WI/SI | Group | Objective | TR/TS |
|---------|-------|-------|-----------|-------|
| Release 16 | WI: 5G V2X with NR sidelink | RAN1, RAN2, RAN3, RAN4 | To specify radio solutions that are necessary for NR to support advanced V2X services (except the remote driving use case which was studied in TR 38.824) based on the study outcome captured in TR 38.885 | TS 36.133, 36.300, 36.304, 36.306, 36.321, 36.331, 36.413, 36.423, 37.324, 37.340, 38.101, 38.104, 38.133, 38.201, 38.202, 38.211, 38.212, 38.213, 38.214, 38.215, 38.300, 38.304, 38.306, 38.321, 38.323, 38.331, 38.413, 38.423, 38.460, 38.463, 38.470, 38.473, 38.886 |
| Release 17 | WI: NR Sidelink enhancement | RAN1, RAN2, RAN4 | To specify radio solutions that can enhance NR sidelink for the V2X, public safety and commercial use cases, with special focus on power saving, enhanced reliability and reduced latency | [none yet] |
| Release 17 | SI: Study on NR Sidelink relay | RAN2 | To study single-hop NR sidelink-based relay | TR 38.836 |
| Release 17 | SI: Study on enhancements to application layer support for V2X services | SA6 | To study enhancements to the application architecture to support V2X services specified in 3GPP TS 23.286 | TR 23.764 |
| Release 17 | WI: Enhanced application layer support for V2X services | SA6 | To define enhancements to the application architecture to support V2X services specified in 3GPP TS 23.286 | TS 23.286, 23.434, 27.007 |
| Release 17 | WI: Band combinations for concurrent operation of NR/LTE Uu bands/band combinations and one NR/LTE V2X PC5 band | RAN4 | To specify band specific RF requirements for the concurrent operation of NR Uu and NR PC5, LTE Uu and NR PC5, NR Uu and LTE PC5 | TR 37.875, TS 38.101 |
| Release 17 | SI: Study on V2X services - Phase 2 | SA2 | To study procedures for V2X authorization and V2X communication | TR 23.776 |

A new SCI format, format 1, is employed. In Modes 3 and 4, a transport block can be sent either once or twice. In case of two transmitting attempts, the information is sent over another subframe, with the same structure: two SCIs and their corresponding PSSCH transport block. In this case, the receiver also implements HARQ combining. Vehicles select their sub-channels in Mode 4 using the sensing-based Semi-Persistent Scheduling (SPS) scheme specified in Release 14. Thanks to the semi-persistent reservation of resources and the inclusion of the reselection counter and packet transmission interval in the SCI, other vehicles can estimate which subchannels/subframes are free when making

**Table 2.4:** Sub 6 GHz NR V2X bands

| V2X operating bands | Sidelink (SL) Tx/Rx operating band $F_{\text{low}}$ - $F_{\text{high}}$ [MHz] | Duplex Mode | Sub carrier spacing [kHz] | Supported bandwidth [MHz] |
|---|---|---|---|---|
| n38 (Licensed) | $2570 - 2620$ | TDD | 15, 30, 60 | 10, 20, 30, 40 |
| n47 (Unlicensed) | $5855 - 5925$ | TDD | 15, 30, 60 | 10, 20, 30, 40 |

their own reservation, which reduces packet collisions. However, it comes at a cost of higher energy consumption due to continuous sensing.

### 2.3.4 NR V2X (Release 16)

To support a wide range of V2X applications with different quality of service requirements and support scenarios with high vehicular density, 3GPP has continued the standardization efforts on V2X communications through NR V2X in Release 16 and 17.

The requirements agreed for 5G V2X services and to be met by 3GPP standards are described in [85], where design requirements for 25 different 5G V2X use cases are presented. Thanks to the flexibility provided by 5G NR and the recent progresses envisioned in NR V2X, the support for a wide range of applications is feasible with NR V2X technology. The initial NR V2X design was developed in NR Release 16 SI [92], and was then included in the NR Release 16 specification based on the NR V2X WI [93]. Like IEEE 802.11bd and 5G NR, NR V2X also considers the use of mmWave bands for V2X applications, particularly for applications that require a short range and high to very high throughputs. However, considering the limited timeline of 3GPP Release 16, NR V2X mmWave operations were deprioritized in the 3GPP WI [93]. In this line, TR 38.885 conducted a limited study on beam management and concluded that it is beneficial for sidelink, but also that in sub 6 GHz bands it is feasible to support V2X use cases without sidelink beam management. Table 2.4 presents the sub 6 GHz operating bands for NR V2X, which is obtained from tables in [97] and [98].

The services specified for NR V2X range between 25 Mbit/s and 1 Gbit/s for data rate, 90 % to 99.99 % for reliability, and 5 ms to 100 ms for latency, depending on the use case [90]. Those latency requirements can not be met by Release 14 LTE C-V2X, but they can be improved considering higher numerologies in NR. Also, the reliability requirement of 99.99 % requires that NR V2X standardizes new enhancements at both resource allocation and scheduling. Extensive details of NR V2X will be given in the Chapter 7, Section 7.2.

## 2.4 Conclusions

In this chapter, we have discussed the technical background of the RAN technologies and the procedures that are the main focus of this thesis. In Section 2.1, we have provided an overview of the intra-LTE handover procedure, highlighting different handover events, measurement quantities, and the parameters that control the handover initiation. We

have shown that irrespective of the handover event type that triggers UE handover measurements, the criterion to select the next eNB is solely based on the strongest signal strength. This approach is maybe easy to implement in actual products. Still, it is not enough to handle challenging propagation scenarios, e.g., the one presented and deeply investigated in Chapters 3 and 4 of this thesis.

Then, Section 2.2 has covered the background of the unlicensed technologies, specifically, the LAA and LTE-U. Between them, LAA is believed to be a technology whose channel access mechanism is like WiFi; therefore, it could provide better fairness. However, based on the details of the LBT procedure of LAA provided in this chapter, one could easily spot the differences between the LAA and WiFi LBT procedures. For example, LAA's method to increase its contention window size could hamper its fair coexistence with WiFi, as this thesis explains in detail in chapter 5.

Lastly, Section 2.3 has introduced V2X technologies, emphasizing mainly those based on the 3GPP Sidelink/PC5 interface. We have highlighted that a V2X UE, based on release 14 and onwards, continuously senses the medium to perform sensing-based semi-persistence scheduling. However, it consumes more energy than random resource selection that does not perform sensing. Moreover, given real-world dynamic scenarios, where the traffic density may vary over time, a UE might not always need to sense the medium continuously. Therefore, it leaves space for a solution that could dynamically regulate the sensing based on the surrounding environment of the UE. One such solution is presented in Chapter 8 of this thesis.

Technology
complexity axis

IMT

MT

IST

LIST

ST

I

LI

Infrastructure support
complexity axis

# Chapter 3

# Machine Learning Based Handover Management

This chapter serves as an initial proof of concept for the wider research conducted on the HO management in chapter 4 of this thesis. In particular, in this chapter, we focus on the HO management problem in cellular networks, and we argue that an ML-based solution would be a beneficial approach to follow, since in this kind of problem, the experience of other users in similar propagation conditions, could be highly useful to make the right decision. We address the validation of the concept first in a simple scenario with 3 cells, and when this is tuned and validated, we extend the solution to a wider multi-cell more realistic cellular network. We also extend the neural network model to address scalability problems in wider scenarios. The more general solution is discussed in chapter 4. In this chapter, we focus on the simple scenario, providing some interesting simulation results, which makes it evident that the SOTA handover algorithms are not sufficient to tackle challenging propagation scenarios. To solve this problem, we present a smart handover management solution, which could enhance the target cell selection capability of a handover algorithm, taking into account the user's perceived QoE. In particular, the handover algorithm learns from its past experience, by using machine learning techniques, how the handover decision to a specific cell influences the QoE of the user. According to our approach, the serving eNB gathers some measurements reported by the UE, which provide information about the radio link conditions of the serving and neighbour eNBs, as well as the QoE of the UE resulting from the past handover decisions. We use a supervised learning approach based on a simple neural network to predict the most appropriate cell for handover. After training has been accomplished, the handover algorithm is able to select a target cell for handover that could provide a better QoE

despite an initially weaker signal upon handover decision. As already mentioned, it is worth repeating that the results presented in this chapter serve as the basis for our final and more complete HO scheme detailed in the next chapter.

The rest of this chapter is organized as follows: In Section. 3.1 we discuss the related work. Section. 3.2 presents the details of the new features developed in *ns-3* and the simulation scenario built using these features. Section. 3.3, presents the system description and the technical specification of the proposed scheme. Then, Section. 3.4, provides the details of the neural network setup and the results of the performance evaluation of the proposed scheme. Finally, Section. 3.5 concludes the chapter.

## 3.1  Related work

Recent surveys on the application of ML learning in mobile networks [99] and their self-organization [100] show that the ML-based solutions would play an essential role in the management of 5G and beyond networks. In the context of HO management, we identify three high-level potential ways to optimize its functionality. The first approach to optimize the HO process is to use model-based solutions based on Markov Decision Process (MDP). The objective, in this case, is to find the probability distribution of taking optimal HO decisions given the input state, which corresponds to a UE state before the HO. In [101], the authors proposed a Viterbi algorithm to find an optimal HO policy to maximize the UE average capacity. The algorithm works under the assumption that the position of the eNBs, the UE trajectory, and the channel characteristics are known a priori. The work in [102] proposes a cell selection procedure based on the Partially Observable Markov Decision Process (POMDP). Specifically, the POMDP predicts the neighboring cells' loading information to optimize the HO rate while maintaining the system throughput. Authors in [103], similarly to [101], proposed a context-aware HO policy, which optimizes the Time-to-Trigger parameter for HO by assuming the knowledge about the UE trajectory. These proposed solutions are based on the assumptions of having strong knowledge about network dynamics, which in turn are hard to capture in real networks. Therefore, model-free solutions which optimize the HO process without this previous and complete information are worth investigating.

The second approach considers then model-free solutions for HO parameter tuning. The idea is to adaptively fine tune the HO parameters defined in the standard to identify the strongest target cell, e.g., Hysteresis, Time-to-Trigger, HO Margin, and Cell individual Offset, by employing ML algorithms. In [21], a Q-learning approach is proposed to optimize the HO parameters. In particular, the model finds the optimal values of Hysteresis and Time-to-Trigger parameters to reduce the radio link failure and ping pong effects. In [22], a method to adaptively select a Hysteresis value to reduce the number of unnecessary HO is proposed. Specifically, it uses a predefined threshold value of RSRQ to adapt the Hysteresis value as per the UE measurements. Authors in [23], proposed a fuzzy logic controller, which finds an optimal value of the HO Margin parameter to reduce the signaling cost caused by HO.

These approaches that aim to select the strongest cell, based on the optimal tuning of HO parameters, have the shortcoming of considering the strongest signal for target cell

selection before the HO. Furthermore, these schemes do not consider a long-term vision of performance indicators in the decision, in terms of, e.g., the perceived QoE, after the HO. For example, in urban scenarios where the HO to the strongest neighbour cell is successful, but shortly after the transmission is deeply affected by the presence of an outage, these HO approaches could fail to provide a satisfactory solution. Thus, they are likely to severely degrade QoE performance, due to the unpredicted cell outage.

As a result, the third approach to HO management, which is also the one considered in this thesis, is a data-driven approach. It aims at using experience extracted from network data to include the vision of long-term optimization in the HO management decision. In [24], the authors proposed a hybrid HO controller based on deep reinforcement learning to minimize the HO rate while maintaining a certain level of system throughput. In particular, the work uses a Deep Neural Network (DNN), composed of LSTM units, which are trained following the supervised learning approach to predict the probability of selecting a target cell. It uses a dataset consisting of RSRQ measurements by simulating a standard compliant HO algorithm before executing the reinforcement learning approach. Similarly, in [25], a DNN is trained to solve the multiclass classification problem. In particular, it uses the RSRP measurements reported by the UEs to their serving next-Generation Node B (gNB). Then, using the softmax activation function for the output layer of the trained model, it computes the probability for a neighbouring gNB to become the next serving gNB. The smart HO approach to select the next serving eNB presented in this thesis, is different from what is proposed in [24], [25]. Specifically, it uses the regression to predict the perceived QoE (i.e., file download time) for each potential target eNB, and it handovers to the one, which could provide a better QoE. Moreover, the inputs to our model, thanks to deep ML architecture, include not only RSRP and RSRQ, but also many other measurements from the whole protocol stack, as it will be discussed in the next chapter.

## 3.2   Simulation scenario implementation

To support the design of these advanced mobility management solutions, it is important to be able to simulate handover scenarios with the complex propagation conditions to explore better the limitations of SOTA HO algorithms and to evaluate the performance of alternative candidate algorithms. We base this study on the LTE module of *ns-3*. It includes key aspects such as handover, fractional frequency reuse and support for simulating the buildings in a scenario. However, when conducting our research, it was not possible to use LENA to simulate coverage holes. Additionally, the preliminary evaluation of a machine learning approach to mobility management would be more easily performed by evaluating offline a large set of different alternative handover decisions; unfortunately, this is not possible with the handover algorithms currently implemented in LENA, which always select the same target base station when facing the same handover conditions. To overcome these limitations, in this thesis, we propose the following contributions:

1. A model for the simulation of obstacles potentially blocking the propagation of radio signals.

**Figure 3.1:** Sequence diagram for calculating obstacle Path Loss

2. A deterministic handover algorithm that can be used for the Offline-Evaluation of learning-based handover algorithms.

3. Two simulation scenarios for handover in the presence of a coverage hole due to an obstacle.

The source code of this module and simulation is publicly available at[1] for use and reproduction of results.

## 3.2.1 New Developed Features in ns-3

This section describes the newly implemented features in *ns-3* for this study.

### 3.2.1.1 Obstacle Model

Our approach to implement the obstacle model and to achieve the desired behavior of coverage holes is to use the existing "Buildings" module of *ns-3* [104]. The structure of these buildings is defined by a 3D axis-aligned box defined by the "Box" class. By doing so, we take advantage of their current functionalities, e.g., maintaining the list of all the obstacles in a simulation, assigning the unique id to an obstacle and using the function "IsInside", provided by the Box class, which indicates whether a node is inside the box or not.

The coverage holes are simulated as follows: if the line segment between the two nodes intersects the box, or any node is located inside the box, then the transmission is attenuated by adding a significant propagation path loss value. To implement this

---

[1]https://github.com/ZorazeAli/ns-3-dev-obstacle

**Figure 3.2:** Logic to trigger deterministic handover

**Figure 3.3:** REM for simple scenario

approach, we need to check whether the signal between the two nodes is blocked by building or not, but *ns-3* did not provide any function for this. Therefore, we extend the functionality of the "Box" class to check the intersection of a line segment between two nodes and the box, based on [105], considering a 3D box. Based on this new feature, we implemented an obstacle path loss model in this thesis that is inherited from the "propagation-loss-model" class to chain our path loss model with existing path loss models. This new path loss model iterates through the list of all the obstacles created during the simulation to check for an intersection between the line segment and the box. Upon the intersection, it adds the path loss and returns the received power, as shown in Fig. 3.1.

## 3.2.2 Deterministic Handover Algorithm

In this thesis, we have implemented the A2 event-triggered deterministic handover algorithm. The purpose of this deterministic handover algorithm is to allow the offline evaluation of the performance of different handover algorithms for each possible target eNBs. The important point while implementing any handover algorithm in *ns-3* is the configuration of the UE measurement report for the handover. Our handover algorithm is based on the Reference Signal Received Power (RSRP), which is configurable by setting the *triggerQuantity* attribute of the UE measurement report. The RSRP threshold value, for which the A2 event is triggered, is also configured through this report. When the algorithm receives the A2 event-based measurement report from the UE, it calls the function *EvaluateHandover*. This function first checks the availability of neighbor cells for handover. If it finds any neighbor and its information e.g. the cell id and RSRP, it checks if the *TargetCellId* is assigned by the user or not. If it is, the *BestNeighborCellId* is set equal to the user-defined target cell id, and the handover is triggered. In case the user does not define the target cell id, it triggers the handover in a non-deterministic way toward the strongest neighbor, as shown in Fig. 3.2.

**Table 3.1:** Simulation network parameters for the simple scenario.

| Parameter | Value |
|---|---|
| System bandwidth | 5 MHz |
| Inter-site distance | 500 m |
| Handover algorithm | A2-RSRP |
| Adaptive Modulation & Coding Scheme | Vienna [106] |
| SINR computation for DL CQI | Control method [106] |
| eNBs antenna type | Isotropic |
| Number of macro eNBs | 3 |
| eNBs Tx Power | 46 dBm |
| Number of UEs in the system | 3 |
| Mobility model | RandomWalk2dMobilityModel Mode: Time, Time: 100 sec, Distance: 4000 m |
| Path loss model | Cost231 |
| eNB Antenna height | 30 m |
| Obstacle height | 35 m |
| Traffic | TCP Bulk File Transfer |
| File size | 15 MB |
| Simulation time | 100 sec |

## 3.2.3 Simulation scenarios

This section provides the details to configure the essential parts of our simulation to achieve the desired scenario topologies. Specifically, we build two scenarios, 1) a simple scenario that is used for a proof of concept study related to the HO management, presented in the next section, and 2) an extensive simulation scenario that is used to build a wide and complete database, which we consider the basis of the experience that a smart network management solution able to construct to take smart HO decisions in such a scenario. As mentioned above, in this chapter, we focus on the research work performed using the simple scenario; therefore, in the following, we present only the details about this scenario and leave the extensive scenario to be discussed in the next chapter.

### 3.2.3.1 Simple scenario

The simple scenario topology consists of three macro eNBs, 3 UEs and an obstacle partially obstructing the coverage by eNB2, as shown in the REM in Fig. 3.3. Each eNB is serving one UE performing TCP file download from the remote host, where UE1 (initially attached to eNB1) is moving around, and the other 2 UEs are stationary. The mobility of UE1 is generated by tuning the parameters (see Table. 3.1) of the *RandomWalk2dMobilityModel* in *ns-3* in such a way that it picks a fixed starting point close to eNB1 and a random angle to move away from the source eNB following the straight line. The simulation consists of 200 runs of a deterministic handover. Each run is repeated twice, first targeting eNB2 and then eNB3 to measure the QoE. For every simulation run, UE1 picks a fixed starting point close to eNB1, and a random angle in the range of $[+X, -2X]$ to move away from the source eNB following the straight line, where $X$ is the angle from eNB1 to eNB3. The rationale behind selecting TCP file download is

**Figure 3.4:** The proposed two level Neural Network Scheme

that TCP is the widely used transport protocol for many interactive internet applications, e.g., the Web. Secondly, investigating the QoE of the user under challenging handover scenarios while introducing a significant load to the network through file download could yield true insights of the mobility management capabilities of the network. The complete set of simulation parameters is described in Table. 3.1.

## 3.3 FFNN based HO management

### 3.3.1 System description

As we stated above, our primary objective is to design a scheme which can enable the handover algorithm to understand whether the target eNB for handover would be able to provide a consistent QoE to the user or not. As a result, the handover algorithm could identify those eNBs which are affected by the undesirable radio propagation scenarios in the network, e.g, coverage unavailability of an eNB caused by an obstacle.

Our proposed scheme, which is depicted in Fig. 3.4, consists in the following: the source eNB gathers the time series of UE measurement reports before the handover, which contains the RSRP and RSRQ of the source and neighbour eNBs. The eNB also collects the information on the QoE of the user as a result of past handover decisions. In our scheme, this QoE is quantified by two metrics, 1) the probability of successfully downloading a file and 2) the file download time for completed downloads. Therefore, we propose to use a two level neural network model to estimate these metrics, as shown in Fig 3.4. At level 1, the first neural network (NN1) is trained using UE measurements as input, and the past QoE in terms of download complete/not complete as output. On the other hand, at level 2, the second neural network (NN2) is trained using as input only those UE measurements for which the file download was completed. Finally, the file download time is the output of the NN2. We propose to use two single-output neural

networks instead of one multiple output neural network, as it is proven that this leads to better results [107]. Moreover, to train both the models , i.e., NN1 and NN2, we use the offline-training approach explained in Section. 1.3.2. Once the training is completed, the handover algorithm of the source eNB uses these two trained NNs to determine the expected QoE to be achieved through all the potential target eNBs. The handover algorithm then triggers the handover to the target eNB for which the file download is expected to finish successfully, and in case, there are two or more potential target eNBs, it handovers to the eNB with the lowest value of the estimated file download time. Related to the selection of time to download as a QoE metric, it is motivated by the fact that it is one of the standardized Key Quality Indicator (KQI) for the file transfer service in mobile networks [108] [109]. In the literature, some models are based on subjective Mean Opinion Score (MOS) to derive the QoE; however, all of them depend on end-to-end throughput perceived by the users [110]. On the other hand, following a similar methodology taken in [111], we take a more generalized approach, i.e., instead of using a specific MOS model for file transfer service, time to download has been used as an indicator of the QoE perceived by the users.

### 3.3.2   FFNN design

In this subsection, we provide a brief overview on important technical specifications of our neural network. For a more detailed description, the reader is referred to the vast available literature on neural networks (for instance, [107], [67], [66]).

For the implementation of our proposed scheme we use a FFNN with single-hidden layer [67], also known as two-layer FFNN, where the number of layers refers to the number of layers with adaptive weights. We choose FFNN because of its ability to model both linear and non-linear functions between inputs and outputs. Additionally, the model obtained with FFNN is more compact and fast to evaluate than other machine learning techniques such as, support vector machines, with the same generalization performance [66]. In general, when working with FFNN with supervised learning [67], such as in our case, one has to build a training database of input and output vectors stored in rows and columns, also known as dataset. This dataset can be represented as,

$$D = \{(X_1, Y_1), \dots, (X_M, Y_M)\} \tag{3.1}$$

where $X$ and $Y$ are the input and output vectors of FFNN and $M$ is the total number of rows in a dataset. Let the index $r$ denote the row number of our dataset and let $t$ denote the time at which the UE measurement report was received. The input vector $X$ can be written as,

$$
\begin{aligned}
X_{r=1} &= [x_1(t), x_1(t-1), \dots, x_1(t-L+1)] \\
&\qquad\vdots \\
X_{r=M} &= [x_M(t), x_M(t-1), \dots, x_M(t-L+1)]
\end{aligned}
\tag{3.2}
$$

where $L$ is the memory of the Neural Network. In Eq. 3.2, $\text{x}_r(t)$ is one UE measurement report received at time $t$, which can be formulated as follows,

$$x_r(t) = [P_1(t), Q_1(t), P_2(t), Q_2(t) \dots, P_N(t), Q_N(t)] \tag{3.3}$$

**Figure 3.5:** Architecture of Neural Network with 12 inputs, 1 hidden layer, 2 bias and 1 output

$P_i(t)$ and $Q_i(t)$ are the RSRP and RSRQ of cell $i$, respectively, and i = 1, ...,N, where N is the total number of cells (i.e, serving cell plus the number of neighbour cells). Finally, the output vector $Y$ for NN1 and NN2 contain the values of the QoE metric. For NN1, these values are stored in the form of logical values of 0 and 1, where 0 indicates "download not complete" and 1 indicates "download complete". On the other hand, for NN2, these values are the download time in seconds for the completed downloads. The task faced by NN1 is a classification problem where the FFNN estimates to which class (0 or 1) the given input belongs. Therefore, we choose the softmax function [66], as an activation function of the output layer of NN1. By using softmax activation function, we force the coupled output of the FFNN to sum to 1, so that they represent a probability distribution across discrete mutually exclusive alternatives. On the other hand, the task faced by NN2 is a regression problem, to estimate the file download time. In this case, we use a logistic function also known as logistic sigmoid activation function [66].

The purpose of using the UE measurements as input is due to the fact that, these timely reported UE measurements change according to the UE position. So, if the UE is moving towards any of the available target eNBs they can provide the information about the possible UE trajectory. Therefore, by training the FFNN with these measurements as input and the QoE metric as output, the FFNN will learn about those mobility patterns which caused the degradation of QoE after the handover is executed. We note that 3GPP standards already contemplate the upload of these UE measurements, as specified for the Minimization of Drive Test (MDT) [69].

**Figure 3.6:** Average Download Time for Completed Downloads (a), Number of Completed and Uncompleted Download Attempts(b) vs Range of Angles.

## 3.4 Performance Evaluation

### 3.4.1 Implementation of the Neural Network

For the implementation of FFNN, we used a publicly available nnet package of R [112], which is a single hidden layer FFNN. The dataset D, containing the UE measurements, is randomly divided into a training set (containing 75 % of the data) and a testing set (containing 25% of the data). All the input and output values are normalized in the range [0,1]. This normalization allows for a faster training process and more accurate estimations [113]. Fig. 3.5 shows the implemented FFNN. The structure is based on a single hidden layer of 4 neurons and 12 neurons in the input layer. Referring to Eq. 3.2 and 3.3, here we consider N=3 and L=2. We fix the maximum number of iterations to 1000 for the training phase. We note that, depending on the complexity of the FFNN one should choose a maximum number of iterations to avoid the early stop of the training process before the algorithm converges. Moreover, we fix the weight decay parameter to 0.0001 and perform a 3-fold cross-validation test that uses Normalize Root Mean Square Error (NRMSE) to select the final FFNN model. It prevents over-fitting, i.e., when FFNN achieves the ideal minimization of the error between the estimated and the actual output of the training set. In this situation, the FFNN loses its generalization property and fails to predict the output of the testing dataset. Another important factor to keep in mind while using FFNN with R is the randomness. For every new seed, the weights of the FFNN take random initial values, so the performances may vary. Additionally, every seed results in different partitions of the dataset into training and testing sets. To account for this randomness, we average the results over 100 seed values to attain a statistically significant evaluation.

**Figure 3.7:** Performance comparison of FFNN vs. SOTA and an optimum scheme for $[-60°, 30°]$

## 3.4.2 Results

We first present some preliminary quantitative results from the deterministic handover campaign as described in subsection 3.2.3.1. The reason for providing these results is to support our argument for the need of a smarter handover algorithm with respect to SOTA approaches in challenging propagation scenarios. To evaluate the performance of the handover algorithm, as a function of the angle with which the UE crosses the affected outage region, the simulation area is divided into three ranges of angles, as shown in Fig. 3.3. Fig. 3.6 (a) and (b), respectively, show the results of the average download time for completed downloads and the number of completed/uncompleted downloads for the complete range of angles. Here we assume that the download time of incomplete downloads is equal to the maximum simulation time, i.e., 100 sec. From Fig. 3.6(a), it can be observed that for the range $[+30°, 0°]$, none of the downloads get completed when the handover is done to eNB2, as the UE experiences poor channel quality due to its distance from the source cell and the high exposure to the affected coverage zone. On the other hand, in the range of $[-30°, -60°]$, none of the downloads is completed when handover is done to eNB3, as in this case as well, the UE experiences poor channel quality due to its distance from the source cell and the small coverage outage area due to the obstacle between eNB3 and UE. As we can notice, the angle range $[0°, -30°]$ is the range where the handover decision has a very high impact, as SOTA handover algorithms based on the A3 event would provide eNB2 as the strongest candidate for handover, and doing so may cause an increase in average download time and the number of incomplete downloads. On the other hand, a handover to eNB3 would not only decrease the average download time but would also decrease the number of incomplete downloads.

We present the performance evaluation of our machine learning-based handover scheme, compared to the SOTA approach and to an optimal handover scheme, which always selects the best eNB to download the file successfully, with the lowest possible delay. In this way, we can see how close our handover algorithm's performance is to the optimal

**Table 3.2:** Comparison of handover schemes for $[-30°, 0°]$

| Handover Scheme | Completed Downloads(%) | Avg.Download Time(Sec) |
|---|---|---|
| SOTA | 54.48% | 50.51 |
| Proposed | 95.37% | 42.51 |
| Optimum | 100% | 42.39 |

one and how much improvement has been achieved by our scheme over SOTA handover algorithm. We note that, here we follow the Offline-Evaluation procedure described in Section. 1.3.2. As shown in Fig. 3.7, the performance of all the schemes in terms of file download time follows the same trend till angle $-24°$. After this angle, the download time for the SOTA handover scheme starts increasing, as it keeps seeing eNB2 as the strongest neighbour. The UE, though, gets more exposed to the affected coverage zone and experiences loss of data and huge delays due to long TCP timeouts, which finally leads to a high percentage of incomplete downloads. On the other hand, the machine learning-based scheme performs very well in the same range of angles, with similar trends to those shown by the optimal handover scheme. We only appreciate sporadic incomplete downloads due to some loss in accuracy of the FFNN. From Fig. 3.7, we can also observe that the divergence in the performance between the handover schemes occurs in the range of $[-30°, 0°]$. Therefore, in Table 3.2, we show results in terms of completed downloads and average download time only for this range of angles. From the results, we can observe that using the ML-based handover scheme, we achieve a 75% increment in the number of completed downloads and a decrease of 15.84% in file download time, with respect to the SOTA handover scheme. To summarize, the performance of our handover scheme is better than the SOTA handover scheme in the challenging propagation scenario presented in this study, thanks to its ability to learn all the UE mobility patterns that affect the QoE of the user.

## 3.5 Conclusions

This chapter presented a proof-of-concept study targeting the complexity encompassing infrastructure-based single-technology axes. Specifically, we have presented a simple machine learning-based handover scheme for improved QoE in LTE scenarios with challenging propagation conditions, e.g., in the presence of obstacles in the coverage area of eNB. Although using a simpler scenario compared to the one presented in the next chapter, this scheme was proposed to answer our first research question (RQ1): "How to use AI/ML in mobility management to achieve better QoE?". The proposed solution uses a two-level FFNN for the implementation of learning capabilities. Using our scheme, the handover algorithm can select the eNB that is expected to yield better QoE, based on the experience gained from past handover decisions. Our performance study showed that our scheme could achieve performance close to the optimal one in challenging scenarios. Therefore, it substantially improves QoE in terms of the number of successful downloads and average download time with respect to SOTA handover schemes, which make decisions based on signal strength, e.g., A3 event-based handover algorithms.

# Chapter 4

# Multi-Task Learning for Efficient Management of Beyond 5G Radio Access Network Architectures

In the previous chapter, we focused only on the single task HO management use case, and used a simple FFNN to solve the regression problem and estimate the QoE of a user. The obtained results proved that the learning approach outperforms traditional HO solutions. This chapter further extends our research horizon, in which we use the multi-task learning paradigm. To this extent, the ML-based models proposed in this chapter address the HO management use case and also tackle the second use case of the initial MCS, targeting a more realistic and heterogeneous simulation scenario. In particular, our objective is to prove the potentiality of MTL to address RAN automation of multiple tasks, which have to function in parallel during the regular operation of the RAN. We propose different deep architectures to address a set of tasks through individual or shared models. Among them, we study the effectiveness of Auto Encoder (AE) [114] to reuse the compressed representation of the data for multiple heterogeneous use cases [115] [116]. This approach can significantly reduce the implementation and computational complexity of the learning architectures. To prove this concept, we propose to target, without loss of generality, two RAN use cases: 1) HO management and 2) the selection of the optimal initial MCS.

We address both use cases, first through *single-task individual* and then through *multi-task shared* models for the sake of complete comparative study. To address the use cases individually, we use an LSTM RNN to take advantage of the temporal characteristic

of the data extracted from several and extensive simulation campaigns using the LTE module of *ns-3* [52]. The LSTM is designed to solve a regression problem to estimate the QoE of the users. We obtain excellent prediction errors, and with these results, we can prove that the learning approach outperforms traditional HO solutions and conservative approaches to select the initial MCS. Successively, we build a different architecture where the two RAN tasks share and train a common AE, based on MTL principles. The same compressed data output of the AE is used as input to two different MLPs, implementing the regression of the particular parameter that we want to estimate for the two use cases (i.e., the HO management and initial MCS). Both MLPs offer excellent regression results similar to the one obtained using LSTM. This means that the AE successfully reduces the dimensionality of the data without losing meaningful information and network performance, which consequently facilitates the sharing of knowledge between different tasks. Therefore, the same architecture can be used to address multiple parallel RAN uses cases.

In addition to that, we go more deeply into the study by comparing two different ways of learning. In the first case, the *parallel MTL* case, we learn the shared model by building a shared database for all the use cases we plan to address. This type of database is viable when we know the needed use cases beforehand. However, RAN management problems can be more complex and continuously require adding new use cases and on-demand tasks to the design without retraining previous tasks from scratch, or compromising their performance. As a result, there is a need for the MTL shared model to be flexible and be able to gradually add more tasks to its knowledge without forgetting previously known tasks. For that, we also propose a second *incremental MTL* scheme, based on the continual learning paradigm [117], where the training database is not built beforehand, but a new task can be incorporated separately while previous task knowledge is preserved. This approach is much more flexible and adequate for real networks and provides clear implementation advantages [118]. Finally, the contributions of this chapter are the following:

- Design of ML models based on single-task and multi-task paradigms to address two RAN use cases[1]. In particular, we propose two models based on two different multi-task techniques, i.e., parallel and incremental learning.

- Performance evaluation of the proposed solutions by comparing the results with 3GPP standardized HO and initial MCS selection schemes.

- To encourage the reproducibility of the proposed models and results, we provide in-depth details of the simulation scenario and steps to create the databases using an open-source simulator *ns-3*.

Last but not least, building upon the above contributions, we are advancing the SOTA and related work in the following aspects:

---

[1]In this study, we select the two RAN use cases to demonstrate the applicability of the proposed ML models in a concise manner. Therefore, these models can be used to address other or more than two RAN use cases.

- We present a holistic solution to handover based on download time that is not limited to adjusting typical HO parameters but considers previous experience to select a target eNB to improve users' QoE.

- Extending our previous work in chapter 3, we introduce an additional RAN use case, i.e., initial MCS selection. Compared to the studies in the literature that lean towards treating these two use cases separately, in this chapter, we study the solutions that allow learning concurrently these RAN tasks based on the MTL learning paradigm.

- We study two possible solutions for MTL in the context of solving RAN problems, one based on parallel learning and another based on incremental MTL, which increases the learning pace and minimizes the delay to deploy RAN solutions.

- We study the performance based on the proposed solutions in comparison to 3GPP standard aligned baselines.

The rest of this chapter is organized as follows. In Section 4.1 we discuss the related work. Section 4.2 introduces the system overview and the scenario that we use for synthetic data generation. Section 4.3, presents the procedure we adopted to generate the synthetic data. Section 4.4 proposes the RNN models for single-task and multi-task learning. Section 4.5 discusses the training of the proposed architectures and the system level performance results. Finally, Section 4.6 concludes the chapter.

## 4.1 Related Work

In Section 3.1 of the previous chapter, we already extensively reviewed the literature on HO management and classified it into three categories. On the other hand, to the best of our knowledge, related work for the initial MCS selection is limited to the usage of reinforcement learning. For example, [119] proposed a solution, that learns the best MCS given the SNR at a specific channel state. Differently from [119], the authors in [120] present a solution that uses Channel Quality Indicator (CQI) as a metric for the state representation. In this case, the authors argue that fine discretization of SNR with discrete MDP leads to a higher state space, which increases the convergence and exploration time. The authors in [121] proposed a deep reinforcement learning approach to overcome the issue of a large state space highlighted by [120]. They used SNR, SINR, the previous action, and its immediate reward for the state representation.

In this thesis, differently from the above presented solutions in [21–23] and [24, 25, 101–103, 119–121], we aim to demonstrate that different RAN use cases, such as HO and initial MCS selection, can be considered as related tasks. Therefore, these tasks can be jointly trained through shared models so that each task can benefit from other auxiliary tasks. Such an approach offers multiple implementations and learning advantages like reduced training effort, improved data efficiency, reduced overfitting through shared representations, and fast learning by leveraging auxiliary information. MTL has already been recently considered in mobile communications literature. In [122], multi-task Sparse Bayesian Learning (SBL) is applied for learning time-varying sparse channels in the uplink

for multi-user massive MIMO systems. Results show that it is possible to considerably reduce the complexity and the required time for the convergence with a negligible sacrifice of the estimation accuracy. In [123], MTL is used to train a shared model for both traffic classification and prediction at the edge of the network. Classification accuracy and prediction error benefit from the shared model and return better performance with respect to single-task neural network architectures. In [124] multi-task DNN framework for Non-Orthogonal Multiple Access (NOMA), namely DeepNOMA, has been proposed to treat non-orthogonal transmissions as multiple distinctive but correlated tasks. To the best of the authors' knowledge, this work is the first one in literature deeply discussing and proving the concept of MTL for the efficient automation of the RAN in future mobile networks, improving so the ability to make connections between facts, observations, patterns, and other tasks from which they learn.

## 4.2 System Overview

This section first describes the RAN use cases that we handle using deep learning solutions. Then, we introduce the target simulation scenario, which is depicted in Fig. 4.1 [125].

### 4.2.1 Target RAN use cases

Considering the high-level objectives, we have selected two RAN use cases to be addressed using a supervised learning approach: 1) HO management and 2) initial MCS selection. In a traditional network, these use cases are handled at Layers 3 and 2 and are treated separately. For this reason, one of the motivations behind selecting the second RAN use case (i.e., initial MCS selection) is to evaluate the efficiency of MTL-based models that effectively use our proposed dataset, spanning multiple layers, to learn these tasks jointly.

- *HO management*: We propose a HO management approach, which allows HO to the cell suggested by a supervised learning algorithm capable of predicting a QoE indicator through a regression procedure. The supervised learning algorithm exploits the experience extracted by data already available in the network (e.g., the Minimization Drive Test database [69]). Based on this, it detects the most appropriate cell to HO, as a function of the future expected QoE perceived by the user, instead of the RSRP or the RSRQ as the standard suggests. We model the problem as a regression problem, where we aim to estimate the necessary time to download a file transmitted over a Transmission Control Protocol (TCP) transport, while the users move around in a realistic multi-cell scenario challenged by deep outage zones. Finally, it is to be noted that this solution has to be considered a component of a more sophisticated HO algorithm that also includes other aspects, e.g., load balancing, QoS requirement of a UE. However, we think that including these additional components is out of the scope of this thesis and, undoubtedly, would be interesting for future work.

**Figure 4.1:** REM for extensive scenario

- *Intial MCS selection*: A standard approach in cellular networks is that when the UE first switches to the CONNECTED RRC state, the initial selected MCS follows a conservative approach that guarantees that the initial transmissions go through. As a result, the lower MCS is usually selected (i.e., 0). We propose using knowledge from data reported by the users to choose an initial MCS in an optimal and non-conservative way to avoid wasting radio resources in the initial data exchange.

## 4.2.2 Simulation Scenario

Generally, a more realistic outdoor cellular scenario is more complex because it consists of several eNBs and UEs, and maybe, more than one obstacle when compared to the simple scenario that we presented in the previous chapter (Subsection. 3.2.3.1). However, an extensive simulation scenario can help build the ground to justify the use of AI/ML-based RAN solutions in real networks. Therefore, in this chapter, we consider a complex and realistic simulation scenario.

In this scenario, we consider a macro cell outdoor scenario, but different from the simple scenario, it is a network consisting of three-sectorial eNBs. A cluster of UEs is placed in each sector at a fixed distance from the centre of a cell, in which the UEs are dropped at random positions. Since, in this scenario, we use TCP as the transport protocol, such deployment of the UEs guarantees to establish a TCP connection between the remote host and the UEs. The UEs start moving after receiving the first packet, following a mobility pattern resulting from configuring the parameters of the *RandomWalk2dMobilityModel* in *ns-3*, as we did in the simple scenario. In particular, for every simulation run, a UE picks a random starting position in the cluster and a random angle in the range of [0° to 360°] to move away from the source eNB following a straight line. To increase the communication challenges in the scenario and to generate more random coverage

**Table 4.1:** Simulation network parameters for the extensive scenario.

| Parameter | Value |
|---|---|
| System bandwidth | 5 MHz |
| Inter-site distance | 500 m |
| Handover algorithm | A2-RSRP |
| Adaptive Modulation & Coding Scheme | Vienna [52] |
| SINR computation for DL CQI | Control method [52] |
| gNBs antenna type | Parabolic |
| gNBs antenna Beamwidth | 70 degrees |
| gNBs antenna max attenuation | 20 dB |
| Number of macro gNBs | 21 (7 cells) |
| gNBs Tx Power | 46 dBm |
| Numerology | 0 |
| Distance between the center points of the UEs cluster and the cell | 100 m |
| UEs Cluster diameter | 50 m |
| Number of UEs in the system | 210 (30 per sector) |
| Mobility model | RandomWalk2dMobilityModel Mode: Time, Speed: 10 m/s Time: 40 sec, Distance: 4000 m |
| Path loss model | Cost231 |
| gNB Antenna height | 30 m |
| Obstacle height | 35 m |
| Traffic | TCP Bulk File Transfer |
| File size | 1.5 MB |
| Simulation time | 40 sec |

patterns, we introduce obstacles in the scenario, which create multiple coverage holes, as shown in Fig. 4.1 [126]. Each UE performs a TCP file transfer to a remote host in Downlink (DL) and Uplink (UL) direction. The complete set of simulation parameters is described in Table. 4.1 [125]. The above simulation scenario is then used to conduct three extensive simulation campaigns, two for the single-task approaches, i.e., the HO management and the initial MCS, and one for the multi-task approach jointly targeting both use cases. Each of them is repeated a specific number of times, which depend on the values of the parameters, i.e., the number of independent simulation runs, the maximum number of neighbours to HO, and the number of initial MCS values evaluated. The data obtained from these campaigns for each UE are stored in the form of a dataset, according to the format described in the next section (Section 4.3). We will explain in detail the use of this simulation scenario to build the databases for single and multi-task learning, targeting the two use cases.

## 4.3 Data Generation

This section describes the characteristics of the collected dataset that we use as input to our proposed deep learning solutions. We constructed this dataset by conducting extensive simulation campaigns in the scenario presented in Subsection 4.2.2. As mentioned at the beginning of this chapter, we model the HO management and initial MCS problems as regression problems, where we need to estimate, respectively, the QoE

**Table 4.2:** List of input and output features used to create the training and testing dataset.

| Input feature | | |
|---|---|---|
| Layer | Measurements | |
| APP | 1. Throughput UL<br>4. Throughput DL | 2. Avg. number of rcvd packets UL<br>5. Avg. number of rcvd packets DL | 3. Avg. number of rcvd bytes UL<br>6. Avg. number of rcvd bytes DL |
| RRC | 7. Cell ID of serving cell<br>10. Cell ID of neighbour 1<br>.<br>.<br>31. Cell ID of neighbour 8<br>34. Total number of radio link failures | 8. RSRP from serving cell<br>11. RSRP from neighbour 1<br>.<br>.<br>32. RSRP from neighbour 8<br>35. Total number of handovers | 9. RSRQ from serving cell<br>12. RSRQ from neighbour 1<br>.<br>.<br>33. RSRQ from neighbour 8<br>36. First target cell ID to handover |
| PDCP | 37. Total number of txed PDCP PDUs DL<br>40. Avg. PDCP PDU delay DL<br>43. Min. PDCP PDU size DL<br>46. Total number of rcvd PDCP PDUs UL<br>49. Min. value of the PDCP PDU delay UL<br>52. Max. PDCP PDU size UL | 38. Total number of rcvd PDCP PDUs DL<br>41. Min. value of the PDCP PDU delay DL<br>44. Max. PDCP PDU size DL<br>47. Total bytes txed UL<br>50. Max. value of the PDCP PDU delay UL | 39. Total bytes txed DL<br>42. Max. value of the PDCP PDU delay DL<br>45. Total number of txed PDCP PDUs UL<br>48. Avg. PDCP PDU delay UL<br>51. Min. PDCP PDU size UL |
| RLC | 53. Total number of txed RLC PDUs DL<br>56. Total number of bytes rcvd DL<br>59. Max. value of the RLC PDU delay DL<br>62. Total number of txed RLC PDUs UL<br>65. Total bytes rcvd RLC PDUs UL<br>68. Max. value of the RLC PDU delay UL | 54. Total number of rcvd RLC PDUs DL<br>57. Avg. RLC PDU delay DL<br>60. Min. RLC PDU size DL<br>63. Total number of rcvd RLC PDUs UL<br>66. Avg. RLC PDU delay UL<br>69. Minimum RLC PDU size UL | 55. Total number of bytes txed DL<br>58. Min. value of the RLC PDU delay DL<br>61. Max. RLC PDU size DL<br>64. Total bytes txed RLC PDUs UL<br>67. Min. value of the RLC PDU delay UL<br>70. Maximum RLC PDU size UL |
| MAC | 71. Initial MCS<br>74. Avg. MCS UL<br>77. Avg. RB occupied DL<br>80. UL CQI | 72 Avg. TB size UL<br>75. Avg. MCS DL<br>78. DL CQI inband | 73. Avg. TB size DL<br>76. Avg. RB occupied UL<br>79. DL CQI wideband |
| PHY | 81. Avg. SINR DL<br>84. Avg. number of UL HARQ NACKs | 82. AVG. SINR UL | 83. Avg. number of DL HARQ NACKs |
| Output feature | | |
| APP | 1. File download time [sec] | 2. Initial DL throughput over 100 msec when a new RRC connection is established after the second handover | |

expected from performing HO to a certain target cell, and the initial throughput obtained by the UEs over a certain window. In general, when working with supervised learning, such as in our case, one has to build a DataBase (DB) with enough data to train, test, and evaluate the model. This dataset consists of input and output features stored in rows and columns. For this purpose, we have identified features at the multiple layers of the simulator protocol stack. These features can bring information to address not only the targeted RAN use case, but also other RAN use cases that could be later considered. In particular, we have organized these features per layer of the 3GPP protocol stack, and presented them in Table 4.2 [125]. 3GPP already contemplates uploading a part of these measurements, e.g., UE measurements, under the Minimization of Drive Test (MDT) functionality [69]. All these measurements are gathered in the simulator, by leveraging the *ns-3* "tracing system", which enables us to write them in text files as an output of the simulation program.

Successively, we run multiple independent runs of the simulation scenario and then post-process all the generated text files to build a unique DB in *csv* format. The rest of this section describes the procedures to construct the DB for training and testing.

## 4.3.1   Procedure to build the database

For the purpose of evaluating and comparing single-task versus multi-task learning performances, we build four databases, two for the single-task approaches, i.e., targeting the two use cases individually, and two for the multi-task approach, considering the parallel and incremental MTL possibilities. In the following, we explain the pseudocode procedure to generate these databases.

### 4.3.1.1 Single-task HO management database (DB1)

In a real-world scenario, a UE served by a gNB could HO to different potential neighbor gNBs. It depends on the HO criterion. Examples of such criteria are the signal strength before the HO, reported using UE measurements, as traditionally proposed in standards, or the QoE after the HO, as proposed in this thesis. This decision is usually affected by the UE's mobility pattern. However, it may also happen that different mobility patterns lead to the selection of the same target neighbour, because it is in all the cases identified as the most suitable neighbour to HO to. A QoE oriented HO algorithm must take these aspects into account. Therefore, the simulation campaigns to build the first DB (**DB1**) consists of several deterministic HOs to learn the QoE, i.e., a file download time for each UE, for the possible mobility patterns. The procedure to generate **DB1** is illustrated with the help of Pseudocode 1. Specifically, to consider both aspects discussed above, the number of deterministic HOs to be performed by a UE of a gNB would depend on the number of independent runs used to generate different mobility patterns of this UE for each HO (first "for" loop of Pseudocode 1), and on the maximum number of neighbours this UE manages to see (last "for" loop of Pseudocode 1). Then, these deterministic HOs have to be simulated for every gNB (second "for" loop of Pseudocode 1) and every UE attached to a gNB (third "for" loop of Pseudocode 1) in our simulation scenario. The measurements resulting from these HOs will assist the proposed architecture in learning the most reasonable neighbour to HO to. For this DB, we collect the data focusing only on the HO management use case.

It is also worth mentioning that we engineered this deterministic HO procedure to collect a synthetic DB in a reasonable time. However, in a real network, it would be possible to collect real online measurements based on the realistic mobility of the UEs during their lifetime while the network is normally operating.

### 4.3.1.2 Single-task initial MCS database (DB2)

The second DB (**DB2**), targets the initial MCS (i.e., DL) use case. The logic to construct **DB2** is somewhat similar to **DB1**. The steps to generate this DB are presented in Pseudocode 2. In particular, in this case, we have to evaluate all, or a set of initial MCS values a gNB could use for a newly connected UE. Moreover, since the MCS depends on the SINR, which depends on the mobility of a UE, this DB should also consider this aspect. Therefore, a simulation for each MCS value (second "for" loop of Pseudocode 2) should be repeated for a number of independent runs (first "for" loop of Pseudocode 2) to record the QoE resulting from different mobility patterns of a UE. Similar to **DB1**, this has to be simulated for all the gNBs and their UEs in our simulation scenario (see, third and fourth "for" loop of Pseudocode 2). For this use case, the selected QoE indicator is the throughput achieved over a certain time window after a successful RRC connection establishment. This window's duration should be smaller than the configured DL CQI reporting interval of a UE, which is typically 200 ms, after which a gNB adapts the MCS based on the reported CQI.

---

Pseudocode 1: Pseudocode to generate **DB1**

---

initialization
$numGnbs \leftarrow 21$ ; $numUesPerGnb \leftarrow 10$
$numNeighboursPerUe \leftarrow 8$; $numRuns \leftarrow 20$

1  **for** $r \leftarrow 1$ **to** $numRuns$ **do**
2     **for** $e \leftarrow 1$ **to** $numGnbs$ **do**
3        **for** $u \leftarrow 1$ **to** $numUesPerGnb$ **do**
4           **for** $n \leftarrow 1$ **to** $numNeighboursPerUe$ **do**
             Start moving away from the serving gNB, based on a random direction.
             HO to neighbour cell $n$ as per event A2.
             Run for simulation time and collect stats with a configured measurement periodicity.

---

Pseudocode 2: Pseudocode to generate **DB2**

---

initialization
$numGnbs \leftarrow 21$ ; $numUesPerGnb \leftarrow 10$
$numMcsValues \leftarrow 3$; $numRuns \leftarrow 63$

1  **for** $r \leftarrow 1$ **to** $numRuns$ **do**
2     **for** $m \leftarrow 1$ **to** $numMcsValues$ **do**
3        **for** $e \leftarrow 1$ **to** $numGnbs$ **do**
4           **for** $u \leftarrow 1$ **to** $numUesPerGnb$ **do**
             Fix the Initial MCS of all the UE to $m$.
             Start moving away from the serving gNB, based on a random direction.
             Run for 40 seconds and collect stats every 200 msec.

---

Pseudocode 3: Pseudocode to generate **DB3**

---

initialization
$numGnbs \leftarrow 21$; $numUesPerGnb \leftarrow 10$
$numMcsValues \leftarrow 3$; $numRuns \leftarrow 8$
$numNeighboursPerUe \leftarrow 8$;

1  **for** $r \leftarrow 1$ **to** $numRuns$ **do**
2     **for** $m \leftarrow 1$ **to** $numMcsValues$ **do**
3        **for** $e \leftarrow 1$ **to** $numGnbs$ **do**
4           **for** $u \leftarrow 1$ **to** $numUesPerGnb$ **do**
5              **for** $n \leftarrow 1$ **to** $numNeighboursPerUe$ **do**
                Fix the Initial MCS of all the UE to $m$.
                Start moving away from the serving gNB, based on a random direction.
                HO to neighbour cell $n$ as per event A2.
                Run for 40 seconds and collect statistics every 200 msec.

---

**Table 4.3:** Database parameters.

| Parameter | Value |
|---|---|
| Total number of input features | 84 |
| Maximum neighbours | 8 |
| UE measurement periodicity | 200 ms |
| MCS values considered | 0 (QPSK), 14 (16 QAM), 28 (64 QAM) |
| Total simulation runs | HO use case : 20<br>Initial MCS use case : 63<br>Multi-Task use case : 8 |

#### 4.3.1.3   Parallel MTL database (DB3)

The procedure to generate the third DB (**DB3**) is the combination of the HO use case (see Pseudocode to generate **DB1**), which is then extended to repeat for all the potential initial MCSs for each gNB. This DB is generated using Pseudocode 3.

#### 4.3.1.4   Incremental MTL database (DB4)

For the evaluation of MTL, we also consider a fourth alternative DB (**DB4**). This DB is built incrementally based on the previous availability of **DB1** and **DB2**. In particular, it starts from **DB1**, and it incrementally adds data from **DB2**. With this DB we aim to evaluate the capability of the proposed architecture to incrementally learn a new task, once it has already been trained for other tasks. It would allow scalability in the RAN management, since new RAN tasks could be incrementally added to the architecture without additional implementation costs. All these databases could be intuitively expressed in the form of a dataset, as detailed in the next subsection.

### 4.3.2   Resulting database

A DB generated using any of the aforementioned pseudocode can be expressed as a matrix $\overline{\mathbf{X}}$.

$$\overline{\mathbf{X}} = \begin{bmatrix} \overline{x}_{1,1} & \overline{x}_{1,2} & \cdots & \overline{x}_{1,m} \\ \overline{x}_{2,1} & \ddots & \cdots & \overline{x}_{2,m} \\ \vdots & \vdots & \overline{x}_{i,j} & \vdots \\ \overline{x}_{n,1} & \overline{x}_{n,2} & \cdots & \overline{x}_{n,m} \end{bmatrix} \tag{4.1}$$

where the feature vector of size 84 (i.e., the total number of input feature for our proposed model) is $\overline{x}_{i,j} \in \overline{\mathbf{X}}$, $1 \leq i \leq n$, and $1 \leq j \leq m$.

The parameter $m$ defines the duration of the time series to be analyzed (i.e., the number of samples in the total simulation time, sampled with UE measurement periodicity), which corresponds to the number of time steps that the LSTM processes to perform the prediction. This number of time steps is the same for all the databases since we used the same periodicity to collect the measurements. On the other hand, the upper limit of $n$,

i.e., the total number of entries in each database is different for all the 4 databases, and can be computed by multiplying the total number of UEs with the maximum neighbor BSs and/or the initial MCS to explore, and the total number of simulation runs. During the simulations, it may happen that some of the data are not available or are not valid. For example, in the simulations used to build **DB1** UEs might experience a RLF when forced to HO to a BS with poor channel conditions. When this happens, we do not have data since the user is not connected. On the other hand, for **DB2** it might happen that the initial throughput is not available due to the fact that the download is concluded before the measurement could be taken (see Section 4.5.1.1 for more details). In these cases, after removing the affected entries from the databases for the overall simulation scenario, the total number of entries, i.e., the parameter $n$, are: 33,500 for **DB1**, 29,648 for **DB2**, 33,662 for **DB3**, and 31856 for **DB4**. Moreover, the total number of simulation runs for each database are selected such that each database contains approximately equal number of entries. The parameters, which dimension these databases are listed in Table 4.3.

Each simulation, where by one simulation we mean the individual run considered for 1 HO to a deterministic target cell and 1 initial MCS, lasts approximately 4,5 hours on an Intel(R) Xeon(R) CPU E5-2650 v3 @ 2.30GHz platform. We have parallel processing capabilities with 40 cores. The raw ns-3 traces occupy 250 MB per simulation. ns-3 traces have been merged in a unique file, where 1 row corresponds to 1 UE data, and the columns are the features previously introduced. Each simulation, once post-processed, occupies 11 MB.

## 4.4 RNN models for Single-Task and Multi-Task learning

In this section, we discuss the RNN models to solve the proposed RAN use cases through individual and joined deep learning models, following a traditional single-task learning or an MTL approach. A RAN efficient management involves several RAN use cases, which are usually handled by ML independent control loops. It means that a separate model is optimized for each task, which results in several task-specific models. However, the single-task approach presents many limitations in terms of coordination of the different tasks, negatively interfering among them, and is challenging from the implementation and computational perspectives. Specifically, we need models to perform multiple tasks in parallel without significantly compromising each tasks' performance. When it comes to learning multiple tasks under a single model, MTL techniques have been proposed in the literature as the solutions.

As mentioned in Section 4.3, the dataset consists of the measurements and traces extracted with a certain periodicity from each layer of the 3GPP protocol stack, which generates a time series of multivariate features. We believe that, by exploiting the temporal characteristic of this data one could understand the impact of HO decisions or select an appropriate initial MCS. Therefore, we propose different architectures, employing RNN with LSTM units [127]. LSTM is a special kind of RNN, which outperforms other ML approaches for time series analysis [128] [129], and solves the

**Figure 4.2:** Single-Task: Many to one LSTM architecture.

problem of long-term dependency issue found in vanilla RNN [130]. Specifically, we propose to model both target use cases as regression problems using LSTM-based architectures where we aim to estimate, respectively, the time to download the file, and the initial throughput over a certain time window.

### 4.4.1 Single-Task Learning

Fig. 4.2, shows the proposed multi-layer many-to-one LSTM architecture for single-task solutions, which is individually designed and trained to address the two selected use cases for study [125]. This model takes all the 84 features as input to infer the time to download for the HO management and the throughput for the initial MCS selection use cases. It processes them in a lag of 16800 data (i.e., 84 features x 200 time steps) samples with multiple batches of fixed size. Moreover, during the course of this study we found this single-task architecture to be very effective in handling the HO management use case, in comparison to other options. Therefore, in this study, we leverage the same model for the initial MCS selection use case but after fine-tuning its hyperparameters, as discussed in the next Section 4.5.1.

### 4.4.2 Multi-task Learning

In MTL, multiple tasks, each of which can be a general learning task, i.e., supervised, unsupervised, semi-supervised, or reinforcement learning tasks, are simultaneously learned through a shared model. It is found that concurrently learning these tasks can lead to performance and/or computational improvement compared to learning them individually. MTL is inspired by human learning activities where people often apply the knowledge learned from previous tasks to help learn a new task. It helps alleviating well-known weaknesses of deep learning, like the large-scale data requirements and

computational demand. We believe that it also brings an added value to the design of an intelligent RAN, where multiple correlated tasks have to be executed concurrently. The setting of multi-task learning is similar to transfer learning. The main difference, though, lies in the fact that in MTL, there is no distinction among different tasks, and the objective is to improve the performance of all the tasks, or reduce the computational component of all the joined tasks together. On the other hand, in transfer learning, the objective is to improve a target task with the support of source tasks. Learning separately multiple tasks brings difficulties that are not present in multi-task learning. It may happen that different tasks have conflicting needs. This may easily happen during the optimization of the RAN, where different tasks may intervene over the same parameters to optimize their functions independently. When the increasing performance of a model of one task hurts the performance of another task with different needs, we talk about *negative transfer*. There are many different factors to consider when creating a shared architecture, such as the portion of the model's parameters that will be shared between tasks. Many of the proposed architectures for MTL play a balancing game with the degree of information sharing between tasks: Too much sharing will lead to negative transfer and can cause the worse performance of the multi-task than the single-task model. At the same time, too little sharing does not allow the model to leverage information between tasks effectively. One commonly used multi-task architecture in computer vision follows the general vision of a global feature extractor made of convolutional layers shared by all tasks, followed by an individual output branch for each task. This architectural approach is usually referred to as *shared trunk*. Other architectures can follow alternative methods, for example, based on having a separate network for each task, with information flows between parallel layers in the different task networks. In the rest of this section, we discuss the architecture and the different options for learning that we propose to implement the MTL vision for efficient RAN management. In particular, we propose an architecture that follows a shared trunk architecture with hard parameter sharing [131] to enhance parameters estimation; however, in our study, we are more interested in its computational efficiency, since it allows us to share the AE training phase among the different tasks [132].

The architecture is based on a multi-layer LSTM AE [133], in charge of performing the shared feature extraction in conjunction with a task specific MultiLayer Perceptron (MLP) neural network, as shown in Fig. 4.3. An AE is an unsupervised ML algorithm, which learns a function to approximate an output identical to the input. Since it is based on the *encoder-decoder* paradigm, the input is transformed into a lower-dimensional space, also known as codeword, to more efficiently model highly non-linear dependencies in the inputs. The compression operation manages to extract more general and useful features, which retain essential aspects of a dataset [134]. Our goal is to smartly reduce the data to be used for inferring the time to download and the initial throughput. We use a single AE whose codeword is shared among the tasks but independent MLPs to estimate the specific QoE indicator of interest for each RAN use case. We opt for this LSTM based architecture, for the same reason already discussed for the single-task case, which is to take the best advantage of the temporal characteristic of the collected data.

In this line, using the model shown in Fig. 4.3 we propose the following two different methods for MTL learning:

**Figure 4.3:** Multi-Task: A combined HO management and initial MCS architecture consist of AE + MLP neural network.

- *Parallel MTL*: We explore learning behaviour when no task is given priority with respect to the others, but all tasks are concurrently learned. The reference database for training, in this case, is **DB3**.

- *Incremental MTL*: In this case, we analyze the learning behaviour when the learning is inherently incremental, meaning we first learn for one task and then incorporate information from new tasks. The advantage of this approach is that once we have trained the shared trunk architecture, we can progressively introduce more tasks to the design of the intelligent RAN, without further implementation costs. The risk, on the other hand, is that of the *catastrophic forgetting* [135], while we aim to incorporate information from new tasks without forgetting the previously learned. The reference databases to train in this case is **DB4**.

All in all, in this section, we proposed different RNN models for single-task, and multi-task approaches that address two RAN use cases, i.e., HO management and initial MCS selection, as shown in Table 4.4.

**Table 4.4:** Proposed approaches and RNN models.

| RNN models | Approach | RAN use case | DB type |
|---|---|---|---|
| Multi-layer many-to-one-LSTM | Single-Task learning | HO | DB1 |
| | | Initial MCS | DB2 |
| Auto Encoder + MLP neural network | Parallel multi-task learning | HO | DB3 |
| | | Initial MCS | DB3 |
| | Incremental multi-task learning | HO | DB4 |
| | | Initial MCS | DB4 |

## 4.5 Performance evaluation

This section discusses training and system-level performance evaluation using the proposed models for the single, and the two multi-task approaches.

### 4.5.1 Training the proposed architectures

The implementation of the models is done in Python, using Keras and Tensorflow as backend. In particular, to speed up the training, testing, and evaluation of these models, we use fast LSTM implementation with Nvidia CUDA Deep Neural Network (CuDNN) library for GPUs [136]. The DBs for each of the proposed approaches have been randomly divided into training and validation sets, using a split ratio of 0.75 and 0.25, respectively. We train and validate the models using the training and validation sets to minimize the reconstruction error over 200 epochs, in case of the AE, or prediction error, in case of the single-task LSTM and MLP. The loss function used to train the models is the Mean Square Error (MSE), and the *RMSProp* algorithm is used to optimize the learning process. Moreover, a linear activation function is used for the output layer of the LSTM (see Fig. 4.2) and MLP, while the *Leaky ReLU* activation function is used for the hidden layers of MLP. We discuss in the following the details of the training and selected architecture for the single task and multi-task architectures.

#### 4.5.1.1 Single-task architectures

We have trained the LSTM architecture shown in Fig. 4.2 for the single task learning based on **DB1**. To select the hyperparameters of this model, i.e., the number of layers (i.e., the values of K and L) and the number of LSTM units (blue LSTM blocks in Fig. 4.2) in each hidden layer, we have tested nine different combinations. Then, we have selected the hyperparameters resulting in the lowest average MSE (over 200 epochs). Fig. 4.4.(a) shows the MSE per epoch of the selected single-task model trained to address the HO management task, using 3 (i.e., K = 2 and L = 1) layers of LSTM nodes, where the numbers, i.e., [84x62x42] separated by "x" in the legend represent the number of hidden LSTM units in each layer. We observe that, after 140 epochs, this model is able to achieve and maintain very low testing loss independently from the number of layers and cells per layer.

We follow a similar approach to train another model based on the single-task architecture in Fig. 4.2, to tackle the initial MCS use case using **DB2**. As shown in Table 4.3, we focus

(a)



(b)

**Figure 4.4:** Single-task training MSE per epoch. (a) HO management use case LSTM 84×62×42 (b) Initial MCS use case LSTM 84×42×22.

on three MCSs (0, 14, and 28), representative of the three main available modulations (e.g., QPSK, 16QAM, 64QAM). The output to be estimated in the regression problem is the initial throughput computed over a window of 100 ms when a new RRC connection is established. The initial throughput at the beginning of the session cannot be considered because, during the initial window of a TCP connection, we are only able to capture messages from its initial handshake, which results in the same initial throughput for all UEs. Moreover, we also avoid taking the measurement after the first HO because the first HO is deterministic as for the HO management use case. The reason is that the **DB2** should be of similar nature as of **DB1** to be combined to construct **DB4**, which is later used for MTL. As a result, we focus on the initial throughput after the second HO.

**4.5.1.2 Multi-task architecture trained with parallel MTL**

For the multi-task architectures, we use a similar approach to the single-task one. The architecture now is based on a shared trunk approach, which first considers a multi-layer LSTM AE in charge of performing the shared feature extraction, and then a per task MLP is used to perform the regression. To train and validate these models, we used **DB3**, obtained to target both initial MCS and HO management use cases. It is also worth mentioning that we consider 8 runs to build **DB3**, in order to maintain a similar dimension for **DB1** and **DB2**. In this training process, first we select the code-word length of the AE, among five different code-word lengths of 50, 100, 200, and 300. In particular, we select a code-word equal to 100 to take into the account the trade-off between the length of the code-word and the MSE of the decoder. Fig. 4.5 (a) shows the AE reconstruction error, i.e., the MSE between original data and the one after decoding, common to the two use cases, using code-word length of 100. Then, using this selected code-word as an input to the MLP neural network, one set of hyperparameters, among 7 (based of the lowest average MSE), is chosen for the two MLPs (see Fig. 4.3). Similarly, Fig. 4.5 (b)-(c) show the regression loss of the AE plus MLP structure for the chosen MLP structure of three layers [80x40x20], to estimate the time to download and the initial throughput, for the HO use case, and for the initial MCS use case, respectively.

**4.5.1.3 Multi-task architecture trained with incremental MTL**

Fig. 4.6 (a), shows the loss of the AE, considering incrementally increasing databases for training (**DB4**). We start from the AE trained with **DB1**, which is indicated in the figure with "0 runs". Then, we progressively add runs from **DB2**, and observe the behavior of the loss of the AE. We consider that the loss is comparable in all cases. Furthermore, Fig. 4.6.(b)-(c) show the loss of the regressors to estimate respectively the time to download the file for the HO use case, and the initial throughput for the initial MCS use case, as a function of the different number of runs. Different behaviours can be observed in the loss of the two regressors. It depends on whether the architecture is first trained for one task, or another task is incrementally learned after the first one, by adding training data to the training DB. In the HO use case, for which the architecture is initially individually trained, we observe that the loss increases when 8 runs are introduced from **DB2**. It is because the architecture has to suddenly adapt to new data coming from a DB built for a different purpose. However, as we add more runs from **DB2**, the loss trend is to get reduced.

On the other hand, the loss in the estimation of the initial throughput for the initial MCS use case, which is the new use case we aim to learn by adding the new data, linearly decreases with the number of runs that we add from **DB2**. This behaviour of the loss is reasonable, since the architecture gradually improves its learning performance, as we add more information related to the MCS use case. We select the combination of AE code-word length of 100 and the MLP of [64x32], which provides us with the lowest average MSE for all the tested run values.

(a)



(b)



(c)

**Figure 4.5:** Multi-task parallel MTL training MSE per epoch. (a) MSE between original data and decoder for the AE trained with DB3 for codeword length 100 (b) HO management use case AE + MLP 80×40×20 (c) Initial MCS use case AE + MLP 80×40×20.

## 4.5.2 System level performance

The performance evaluation of these models is performed in an offline manner using the Offline-Evaluation procedure described in Section. 1.3.2. In particular, to perform this evaluation, we consider two extra simulation campaigns using a Run value which was not used to build the training dataset (i.e., Run 21 for the HO use case and Run 65 for

**Figure 4.6:** Multi-task incremental MTL training average MSE over 200 epochs. (a) Average MSE between original data and decoder as a function of the incremental runs for the AE of codeword length 100 (b) HO management use case AE + MLP 64×32 (c) Initial MCS use case AE + MLP 64×32.

the initial MCS use case). This approach allows us to evaluate the generalization of the models. In the following, we analyse the results obtained using the trained models for the two use cases.

**Figure 4.7:** ECDF of the difference of download time obtained by the benchmark A2-based HO benchmark and the ML-based architectures.

### 4.5.2.1  HO use case

For the HO management, we compare the time to download for each UE, obtained after selecting the target cell providing the lowest predicted time to download, to the one achieved by using a benchmark approach, i.e., A2-RSRP-based HO algorithm. The first campaign aims at gathering the file download time using the benchmark HO algorithm (e.g., A2-RSRP). The second simulation campaign is conducted in a similar way as the one to build the training dataset, i.e., it consists of 8 deterministic HOs. Following this approach, we construct 8 input strings for each neighbour of a UE, which consists of 1 row and 16800 columns (i.e, 84 features x 200 time steps). These strings are used individually as their input to obtain a predicted time to download for all the architectures, i.e., single task and multi-task (parallel and incremental learning). Finally, for each UE, we select the gNB with the minimum predicted time to download for the HO. We compare results of the ML-based and the benchmark approaches for the UEs that successfully finalize the download. In particular, we compare the number of UEs completing the download and the time needed to download the file.

We first compare the performance of the benchmark HO algorithm with the single-task and the parallel MTL architectures. Fig. 4.7 shows the Empirical Cumulative Distribution Function (ECDF) of the difference between the download time observed by these UEs using the benchmark and the proposed models. The results obtained using the benchmark HO algorithm show that there are 63 (i.e., 30%) UEs out of 210, which are able to finalize the download. On the other hand, 77 (i.e., $\approx 37\%$) UEs are able to successfully download the file using the single-task and parallel MTL models. It means that the ML approach manages to increase by 18% the number of UEs able to finalize the download during the simulation time. Moreover, there are 62 common UEs, which were always able to download the file, irrespective of the tested HO solution, i.e., benchmark or ML-based.

**Figure 4.8:** Example of reduction in the duration of radio link failure with proposed ML models.

Out of these 62 UEs, the ECDF trend in Fig. 4.7 on the positive x-axis shows that we can reduce the file download time for 56 UEs compared to the benchmark case using the single-task or the parallel MTL architectures. However, there are 6 UEs that experience marginally higher download time than the benchmark (see the trend on –ve x-axis). We believe that their performance can be improved by increasing the size of the database used to train the models and by further fine tuning their hyper-parameters. Moreover, this evaluation shows that the MLP, fed with the AE code-word of 100 performs similarly to the LSTM. This proves that the AE has efficiently transformed the inputs into a lower-dimensional space without losing the meaningful information of the dataset for the use case of the HO.

We now evaluate the capabilities of incremental MTL offered by the AE-based architecture. In particular, we want to prove that an AE that is trained for a specific use case (e.g., the HO) can be reused for another use case. As mentioned in Section 4.5.1.3, we first consider the AE and MLP model trained with **DB1**. In this DB we removed the entries where the initial throughput is not available for the reasons described earlier, e.g., when the file download finishes before the second HO. We observe that the HO performance based on this architecture is similar to the one obtained with the single-task learning. It is reasonable since with "0 runs" the DB is still purely built to handle the HO use case only. However, using the other four incrementally trained models, we notice that the performance of the HO algorithm is the same for all of them. The reason is that, even while observing some difference, the regression losses of these models are comparable and low enough to provide comparable system performance. This result is further validated when compared to the results achieved using single-task and parallel MTL, as shown in Fig. 4.7. In this figure, to simplify the representation, we only present the results using the incremental MTL model trained with "8 runs" from **DB2**.

73

**Table 4.5:** Summary of HO use case results.

| Approach | % of UEs finalizing the download | % of UEs decreasing the time to download |
|---|---|---|
| A2-based benchmark | 30% (out of 210) | - |
| ML-based | 37% (out of 210) | 90% (out of 62) |

The offline evaluation performance for the HO use case reaches exactly the same results for all the approaches. It allows us to conclude that incrementally introducing runs from a different database adding new information to the system, does not jeopardize the previously learned information, in our case, where the features of the two databases are the same. For our case and the nature of the database, we do not observe any phenomenon of catastrophic forgetting, i.e., the tendency of an artificial neural network to entirely and abruptly forget previously learned information upon learning new information. More research should be conducted to evaluate how different **DB1** and **DB2** can be to maintain the same conclusion that we reach here.

Furthermore, in Fig. 4.8, we present, as an example, 2 UEs out of 56 UEs for which ML reduced the time spent in RLF and, consequently, the time to download the file (there are more UEs in the scenario experiencing the same performance advantage when using the ML technique). We notice that these UEs experience an RLF just before the first HO irrespective of the scheme used, i.e., benchmark or ML. In fact, once the UE is inside a coverage hole generated by an obstacle, all gNBs are unable to offer any service, and there is no coverage from any of the surrounding gNBs. As a result, an obstacle impairs the coverage of all gNBs equally. However, in those challenging situations, with the help of ML models, we can reduce the RLF duration for these UEs by 400 ms (i.e., 400 Transmission Time Intervals (TTIs)), which also improves their time to download. The reason is that ML models, thanks to their capability of learning from past experience, can identify a more appropriate neighbour gNB to HO to provide more extended service than the benchmark, and doing so reduces the RLF duration. Finally, Table 4.5 summarizes the comparative results between the different approaches for the HO use case.

#### 4.5.2.2 Initial MCS use case

We evaluate the initial MCS performance, following an offline strategy, as we did previously for the HO use case, using an extra run, "Run 65". First we consider traces for the three evaluated MCSs, which provide three different input strings for each UE to get the predicted initial throughput for the selected MCS values. Then, for each UE, we choose the MCS, which results in higher initial throughput. At this point, we further filter out some UEs for those cases when the proposed ML models select MCS 0. In particular, the throughput achieved by the UEs using the benchmark scheme, which always selects MCS 0, is the same that we get when the ML-based solutions also consider the same MCS. Therefore, we consider only those UEs for the performance evaluation for which the ML-based models select MCS values different from 0, i.e., 14 and 28.

We observe similar trends for the initial MCS use case to those previously observed for the HO management case. In particular, the gain in the performance is the same when using any of the proposed ML-based models. In total, we obtain 86 UEs out of

**Figure 4.9:** ECDF of the difference between initial throughput obtained with benchmark MCS selection and the ML-based architectures.

210 for which we were able to record the initial throughput after the second HO, and the proposed approaches select a different MCS from 0. Out of these 86 UEs, 44 (i.e., 51.16%) select MCS 28 and 42 (i.e., 48.83%) select MCS 14. From the analysis of the results in Fig. 4.9, we conclude that all the 86 UEs that obtained a initial throughput get better initial throughput than the benchmark, when considering the ML-based approaches. The average initial throughput per UE considering the benchmark with MCS 0 is 0,051 Mbit/s, while the average initial throughput attained using the ML-based models is 0,1944 Mbps. Thus, on average, we obtain a 73.75% increment of initial throughput per UE. Moreover, for the incremental MTL approach we observe no difference in the performance of the selection of the appropriate initial MCS, when using 8, 16, 32, or 48 runs from **DB2**, to train the AE and MLP. In this case, we also believe that the average MSE using the traces only from 8 runs is already low enough (see Fig. 4.6.(c)) to provide the performance similar to the one using the single-task or parallel MTL approaches. Therefore, in Fig. 4.9 we plot only the results obtained using the incrementally trained model using 8 runs.

It is also worth mentioning that in case of incremental MTL, we are able to obtain already acceptable results for both use cases by using the joined DB, which has a similar dimension as of the individual **DB1** or **DB2**. On the other hand, to target the use cases with single-task approaches, we should train two independent architectures with a database of a dimension twice as big as the one we need with the incremental MTL use case. The advantage that we get with the incremental MTL, with respect to the parallel MTL or the separated single-task approaches, is at the implementation level since, at any moment, we are able to add a new use case to our learning architecture by incrementally training the model. It guarantees scalability concerning all the RAN use cases we wish to add to the design. Finally, Table 4.6 summarizes the gains obtained using ML approaches over the benchmark scheme for the initial MCS use case.

**Table 4.6:** Summary of initial MCS use case results.

| Approach | % of UEs selecting MCS 28 | % of UEs selecting MCS 14 | Initial throughput increase per UE |
|---|---|---|---|
| ML-based | 51.16% | 48.83% | 73.75% |

## 4.6 Conclusions

In this chapter, we tackled the complexities along infrastructure-based single-technology axes using two RAN use cases, 1) the HO management and 2) the selection of an initial MCS when UEs establish a new connection with a eNB. In particular, to exploit the network data's temporal characteristics, we proposed an RNN based on LSTM for single-task learning to address the aforementioned two use cases separately. This was done to answer our two research questions, (RQ1): "How to use AI/ML in mobility management to achieve better QoE?" and (RQ2): "How to use AI/ML to select the initial MCS for newly connected mobile devices to achieve better throughput?". For the single-task approach used for HO management, the results proved that the proposed AI/ML models outperform the A2 event-based benchmark HO algorithm in terms of the number of successful downloads and time to download statistics. Additionally, a similar model proposed for initial MCS selection also provided the gain in terms of the increased initial throughput by selecting a better MCS than a benchmark scheme, which always selects MCS 0 upon establishing a new RRC connection.

After the single-task learning solutions, we proposed an AI/ML framework to answer our third research question (RQ3) "How to generalize an AI/ML solution to address diverse RAN use cases?" To do so, we proposed an AI/ML model based on LSTM AE along with an MLP for the MTL approach. The results show that the models based on AE, used for the MTL parallel and incremental learning, perform similarly to the single-task model using only the LSTM. It is proven that the AE could efficiently compress the inputs into a lower-dimensional space without losing the dataset's meaningful information. The MTL solution, which allows sharing training models among RAN tasks, provides then a series of advantages at implementation, coordination, and training levels. Additionally, the model trained by employing the incremental learning approach did not suffer from the phenomenon of *catastrophic forgetting.*

Technology
complexity axis

MT

IMT

ST

IST          LIST

I          LI          Infrastructure support
complexity axis

# Chapter 5

# Contention Window Size Adaptation in LTE-LAA

As mentioned in Chapter 2, Section 2.2.1, 3GPP has standardized LBT as the default channel access scheme for LAA. However, in spite of adopting LBT, the performance of Wi-Fi when coexisting with LAA is highly dependent on how the LBT parameters are configured by LAA. In this chapter, we focus on the adaptation of the CW size parameter of LBT in LAA. This parameter is of key importance to avoiding collisions or to resolve the contention among the colliding stations. Specifically, we propose a CW size adaptation scheme, which could infer the collisions in all the subframes of a TxOP by combining the HARQ feedbacks from LAA UE and the sensing data gathered at the eNB. In particular, the proposed scheme learns from its past experience, through a machine learning approach, how many Negative Acknowledgements (NACKs) per subframe of a TxOP will be received under certain channel conditions.

Similar to Chapter 2, we based our model on FFNN due to its capability of providing a compact and easy to evaluate model as compared to other machine learning approaches [66]. According to our scheme, the eNB gathers the HARQ NACKs reported by the UEs resulting from past TxOPs and the sensing data by using a Wi-Fi listener at the eNB. This provides the information about the radio activity of other Wi-Fi Access Points (APs) and STAs. We note that, already different LTE-U/LAA products include a Wi-Fi listener for similar purposes in their implementation [137] [138].

The proposed scheme is able to predict the number of NACKs for all the subframes of a TxOP, without waiting for any delayed HARQ feedback after the TxOP ends. To summarize, the proposed scheme overcomes the limitations faced by the 3GPP or other

options found in the literature, by predicting the collisions in all the subframes of a TxOP in a timely manner. Furthermore, unlike the schemes discussed in the literature, the CW size is exponentially increased upon the reception of a NACK for each subframe of a TxOP and is not dependent on the information exchange between the LAA nodes or only on the sensing performed at the eNB.

This chapter is organized as follows. The related work around the CW update of LAA is presented in Section 5.1. Section 5.2, presents in detail the limitations of the CW size adaptation scheme proposed by 3GPP. We discuss the proposed scheme and the design of our FFNN in Section 5.3. In Section 5.4, we present a brief explanation of the simulation scenario and the benchmark schemes to which we compare our scheme. Finally, Section 5.5 presents the achieved results and Section 5.6 summarizes the main conclusion.

## 5.1   Related work

According to 3GPP, the CW size is proposed to be increased if 80% of the HARQ feedbacks belonging to the first subframe of the most recent TxOP are NACKs [80]. This scheme has two potential drawbacks. First, since LTE is capable of scheduling multiple users in a single subframe, the 80% threshold may be hard to meet. If a collision happens, but less than 80% of the scheduled users suffer from the collision, the LAA eNB will not increase its CW, and the collision will remain undetected. Second, due to the inherent latencies introduced by the LTE protocol stack, the HARQ feedback associated with a certain subframe is received at least 4 ms after its transmission time (see Fig. 2.3). Therefore, 3GPP proposes only to consider the collisions detected during the first subframe of a TxOP to update the CW with a minimum delay. As a result, the collisions from the rest of the subframes are ignored. Besides what is proposed by the standard, the literature also proposes some other techniques to adapt the CW size of LAA. For example, in [139] an analytical model based on a Markov chain is proposed to find an optimal fixed CW size for the LAA eNBs in the scenario. However, a fixed CW size may increase the chances for an LAA eNB to access the medium at similar times, and the CW size is not updated upon collisions. Similarly, authors in [140] built a model based on the Markov chain to investigate the LAA and Wi-Fi coexistence performance analytically. Furthermore, the authors also propose an updated LBT process for LAA by introducing a Maximum Contention Window Timer Mechanism (MCWTM). The MCWTM mimics the Wi-Fi Retry Limit Mechanism [141] by regulating the count with which the maximum contention window can be consecutively used. Although it is analytically shown that MCWTM reduces the system waiting time under poor channel conditions, the model assumes an ideal HARQ feedback scheme. An enhanced LBT algorithm is proposed in [142] for adapting the CW size in LAA according to the information exchanged among the neighboring nodes. However, the performance of this scheme is dependent on the information exchange among the LAA nodes, and this requires extra signalling to be defined and transmitted. In [143], a sensing-based scheme for LAA eNB is proposed to adjust the CW size by comparing the ratio of busy slots between two backoff periods. This scheme overcomes the limitation of the previously discussed schemes. Still, the CW size is updated only based on sensing performed at the eNB, without considering the

feedback from the user. Therefore, this scheme will easily be vulnerable to hidden node problems.

## 5.2   Limitation of 3GPP approach

As we have highlighted in the previous section, the LAA channel access procedure is very similar to the Wi-Fi's one. However, it still presents some dissimilarities due to the inherent differences in LTE and Wi-Fi technologies. We discuss here the main limitations that we see in the LAA channel access procedure, which may generate issues at the time of fulfilling a fair coexistence with Wi-Fi.

- **80% threshold may be hard to meet even in case of collisions**. Various UEs that are scheduled in the same subframe may experience in general different levels of interference due to the different channel conditions. In order to receive 80% of NACKs, to update the CW in LAA, it is necessary that more than 80% of UEs, which are scheduled in the same subframe suffer from the high level of interference by neighboring Wi-Fi or LAA nodes. However, due to the random nature of the radio propagation environment and different geolocations of the users, it might not always be the case. This limits the application of CW update only to the scenarios with high interference, while in other scenarios this condition might never be met. In such a case, collisions at many UEs may be ignored and the LAA backoff mechanism may not be properly exploited.

- **Length of the TXOP in LAA may happen to be higher than Wi-Fi**. The Wi-Fi STAs in all the coexistence studies of LAA and Wi-Fi comply with the new Wi-Fi standards, such as, IEEE 802.11n or 802.11ac. In this study the Wi-Fi network is composed of IEEE 802.11n complying STAs. According to the MAC layer enhancement proposed under this standard, multiple MAC layer frames are aggregated to form one big Physical Service Data Unit (PSDU). The maximum length of a PSDU or an Aggregated MAC Layer Protocol Data Unit (A-MPDU) is 65535 bytes [144]. Now, let us consider a scenario where both LAA and Wi-Fi have BE traffic in the downlink, since the main use case for LAA was proposed to be the offload of BE and BK traffic [79]. In Wi-Fi (802.11n), the default TXOP for BE traffic Access Category (AC) is equal to 0. The TXOP value of 0 means that only one Physical Protocol Data Unit (PPDU) can be transmitted at a time before competing again for access to the channel, with the maximum CW size of 1023 [145]. The maximum allowed time in which a PPDU containing an A-MPDU can be transmitted is up to 5.484 ms in the Mixed mode, and up to 10 ms if the Greenfield mode is used [144]. For a more detailed description of these modes, the reader is referred to the vast available literature on IEEE 802.11n enhancements (for instance, [145] [146] [147]). The Greenfield mode is not a widely adopted feature and is especially avoided in large-scale networks [146]. Therefore, all the Wi-Fi nodes in this study are configured to use the Mixed mode. Even if the Greenfield mode is used, the transmission time for a PPDU can be short when a Wi-Fi transmitter uses a higher Modulation and Coding Scheme (MCS) and its A-MPDU size is limited to 65535 bytes. Contrarily, according to the LAA LBT priority class 3 [80], LAA is

**Figure 5.1:** Contention window update comparison

allowed to occupy the channel up to 8 ms (or even 10 ms), regardless of the MCS and channel conditions. This could lead to an unfair behavior to Wi-Fi, since LAA after winning the channel may occupy it for a longer time than Wi-Fi. Moreover, in case of collisions, Wi-Fi will spend more time in backoff as compared to the LAA LBT priority class 3. This is due to the difference in Wi-Fi and LAA maximum allowed CW size, which is 63.

- **Collisions in LAA are not always detected**. The HARQ procedure in LTE uses one of the soft combining techniques, i.e., Incremental Redundancy (IR) and Chase Combining (CC), in which the failed transmissions are not wasted, but combined with the retransmissions. Thus, it may happen that an unsuccessful retransmission, due to a collision, does not result in a NACK, because the combined information is enough for the UE to decode the data successfully. On the other hand, Wi-Fi uses the Automatic Repeat Request (ARQ) with an ACK. Unlike HARQ, ARQ always discards the data with errors and asks for a new transmission. Therefore, due to the efficiency of a soft combining technique used in the LTE HARQ procedure, it may happen that for the same collision, Wi-Fi detects more collisions than LAA. In this case, it would be beneficial for LAA to also consider the feedbacks from other subframes of a TXOP to detect collisions.

As a result of the above observations, even if the LAA and Wi-Fi channel access mechanisms are similar, the CW of LAA will not evolve in the same way as Wi-Fi CW, and most of the CW updates will be concentrated around the lower CW values. This behavior is demonstrated in Fig.5.1, where the results are obtained from a simulation of LAA and Wi-Fi coexisting nodes in an indoor scenario described in Sec.5.4.1, and we observe the evolution of the contention window upon collisions. Therefore, the LAA eNB will not only backoff less, but may also take the channel for a longer time compared to Wi-Fi, if it has a longer TXOP. This may cause degradation in Wi-Fi performance when coexisting with LAA.

# 5.3 Proposed scheme

## 5.3.1 CW size adaptation algorithm for LBT-LAA

As discussed in Sec.5.2, if the state-of-the-art (SOTA) LAA CW size adaptation scheme is unable to meet the 80% collision requirement in the first subframe, and also ignores the collisions after the first subframe, this could result in an unfair behavior in terms of channel occupancy towards a coexisting Wi-Fi network.

To tackle this unfair behavior of LAA, we propose to consider the HARQ feedback from all the subframes of a TxOP. However, due to the inherent protocol latency in LTE, all the HARQ feedbacks belonging to the subframes of a TxOP will not be available on time, i.e., when the grant time out occurs, as it is shown in Fig. 2.3. To overcome this limitation, we propose a supervised Neural Network (NN)-based CW size adaptation scheme. Details on the specifications of the NN used, and the reasons for making these design choices are given in the following subsection. This learning scheme infers the possible number of NACKs which could be received during a TxOP. We note that, we increase the NACK counter only once for each subframe due to the possibility of receiving multiple NACKs per subframe. This means, that the maximum number of NACKs the NN can predict is less than or equal to the number of subframes in a TxOP. The eNB builds a profile of each TxOP by storing the number of NACKs received for all the subframes and the additional sensing data belonging to a TxOP, as defined in the Eq.5.2 of the following subsection. We note that, since we use a Wi-Fi listener at the eNB, the sensing data related to the Wi-Fi transmissions can be stored even when the eNB is transmitting. The NN is trained by using this sensing data as input, and the total number of NACKs received for all the subframes as output. Once the training is completed, the eNB, after the grant timeout, uses this NN to predict the expected number of NACKs to be received, without waiting for any pending HARQ feedbacks. This predicted number of NACKs corresponds to the number of CW updates to be performed at the end of a TxOP. For example, if the predicted number of NACKs is 2, the CW size of LAA is increased twice. Then, if the eNB has more data to transmit, it initially senses the channel for $T_d$ period and chooses a random backoff value between 0 and the updated *CW* value. Once the backoff counter reaches zero and the eNB is allowed to transmit, i.e., before the start of the next TxOP, we reset the CW size to $CW_{min}$. With the help of the NN, our scheme is able to predict all the collisions happening in each subframe of the TxOP, so that the CW of LAA can expand at a faster pace than the other SOTA approaches. We consider that this approach is reasonable, given the fact that in general LAA, due to its longer TxOP, will occupy more channel as compared to Wi-Fi.

However, it would be unfair for LAA to keep the same CW size also for the following TxOP, for which it has already performed the backoff at the end of the previous TxOP. The proposed scheme is further illustrated with the help of Algorithm.1.

---

**Algorithm 1** LAA NN-based CW size adaptation scheme

---

1: $PN$ $-$ Predicted NACKS by NN
2: $BC$ $-$ Backoff counter
3: $CI$ $-$ Channel idle time
4: $GT$ $-$ Grant timeout
5:                                              $\triangleright$ Initialization
6: $CW_p \leftarrow CW_{min}$
7: **if** $PN > 0$ && $GT == 1$ **then**                $\triangleright$ End of current TxOP
8:     **for** $i \leftarrow 1$ **to** $PN$ **do**
9:         **if** $CW_p == CW_{max}$ **then**
10:            **break**
11:         **else**
12:            $CW_p \leftarrow (2 * CW_p)$
13:         **end if**
14:     **end for**
15:     $BC \leftarrow \mathbf{rand}(0, CW_p)$
16:     **function** DEFER()                   $\triangleright$ Sense the channel for $T_d$
17:     **while** $BC \neq 0$ **do**                       $\triangleright$ Start Backoff
18:         **if** $CI == 9$ $\mu$s **then**         $\triangleright$ Channel idle for slot duration
19:            $BC \leftarrow BC - 1$
20:         **else**
21:            **go to** 16
22:         **end if**
23:     **end while**
24:     **if** $BC == 0$ **then**                     $\triangleright$ Before next TxOP
25:         $CW_p \leftarrow CW_{min}$
26:     **end if**
27: **else**
28:     $CW_p \leftarrow CW_{min}$                       $\triangleright$ No NACKS
29: **end if**

---

### 5.3.2 Implementation of the proposed scheme

To implement our aforementioned scheme, once again, we choose a two-layer FFNN because of its ability to model both linear and non-linear functions between inputs and outputs [67]. Moreover, the model obtained with FFNN is more compact and faster to evaluate than other machine learning techniques such as, Support Vector Machine (SVM), with the same generalization performance [66]. And, since we are employing supervised learning using offline-Training procedure explained in Section. 1.3.2 to build the FFNN model, we need to build a training and testing dataset [67]. This dataset can be represented as,

$$\mathbf{D} = \begin{bmatrix} (X_1, Y_1) \\ \vdots \\ \vdots \\ (X_N, Y_N) \end{bmatrix} \qquad (5.1)$$

where $X$ and $Y$ are the input and output vectors of FFNN and $N$ is the total number of rows in a dataset. Each row in the dataset corresponds to a single TxOP, such that, $N$ is equal to the total number of TxOPs, observed during the simulation campaign for building the database. Let the term $k$ denote the TxOP ID of each TxOP in our dataset. The input vector $X$ can be written as,

$$
\begin{aligned}
X_{k=1} &= [w_1, u_1, c_1, d_1] \\
X_{k=2} &= [w_2, u_2, c_2, d_2, p_1] \\
X_{k=3} &= [w_3, u_3, c_3, d_3, p_2, p_1] \\
X_{k=N} &= [w_N, u_N, c_N, d_N, .., p_{N-1}, ..p_{N-L+1}]
\end{aligned}
\tag{5.2}
$$

where $L$ is the memory of the NN and,

- $w_k$ = Wi-Fi transmissions observed during $k$-*th* TxOP

- $u_k$ = Number of UE scheduled during $k$-*th* TxOP

- $c_k$ = NACKs received in $k$-*th* TxOP, before its grant timeouts

- $d_k$ = Duration of the $k$-*th* TxOP

- $p_{k-1}$ = NACKs received in $k-1$ TxOP

Finally, the output vector $Y$ contains the total number of NACKs belonging to all the subframes of a TxOP. For the implementation of FFNN, we use the publicly available *nnet* package of R [67], which is a single hidden layer FFNN. To obtain the above mentioned data for each LAA eNB, we run a simulation campaign with a SOTA CW size scheme for 2000 sec. The details of the simulation scenario are presented in the next section. During this simulation campaign, on average, for every LAA base station, 9000 data samples are stored in the form of the dataset $D$. Then, this dataset is randomly divided into a training set (containing 75% of the data) and a testing set (containing 25% of the data).

Referring to Eq. 5.2, we consider $N$ to be equal to the total number of TxOPs observed during the simulation campaign for each eNB and $L = 3$ in our case. The value of parameter $L$ is selected by experimenting with a range of values, i.e., $1 \leq L \geq 20$, where $L = 3$ resulted in prediction accuracy as obtained using higher values. Then, all the input and output values are normalized in the range [0,1]. This normalization allows for a faster training process and more accurate estimations [113]. After performing the 3-fold cross-validation test, the final FFNN model for each eNB is selected based on the lowest NRMSE. On average, the selected FFNN model for each eNB converges after 830 iterations and gives us more than 85% prediction accuracy. The cross-validation test is performed to prevent overfitting, i.e., when FFNN achieves the ideal minimization of the error between the estimated and the actual output of the training set. If the overfitting happens, the FFNN loses its generalization property and fails to predict the output of the testing dataset.

**Figure 5.2:** BS Corner indoor scenario layout

# 5.4   Performance evaluation setup

## 5.4.1   Simulation Scenario

The simulation scenario implementation has been done using the LAA and Wi-Fi coexistence module of *ns-3* simulator [148]. As shown in Fig. 5.2, we consider an indoor simulation scenario with the base stations placed on the corners of a building with the dimension $120 \times 50$ meters with no walls. This scenario is based on the 3GPP indoor scenario used for evaluation studies presented in [79], but the base stations are placed at the corners of a building to get closer to a real world implementation. Moreover, this deployment helps us better evaluate the performance of our scheme by increasing the number of collisions generated by hidden nodes. We evaluate the fairness according to the methodology used by 3GPP, in which there are two operators, i.e., operator A and operator B. In step 1, both operators deploy Wi-Fi, while in step 2, operator A substitutes Wi-Fi with LAA. There are four base stations with the fixed location and 20 UEs/STAs per operator, which are randomly dropped inside the building. We consider that both LAA nodes and Wi-Fi nodes use the same channel 36 of 20 MHz. Moreover, The Wi-Fi nodes in our simulations are configured with MIMO $2 \times 2$ TX/RX antennas, supporting the rates up to MCS 15 with a long guard interval. Similarly, for LAA nodes we use the *ns-3* MIMO model supporting up to MCS 28.

As for the traffic model, we consider the File Transfer Protocol (FTP) 1 model proposed by 3GPP in [79]. We choose LBT priority class 3 [80] with TxOP of 8 ms and $\lambda=5$ to simulate a level of load that allows both LAA and Wi-Fi to always have data available to fill their TxOPs. We note that, for the purpose of fair evaluation, the maximum CW size of both LAA and Wi-Fi is set to 1023, instead of 63 and 1023. Moreover, for the propagation model, we use 802.11ax indoor model for all the small cells in the scenario. Simulations are run for 2000 sec for all scenarios.

## 5.4.2   Benchmark Schemes

In this subsection, we provide a brief explanation of our benchmark CW size adaptation schemes, implemented to evaluate the performance of our proposed scheme.

### 5.4.2.1 3GPP Sensing scheme

This scheme is implemented as per the proposal submitted in 3GPP by different vendors [149]. According to this proposal, the adaptation of LAA CW size is based on the observation of busy and idle slots at the eNB in an observation window. The observation window is the time between the random backoff counter is drawn and the time when the counter reaches zero. The CW size of LAA is increased if the following condition is met,

$$\frac{Number\ of\ busy\ slots}{Number\ of\ CCA\ slots} > Initial\ backoff\ counter \tag{5.3}$$

Otherwise, the CW size is reset to $CW_{min}$.

### 5.4.2.2 Preamble detection scheme

This scheme leverages the Wi-Fi listener at the eNB to detect the collisions between the eNB transmissions and the captured Wi-Fi signals. Specifically, the Wi-Fi listener calculates the Signal-to-Interference-plus-Noise Ratio (SINR) based on the preamble of the captured Wi-Fi signal. This SINR value is internally passed to the eNB. The eNB marks a collision and increases its CW, if the SINR falls below 0 dB and the current eNB state is in transmission. The CW size is reset to $CW_{min}$, once the eNB is allowed to take the channel for the next TxOP, i.e., after performing the channel access procedure.

### 5.4.2.3 Ideal HARQ scheme

As indicated by its name, it is an ideal CW size adaptation scheme in which the HARQ feedbacks from UEs are available to the eNB without any delay, i.e., at the end of each subframe of a TxOP. The eNB increases the CW once per subframe, if any of the HARQ feedback belonging to a subframe is a NACK. Similar to the preamble detection scheme, once the eNB has performed the random backoff on the basis of updated CW size, the CW size is reset before the next TxOP.

## 5.5 Results

We use user perceived "throughput" and "latency" as our main performance metrics to evaluate the performance of all the schemes presented in this chapter. In *ns-3*, these metrics can be calculated by using the built-in FlowMonitor tool that tracks per-flow statistics at the IP layer including throughput and latency. We note that, the Wi-Fi performance of operator B in step 1 is the baseline performance for a Wi-Fi network, when it coexists with LAA, under all the CW size adaptation schemes. Let us first discuss the performance of SOTA-HARQ and 3GPP sensing-based schemes. Fig. 5.3 (b) and Fig. 5.4 (b), show that Wi-Fi experiences the lowest throughput and latency performance when LAA uses these schemes. This performance degradation in Wi-Fi is

**Figure 5.3:** Throughput (Mbps) performance comparison. (a) LAA, (b) Wi-Fi.

mainly caused by two reasons. First, as mentioned in Sec.5.2, we observe that Wi-Fi experiences shorter TxOP than LAA, and consequently on average it spends more time in contention as compared to what it spends in step 1, i.e., when it coexists with other Wi-Fi networks. As a result, Wi-Fi flows experience higher latencies as compared to the baseline latencies. Second, in case of collisions caused by hidden nodes, the CW size of Wi-Fi is increased to its maximum value more often than the LAA, as shown in Fig. 5.1. This increases the backoff time for the Wi-Fi transmission for the next channel access, which results again in a higher latency and a lower throughput for Wi-Fi.

On the contrary, with the SOTA HARQ scheme the CW size of LAA does not reach its higher values because of not meeting the 80% threshold in the first subframe and also due to the fact that collisions in other subframes of the TxOP, except the first subframe, are ignored. Similarly, for the 3GPP sensing-based scheme, the threshold for increasing the CW size in Eq. 5.3 is not often met, because it depends on the random activity of other nodes during the backoff period of LAA. Moreover, it does not consider the real collisions happening during the TxOP of the LAA node. Therefore, LAA flows with any

**Figure 5.4:** Latency (msec) performance comparison. (a) LAA, (b) Wi-Fi.

of these two schemes experience lower latencies and higher throughput due to its lower backoff time and longer TxOP, as shown in Fig. 5.4 (a) and Fig. 5.3 (a).

With the preamble detection-based scheme, LAA is enabled to guarantee a fair coexistence to Wi-Fi in terms of both throughput and latency. However, this scheme results not being fair to LAA, since it increases the CW size of LAA without evaluating if the overlapping Wi-Fi signal is actually causing a corruption of the UE data.

As a result, LAA backoffs more than necessary, which eventually limits its channel occupancy and result in its low throughput and higher latency.

We now discuss the performance of our NN-based CW size adaptation scheme, in comparison to the proposed benchmark schemes (i.e., SOTA-HARQ, 3GPP sensing and Preamble detection) and an ideal HARQ scheme. We note that, here, the performance evaluation of the proposed NN-based model is performed using the Semi-Online-Evaluation approach explained in Section. 1.3.2. The ideal HARQ scheme serves the purpose of an optimal scheme which detects all the collisions in a TxOP of LAA

on the basis of timely received HARQ feedbacks. In this way, we can see how close the performance of our NN-based scheme is to the optimal one, and how much improvement has been achieved by our scheme over SOTA and other benchmark schemes. As shown in Fig. 5.3 (b) and Fig. 5.4 (b), the NN-based scheme enables LAA to better coexist with Wi-Fi by improving its coexistence performance as compared to SOTA-HARQ and 3GPP sensing-based schemes. The reason is that, the NN-based scheme on the basis of its learning capability can predict the number of NACKs for all the subframes of a TxOP when the grant timeout occurs. Thus, it overcomes the limitations of SOTA-HARQ and 3GPP sensing-based scheme discussed earlier. On the other hand, the preamble-based scheme provides better coexistence performance to Wi-Fi as compared to the NN-based scheme. However, this performance gain is achieved at the cost of extreme degradation in LAA performance, since as already explained, this scheme is unfair to LAA. The NN-based scheme achieves a better trade-off between the degradation in LAA performance and the improvement in Wi-Fi ones, which ensures a better coexistence. Moreover, its performance trends are very similar to those shown by the ideal HARQ scheme.

## 5.6  Conclusions

This chapter focused on the complexity bounded by the infrastructure-based multi-technologies axes. Specifically, we have studied the LAA CW size adaptation approach proposed by 3GPP. Based on the simulation results obtained through a 3GPP aligned LAA module developed in *ns-3*, we show that due to the limitations of this approach, the evolution of the CW in LAA is very different from the Wi-Fi CW. Following this approach, it is complicated to combine HARQ feedbacks from multiple UEs, which may be experiencing different channel conditions, and then extract a decision over the occurrence of a collision based on those combined HARQ feedbacks. Furthermore, due to the inherent latencies in the LTE protocol stack, this approach only considers the HARQ feedbacks from one subframe of a TxOP to reduce the delay between the transmission and the detection of the collision. To overcome these limitations, we propose a solution based on supervised machine learning using a FFNN, which answers our fourth research question (RQ4): "How to use AI/ML to guarantee fair coexistence of LAA in the unlicensed spectrum?" Using our scheme, the LAA eNB is able to predict the collisions for all the subframes of a TxOP, and increases the CW size at the end of a TxOP without waiting for delayed HARQ feedbacks. Our performance evaluation shows that the FFNN-based scheme provides a better coexistence performance to Wi-Fi as compared to the 3GPP approach. Moreover, when compared to the schemes which provide better coexistence to Wi-Fi by degrading the LAA performance, the FFNN-based scheme provides a better trade-off by achieving a similar Wi-Fi performance with minimum degradation in LAA performance.

The fairness evaluation conducted in this chapter strictly follows the 3GPP methodology [79]. That is, we drew the above conclusions with the help of the throughput and latency CDF curves in Fig 5.3 and Fig 5.4. However, interpreting those curves by just looking at them is not so straightforward when they overlap in a particular region of the plot and for others not, e.g., see the curves of "WiFi (Preamble Detection)", "WiFi (NN)", and "WiFi (Baseline)" in Fig 5.3 (b). Therefore, in the next chapter, we present

a statistical framework to understand better the fairness evaluation of LAA and LTE-U towards Wi-Fi.

Technology
complexity axis

IMT

MT

IST          LIST

ST

I                    LI          Infrastructure support
complexity axis

# Chapter 6

# On Fairness Evaluation: LTE-U vs. LAA

In this chapter, we propose a statistical framework to systematically evaluate the fairness offered by different LTE technologies when they coexist with Wi-Fi in the unlicensed band. In particular, we study the coexistence performance of both 3GPP LAA and LTE-U, as specified by LTE-U forum. As mentioned in Chapter 2, the 5G NR technology, the NR-U, is also based on the 3GPP LAA methodology, hence, follows the same fairness evaluation procedure as in LTE [45]. Therefore, the proposed framework, though applied to LTE based unlicensed technologies in this thesis, can be easily used to evaluate the fairness of 5G and beyond technologies, e.g., NR-U. In the rest of the chapter, when we refer to ULTE, we will refer to both LAA and LTE-U unless explicitly specified.

We map the generally accepted 3GPP definition of fairness onto the stochastic dominance concept. Specifically, we use the two-sample one-sided Kolmogorov-Smirnov test (KS-test) to test the specific hypothesis of fairness defined through throughput and latency performance, as proposed by 3GPP. We evaluate throughput and latency using the *ns-3* simulator for LTE and Wi-Fi coexistence. Particularly, a desired hypothesis regarding throughput fairness is that when coexisting with ULTE, Wi-Fi throughput distribution curves should stochastically dominate those throughput curves when coexisting with another Wi-Fi. That is, a Wi-Fi that coexists with ULTE must achieve an equivalent or higher throughput when it coexists with another Wi-Fi. For latency, the hypothesis is exactly the contrary, i.e., the CDF of the latency when Wi-Fi coexists with ULTE, should not stochastically dominate the one, when Wi-Fi coexists with another Wi-Fi, as low latencies are desirable for a fair coexistence. We test these hypotheses through

KS-test because it is a useful tool to test the first-order stochastic dominance, especially with our data, since it does not require any hypothesis on the distribution of the data.

The rest of the chapter is organized as follows. In section 6.1, we present the related work. Next, the statistical framework we propose to use for the fairness evaluation is presented in section 6.2. Section 6.3 discusses the simulation results and the data analysis campaign. Finally, section 6.4 concludes the work with final thoughts on the fairness evaluation and suggestions for future works.

## 6.1 Related work

The use of LTE in unlicensed spectrum generates multiple challenges, since LTE has been designed to work in licensed spectrum on the basis of uninterrupted and synchronous operation. Consequently, the efficient use of unlicensed bands to offload LTE traffic is of main concern for both ULTE and widely deployed Wi-Fi networks. According to 3GPP, for a fair co-existence with Wi-Fi, LAA must not harm a Wi-Fi network more than an additional Wi-Fi network on the same band. Based on this definition, in the literature different contributions try to evaluate the coexistence performance of ULTE and Wi-Fi. In [150], the impact on the fairness of two LAA-LBT-based channel access schemes with the Wi-Fi network has been studied. The performance analysis is done by comparing the mean and the CDFs of latency and throughput under different signal/energy detection thresholds by LAA-LBT procedure. In [151], the authors have proposed two channel sensing schemes for LTE in unlicensed band, and the fairness of both schemes is established on the basis of mean user throughput of the Wi-Fi network. In [152], the performance of various co-existence methods based on LBT and Duty Cycle (DC) have been evaluated using the Monte-Carlo simulations. The fairness of these methods is evaluated on the basis of average number of collisions and average latency experienced by the coexisting Wi-Fi network. We observe that the conclusion on the fairness in these works are mainly driven by qualitative comparisons of average values and CDFs.

Differently, in this work, we focus on statistically evaluating the fairness offered by both LTE-U and LAA paradigms to a coexisting Wi-Fi network. We perform a simulation study of the coexistence behavior by following the 3GPP methodology presented in [79], when Wi-Fi coexists with both technologies. We discuss in detail the behavior of the throughput and latency curves and we try to draw conclusion in terms of fairness. We observe, however, that the coexistence curves may converge in some areas and diverge in others, and it is not so intuitive to claim if ULTE is actually fair to Wi-Fi, or not. In particular, this qualitative approach leaves much space to partial interpretations. Therefore, in this study, our main focus is towards the use of more rigorous approach which allows us to statistically evaluate the concept of fairness and not just only the comparison of the two ULTE paradigms, which has already been discussed in many works in the literature. In this context, statistical analysis offers interesting tools and concepts which can serve the purpose. In the following section, we employ these concepts to propose a statistical framework that maps the 3GPP fairness definition onto the concept of stochastic dominance.

## 6.2    Statistical framework

The fairness as defined by 3GPP [79], *is the capability of an LAA network not to impact Wi-Fi networks active on a carrier more than an additional Wi-Fi network operating on the same carrier, in terms of both throughput and latency.* The same fairness concept also has been proposed to evaluate fairness for LTE-U [137]. We evaluate fairness according to the scenarios defined by 3GPP, in which there are two operators, i.e., operator A and operator B. In the first step, both operators deploy Wi-Fi, while in the second step, operator A substitutes Wi-Fi with ULTE. We claim ULTE is fair to Wi-Fi if, when switching from Wi-Fi to ULTE the performance of operator B is not negatively affected.

However, based on this qualitative definition provided by 3GPP, we choose to compare fairness by comparing the CDFs distributions of throughput and latency. If the curves overlap or the Wi-Fi performance after substitution is improved, we can safely say that ULTE is coexisting fairly with Wi-Fi. Contrarily, if these curves partially overlap or diverge in some areas (this is what we observe in our results), one cannot do much beyond explaining the reasons for the divergence and it is still debatable if the ULTE behavior should be considered as fair or not. Therefore, we believe that there is a need of a more detailed analysis of the data, which can systematically tell us whether, up to a certain tolerable extent, ULTE behavior is fair or not. Statistical data analysis offers tools to compare the statistical behavior of data, and we believe that it could add value to the evaluation of fairness in ULTE and Wi-Fi coexistence. In particular, we rely on the concept of first order stochastic dominance, which assumes that *a distribution X stochastically dominates a distribution Y, if the CDF of X lies on the right side of CDF of Y.* One way for 3GPP definition of fairness to be quantified is to leverage the concept of stochastic dominance. We obtain empirical CDFs (ECDFs) of the key performance parameters and we measure the extent to which one CDF dominates the other CDF, which could be expressed as follows [153]:

$$T_{wl}(x) \leq T_{ww}(x) \qquad \forall x \in [0, \infty) \tag{6.1}$$

Where $T_{wl}(x)$ is the CDF of the throughput of a Wi-Fi network when it coexists with an ULTE network and $T_{ww}(x)$ is the CDF of the same Wi-Fi network, when it coexists with another Wi-Fi network. In practice this means that, for ULTE to be fair, the same or higher throughput must be obtained by operator B, after the substitution. So, when operator A substitutes Wi-Fi with ULTE, the throughput CDF of operator B must stochastically dominate the one, obtained before the substitution of ULTE. In terms of latency, the fairness could be achieved if the CDF of the latency distribution of a Wi-Fi network, when coexisting with ULTE network, does not stochastically dominate the baseline Wi-Fi distribution, since improved latency means smaller latency values. In other words,

$$L_{wl}(x) \geq L_{ww}(x) \qquad \forall x \in [0, \infty) \tag{6.2}$$

Where $L_{wl}(x)$ is the CDF of the latency in the ULTE-Wi-Fi coexistence scenario and $L_{ww}(x)$ is the CDF of the latency in the baseline scenario.

Statistical hypothesis testing is a procedure in which sampled data are employed to test a hypothesis about a single population or the correlation between two or more populations. This hypothesis is either a null hypothesis $H_0$, i.e., a statement about the distribution of observations we want to test, or an alternative hypothesis $H_1$, i.e., an alternative statement in the case a null hypothesis is failed [154]. In the following, we will establish these two types of hypotheses for throughput and latency. There are many statistical techniques which could be applied to test the hypothesis under study and they are mainly divided into parametric and non-parametric statistical tests. Following the method presented in [155], we use a non-parametric two-sample one-sided Kolmogorov–Smirnov test (KS-test) to test the first order stochastic dominance of two Wi-Fi distributions. This kind of tool is particularly useful in our problem, because it does not require any hypothesis on the underlying distribution of the data (e.g. many methods require normal distribution), and it offers no restrictions on the size of samples [156]. Following the steps of hypothesis testing, we state our $H_0$ and $H_1$ hypothesis for throughput and latency as follows,

1. Throughput:

   **$H_0$:** The throughput distribution of Wi-Fi when coexisting with ULTE stochastically dominates the baseline throughput distribution.

   **$H_1$:** The throughput distribution of Wi-Fi when coexisting with ULTE does not stochastically dominate the baseline throughput distribution.

2. Latency:

   **$H_0$:** The latency distribution of Wi-Fi when coexisting with ULTE does not stochastically dominate the baseline latency distribution.

   **$H_1$:** The latency distribution of Wi-Fi when coexisting with ULTE stochastically dominates the baseline latency distribution.

As a result of this test, two values are obtained:

- $D_{max}$: It is the maximum measured distance between the two ECDFs.

In the literature, $D_{max}$ is mentioned as the *test statistic* value for the KS test. Generally, to test the hypothesis in the test, the calculated $D_{max}$ value is compared with a critical D value ($D_{crit}$) obtained from the KS table, at a certain significance level $\alpha$[1]. The null hypothesis is rejected, if the value of $D_{max}$ is greater than $D_{crit}$ value [154]. Following a common statistical practice in the literature, we use $\alpha = 0.05$. In our case, by using the Table G in [154], for m=60 and n=62, the value of $D_{crit}$ is 0.22. Where m and n are the number of IP level flows in our simulation.

- P-value: It is the probability of a null hypothesis being true (with a certain significance level $\alpha$).

---

[1]The significance level $\alpha$ is the error probability for the Type I error in which $H_0$ is rejected when in fact $H_0$ is true. The term confidence interval and $\alpha$ are related such that: confidence interval = 1 - $\alpha$ [154].

**Figure 6.1:** Indoor scenario layout

The interesting point about the P-value is that it quantifies the dominance of the whole distribution and not simply how close/distant they are at a certain reference point or how similar the obtained average values are. If the P-value obtained in the test is less than $\alpha$, it indicates that we have less evidence for our null hypothesis to be true and vice versa. Specifically, if the P-value is below the significance level $\alpha$ we reject the null hypothesis and claim that the results obtained show unfairness to Wi-Fi. Therefore, the lower the P-value the more likely the technology evaluated is unfair to Wi-Fi. On the other hand, for values higher than the significance level $\alpha$, the more likely it is that the technology evaluated is fair to Wi-Fi. To summarize, in this study we propose the following steps to systematically evaluate the fairness:

1. Build throughput and latency ECDFs of the different flows crossing the Wi-Fi network of operator B, when it coexists with Wi-Fi network of operator A, and use it as baseline.

2. Repeat the procedure in step 1, when Wi-Fi network of operator B coexists with ULTE network of operator A.

3. Apply one-sided two-sample KS-test.

4. Use the P-value to accept or reject the null hypothesis.

## 6.3 Results

In this section, after discussing the simulation scenario, we first present LAA, LTE-U and Wi-Fi performance by analyzing the CDF plots of their throughput and latency, when they coexist with another Wi-Fi network in the indoor scenario as defined by 3GPP in [79]. In *ns-3*, we are calculating throughput and latency by using the built-in FlowMonitor tool that tracks per-flow statistics at the IP layer. We then post-process these results to obtain the CDFs. Later, on the basis of the statistical approach discussed in Section 6.2, we perform a statistical analysis on the throughput and latency distributions to evaluate the fairness of LAA and LTE-U networks towards Wi-Fi.

### 6.3.1 Simulation Scenario

As shown in Fig. 6.1, we consider the indoor simulation scenario proposed for coexistence evaluations by 3GPP [79]. In the indoor scenario, two operators deploy four small cells

**Figure 6.2:** Throughput performance: FTP($\lambda$=2.5) over UDP

in a building with the dimension 120x50 meters with no walls. The four base stations for each operator are equally spaced, while the base stations from the two operators are placed with an offset on the $X$ axis. There are 20 UEs/STAs per operator which are randomly dropped inside the building. The unlicensed 20 MHz band at 5.180 GHz is shared by both operators. The licensed band for ULTE networks is not simulated here.

As for the traffic model, we consider the FTP 1 model proposed by 3GPP in [79]. We consider the maximum allowed $\lambda$=2.5. Moreover, for the propagation model, we use 802.11ax indoor model for all the small cells in the scenario. Simulations are run for 52 sec for all scenarios, which is enough to generate ideally over 6000 TxOP of 8ms for ULTE operator to access the level of fairness it could offer to Wi-Fi.

### 6.3.2   LAA vs LTE-U performance analysis

Fig. 6.2 plots the CDF of the throughput achieved in the following scenarios,

1. Wi-Fi over Wi-Fi: Wi-Fi network of operator A coexists with Wi-Fi network of operator B. We present the throughput of Wi-Fi network of operator B, the trend is very similar to the one of operator A, which is not shown to simplify the figure.

2. LTE-U over Wi-Fi: LTE-U network of operator A coexists with Wi-Fi network of operator B. In this scenario, operator A network is substituted by LTE-U and we present the throughput of both LTE-U and Wi-Fi networks.

3. LAA over Wi-Fi: LAA network of operator A coexists with Wi-Fi network of operator B. In contrast to previous scenario, operator A network is substituted by LAA and the throughput of both LAA network of operator A and Wi-Fi network of operator B are presented.

Let us start by analyzing the LAA coexistence behavior. When Wi-Fi coexists with LAA, we observe that in most of the cases the LAA network is coexisting fairly with the Wi-Fi

**Figure 6.3:** Latency performance: FTP($\lambda$=2.5) over UDP

network, allowing it to achieve the same throughput as it has achieved in Wi-Fi over Wi-Fi scenario. However, in some cases Wi-Fi experiences medium and low throughput flows mainly due to:

1) Transmission of reservation and control signals, i.e., periodic transmissions of discovery reference signals (DRS) which occupy one subframe (1 msec), in case they have to be transmitted without data. These control signals occupy more channel than Wi-Fi beacons and interrupt Wi-Fi flows, causing more contention, and increase in latency as shown in Fig. 6.3. This prevents these flows from reaching the maximum throughput.

2) Collisions due to hidden terminals. We observe that 1.27% of signals experience collisions caused by nodes below the Wi-Fi CCA-ED threshold (-62dBm) and LAA ED threshold (-72dBm). These low throughput flows would benefit from further lowering the LAA and Wi-Fi ED thresholds to -82dBm or from LAA supporting CTS2Self [144]. This would enable the Wi-Fi nodes to backoff upon detecting the CTS2Self messages, which happens at -82dBm or below.

As per LAA throughput, the great majority of the flows achieve the maximum throughput, while three very low throughput flows are observed due to the hidden terminal problem, as mentioned above. In terms of latency, the same three flows suffer from high latency, as shown in Fig. 6.3. As discussed above for Wi-Fi, these low throughput flows can also be recovered by implementing CTS2Self functionality in LAA.

When co-existing with LTE-U, Wi-Fi throughput is mainly affected by the fact that LTE-U starts transmitting without first listening to the medium, as explained in chapter 1. This increases the chances to collide with ongoing Wi-Fi transmissions. In particular, we observe 1.27% of collisions in case of LAA and 2.83% of collisions in LTE-U, which represents a 55% increment. On the other hand, LTE-U performance as compared to LAA, is mainly degraded due to three reasons:

1) Increased number of possible collisions with the ongoing Wi-Fi transmissions due to the duty cycle transition from OFF period (T-OFF) to ON period (T-ON), as mentioned earlier.

2) The lack of LBT capability makes LTE-U nodes not backoff to each other and so they may happen to coexist with frequency reuse 1. This increases the spectral efficiency of

**Table 6.1:** KS-test results

| | | Indoor | |
|---|---|---|---|
| | | P-value | D-Max |
| LAA | Throughput | 0.00273 | 0.342 |
| | Latency | 0.00573 | 0.303 |
| LTE-U | Throughput | P≪0.05 | 0.543 |
| | Latency | P≪0.05 | 0.479 |

LTE-U network at the the cost of increased inter-cell interference. As a result, we observe lower MCS values and an increased number of Hybrid-ARQ (HARQ) retransmissions.

3) High latency experienced by the flows which are interrupted by LTE-U OFF periods. In our simulation, with the maximum achievable throughput, one file transfer of 0.5 MB takes ∼30 msec to complete. Consequently, the flows which are unable to complete during one T-ON period, i.e., 40 msec in our case with 50% duty cycle, have to wait for additional T-OFF time of 40 msec to resume. Therefore, these flows experience high latencies as compared to LAA, as it is shown in Fig. 6.3. Contrarily, those flows which in turn are able to complete their transmission during one T-ON period are able to achieve high throughput values comparable to those offered by LAA.

From this analysis of results, we observe that Wi-Fi curves when coexisting with LAA and LTE-U are similar, but they both show some divergence. We intuitively can say that LAA seems more fair, but we are unable to state up to which extent LTE is fair with Wi-Fi. As a result, in the following subsection we proceed to statistically analyzing and comparing the throughput and latency distribution to get more insights on the fairness of the coexistence performance.

### 6.3.3 Statistical analysis of fairness

As per our discussion in Section 6.2, we use a single-sided two sample KS-test to systematically estimate the fairness achieved in both LAA and LTE-U scenarios. In particular, we use the *Stats* package of R [157], which is a widely used open source tool in statistical studies. As shown in Table 6.1, the obtained P-values for throughput and latency are less than $\alpha$ (i.e., 0.05 in our case). This indicates that at significance level of 0.05, KS-test rejects the null hypothesis, i.e., the throughput distributions obtained from Wi-Fi over LAA do not stochastically dominate the throughput distribution of Wi-Fi over Wi-Fi, and the latency distributions obtained in the same scenario do stochastically dominate the throughput distribution of Wi-Fi over Wi-Fi. Similarly, for the KS-test of Wi-Fi over LTE-U, the resulting P-values for throughput and latency are much less than the significant level of $\alpha$, indicating that we have no evidence for $H_0$ to be true. The P-values which were very low and statistically insignificant are represented by P≪0.05 in Table 6.1. We reach the same conclusion when analyzing the $D_{max}$ in Table 6.1. We observe that in all cases $D_{max}$ is higher than the $D_{crit}$ value.

On the basis of these statistical results, we can conclude that for the indoor scenario defined by 3GPP that we evaluated, neither LAA, nor LTE-U pass the fairness test. However, we also conclude that LAA behaves more fairly than LTE-U.

## 6.4 Conclusions

In this chapter, we have presented the comparative study of the coexistence performance of the ULTE (i.e., LTE-U and LAA) network, when coexisting with the Wi-Fi network. We have started from the 3GPP definition of fairness and methodology for evaluating fairness. Specifically, we have followed the same methodology and evaluated the coexistence performance in terms of the throughput and latency in the indoor scenario that is implemented using the *ns-3* LTE-WiFi coexistence models of ULTE. These models allow for full protocol stack, and 3GPP and LTE-U forum-compliant evaluations. After analyzing the results, we conclude that it is not easy to claim in a quantitative way that the behavior of the ULTE network is actually fair/unfair towards the Wi-Fi network. Therefore, we have proposed a formal framework based on statistical data analysis to evaluate the concept of fairness. In particular, we have mapped the 3GPP fairness concept onto the first-order stochastic dominance concept. We have tested the hypothesis of such defined fairness through the KS-test.

Based on the results, ULTE behavior when 3GPP or LTE-U forum specs are followed, needs the support of proprietary solutions, e.g., the one presented in chapter 5, on top to allow a fair coexistence with WiFi technology. As a result, further work is required to improve the management of spectrum resources, and solutions strictly complying with the specifications are not enough to guarantee fairness. On the other hand, from the analysis conducted, it emerges that LAA, which provides access to the channel much similar to Wi-Fi, appears as a preferred technology for fair coexistence compared to LTE-U. This result should be further confirmed after considering different traffic models, and other scenario configurations, in the context of full protocol stack evaluations. Finally, the interesting future works in the area of fairness could be the ones that deal with generalizing the proposed framework by also considering second-order stochastic dominance concepts, which would allow concluding on the fairness even in cases of non-monotonic curves intersecting between each other.

Technology
complexity axis

MT

IMT

IST

LIST

ST

I

LI

Infrastructure support
complexity axis

# Chapter 7

# NR V2X Technology and Implementation

Following the successful use of sidelink in LTE for Proximity Services (ProSe) and C-V2X, the 3GPP has standardized its evolution in NR systems in the context of the so-called NR V2X. This new technology is expected to complement LTE C-V2X for advanced services by offering low latency, high reliability, and high throughput V2X services for advanced driving use cases. To do this, NR V2X is equipped with new features, such as the support for groupcast and unicast communication, a novel feedback channel, and a new control channel design. In this chapter, we provide detailed overview of NR V2X technology, with special emphasis on Mode 2 for out of coverage operation and autonomous resource selection. Furthermore, this chapter presents a system-level NR V2X standard-compliant simulator, as an extension of the popular and open-source NR network simulator 5G-LENA, based on *ns-3*. In particular, we focus on the design, implementation, and evaluation of the sensing-based resource selection in NR V2X Mode 2, in a highway scenario. Through several and extensive simulation campaigns, we test the impact of different NR V2X parameters, such as the numerology, the resource selection window size, the number of retransmissions, the maximum number of resources per reservation, and the probability of keeping the same resources during reselection, in a sensing-based resource selection. Finally, we provide a comparison campaign that shows the gains attained by the sensing-based resource selection, proposed during 3GPP Release 16, over the random selection strategy, considered in 3GPP Release 17 for power saving purposes.

The rest of the chapter is organized as follows. In Section 7.1, we provide a literature review of NR V2X and the existing open source simulators to simulate sidelink

communications. Section 7.2 presents 3GPP NR V2X specifications and reviews in detail NR V2X Mode 2 transmissions. Section 7.3 presents the simulation models and the implementation details. Section 7.4 discusses multiple simulation campaigns. Finally, Section 7.5 concludes the chapter.

# 7.1   Related work

While LTE C-V2X has been widely studied analytically and through simulations by academia and industry [89, 91, 158–160], the studies on NR V2X are comparatively new and growing in numbers day by day. Authors in [161] provide an overview of the standardization activities for vehicular communications at mmWave bands, including IEEE 802.11bd and 3GPP NR V2X specifications. Authors in [162] review the NR V2X design in 3GPP Release 16, with respect to the network architecture, security, and protocol enhancements. Authors in [163] provide a comprehensive overview of 3GPP NR sidelink transmissions, including physical layer structure, resource allocation mechanisms, and synchronization procedures. A more in-depth tutorial of 3GPP Release 16 NR V2X standard is presented in [164], including overview of the PHY layer, resource allocation, quality of service management, mobility management for V2N communications, and coexistence mechanisms between NR V2X and LTE C-V2X. In [165], the impact of the NR numerology on the V2X autonomous sidelink mode (similar to NR V2X Mode 2) is assessed. However, in [165], the evaluation is done over an LTE C-V2X simulator.

As it can be observed from the above review, most of the publicly available papers about NR V2X deal with a 3GPP standard overview, but few of these works discuss simulation studies. In addition, a key challenge to evaluate performance of NR V2X is that, despite the set of simulation results by industry and in literature, the simulators are not publicly available. Normally, simulators used in 3GPP are required to pass through a calibration procedure, but they are private, and not available to the research community. Consequently, the obtained results are neither reproducible, nor comparable, and system performance metrics are presented without much details revealed about the underlying models and assumptions. There are then private commercial simulators that are available after paying an annual license fee for using them. Often, if not in all cases, the license is very restrictive and does not allow modifications or inspection of the source code, which is a clear limitation for the research and the potential innovation. To the best of our knowledge, open source end-to-end simulators for 5G V2X communications compliant with NR V2X Release 16 specifications are not yet available to the research community.

There are five main open source and end-to-end simulators that have been developed to simulate sidelink communications. First, an LTE D2D communication simulation model based on *ns-3* was introduced and validated in [166]. Models are currently available through the *ns-3* App store. Authors in [167] presented the first open-source simulator for LTE C-V2X Mode 4 communications, based on *ns-3*. An open-source 802.11p and LTE C-V2X simulation/emulation tool for *ns-3*, called ms-van3t, has been recently released in [168], which provides integration of *ns-3* with the open-source Simulation of Urban MObility (SUMO) simulator for mobility management and mobility tracking. The work in [169] presents LTEV2Vsim, a simulator for LTE C-V2X Mode 3 and Mode 4 that is written in Matlab, freely available, and which focuses on Medium Access Control (MAC)

and PHY layers procedures. Finally, authors in [170] introduced an *ns-3* simulator for NR V2X at mmWave carrier frequencies. The model in [170] is compliant with 3GPP antenna and channel modeling for NR V2X, but not with NR V2X specifications at RRC and MAC layers. In particular, the model was developed before the finalization of NR V2X specifications. Therefore, at MAC layer it follows Mode 2 (c) for resource allocation, which was proposed as one of the options for study in 3GPP TR 38.885 [92]. Specifically, it uses TDMA to assign resources using the slots in a subframe, i.e., UEs scheduled in a subframe have orthogonal resources to transmit on; hence, they do not collide. Moreover, the error model used at the PHY layer is more suited for LTE but not for NR. All these limitations have been addressed by our model as described below and throughout this chapter.

## 7.2 NR V2X Technology Review

This section presents the main highlights of NR V2X technology in 3GPP, with special emphasis to NR V2X Mode 2 transmissions.

### 7.2.1 NR V2X

#### 7.2.1.1 Communication types

Differently from LTE C-V2X that focused on periodic basic safety messages, NR V2X has been designed to support various use cases, including the transmission of periodic traffic as well as reliable delivery of aperiodic messages. To support a wide range of different new applications, NR V2X goes beyond the only broadcast communications proposed by LTE C-V2X, and provides support for three types of transmissions: broadcast, groupcast, and unicast [85, 92]. In NR V2X unicast transmissions, the transmitting UE has a single receiver UE associated with it. The groupcast mode is used when the transmitting UE wishes to communicate with a specific sub-set of UEs in its vicinity. Finally, broadcast transmissions enable a UE to communicate with all UE within its transmission range. In NR V2X, a single UE can establish communications of multiple types simultaneously. For example, a platoon leader UE can communicate with its platoon member UEs using the groupcast mode, while using the broadcast mode to transmit other periodic messages to UE that are not part of the platoon.

#### 7.2.1.2 Sidelink physical channels and reference signals

Sidelink communications in NR V2X use the following physical channels [171]: 1) the Physical Sidelink Broadcast Channel (PSBCH) for sending broadcast information (like synchronization of the sidelink), 2) the PSCCH for sending control information (1st-stage-SCI), 3) the PSSCH for sending control (2nd-stage-SCI), data and CSI in case of unicast, 4) and the PSFCH for sending HARQ feedback in case of unicast and groupcast modes. The PSFCH is a new channel, which was not previously considered in LTE C-V2X. For these channels, numerologies 0 (SCS=15 kHz), 1 (SCS=30 kHz), and 2 (SCS=60 kHz) are

**Figure 7.1:** Time/frequency frame structure and definition of sidelink resource pool for NR V2X TDD. Example with 2 subchannels of 10 RBs each, using TDD pattern of [D D D F U U U U U U] and sidelink bitmap of [1 1 1 1 1 1 0 0 0 1 1 1].

supported at sub 6 GHz bands, and numerologies 2 (SCS=60 kHz) and 3 (SCS=120 kHz) can be used at mmWave bands [172]. For PSSCH, the supported modulation schemes include QPSK, 16-QAM, 64-QAM, and 256-QAM. Instead, for PSCCH, only QPSK transmission is supported.

Regarding the reference signals, NR V2X uses [163] 1) the Sidelink Primary/Secondary Synchronization Signal (S-PSS/S-SSS) for synchronization. S-PSS/S-SSS are transmitted together with the PSBCH in the so-called synchronization signal/PSBCH block (SSB). The SSB uses the same numerology as the PSCCH/PSSCH on that carrier. 2) Demodulation Reference Signals (DMRS) to estimate the channel and perform data decoding. 3) Phase Tracking Reference Signal (PT-RS) to compensate for phase noise. 4) Channel State Information Reference Signal (CSI-RS) to estimate the channel and report channel quality information, similarly to NR.

### 7.2.1.3 Sidelink resource pool

An important aspect of sidelink communications is the definition of sidelink resource pools. In NR V2X, a UE can be configured by higher layers with one or more sidelink resource pools. A sidelink resource pool can be used for transmission and reception of PSCCH/PSSCH, and can be associated with either sidelink resource allocation Mode 1 or Mode 2 [92]. In the frequency domain, a sidelink resource pool consists of a number of contiguous subchannels [173]. The size of each subchannel is fixed and it is composed of $N$ contiguous RBs. Both the number of subchannels and the subchannel size are higher layer pre-configured, by RRC. NR V2X supports $N = 10, 15, 20, 25, 50, 75$, and 100 RBs for possible sub-channel sizes [174]. In the time domain, the resources (i.e., slots) available for sidelink are determined by repeating sidelink bitmaps. The bitmap is pre-configured and characterized by a certain size. The resource pool parameter from RRC, sl-TimeResource, defines the bitmap size and takes values 10, 11, 12, ..., 160 [175].

In particular, in case of TDD, the resources available for sidelink are given by the combination of the TDD pattern and the sidelink bitmap. We also note that, unlike LTE sidelink specification related to the TDD pattern and the size of sidelink bitmap [176], NR sidelink specification is flexible and any valid NR TDD pattern can be used with any structure of a sidelink bitmap, which has a size specified by the standard [177]. Since NR V2X may be developed both in a carrier dedicated to Intelligent Transport System (ITS) or to cellular services, the standards support both the cases where all the symbols in a slot are available for sidelink, or only a consecutive subset of them [92]. In ITS spectrum, all the symbols are always allocated to sidelink. Within the slots available for sidelink, the specific Orthogonal Frequency Division Modulation (OFDM) symbols used for sidelink transmission/reception are fixed and pre-configured. Two RRC parameters pre-configure the symbol index of the first symbol and the set of consecutive symbols in a slot available for sidelink [173].

In Fig. 7.1 we illustrate the time/frequency frame structure of NR V2X and the definition of sidelink resource pools for TDD systems. The example is shown for the case of 10 MHz bandwidth using numerology 1 (i.e., SCS 30 kHz), and 2 subchannels, each composed of 10 RBs where RB 1 is the starting RB of the first sidelink subchannel. In time, we consider a TDD pattern of [D D D F U U U U U U] (i.e., one downlink slot, followed by a flexible slot[1], and three uplink slots), and a sidelink bitmap of [1 1 1 1 1 1 0 0 0 1 1 1]. As it can be observed, the TDD pattern is repeated in time, and each index of the sidelink bitmap applies to the uplink slots (U) in the TDD pattern, repeatedly, thus indicating the slots available for sidelink. In the frequency domain, a sidelink resource pool consists of a number of contiguous subchannels [173], therefore, as per [175], the last 4 RBs are not available for sidelink. As a result, in the figure we illustrate in green which slots/RBs are available for sidelink communications in the mentioned configuration example. This structure is typically used by an out-of-coverage NR V2X UE using Mode 2 operating in any of the V2X bands listed in Table 2.4. On the other hand, an in-coverage NR V2X UE operating in either Mode 1 or Mode 2, will tailor its time/frequency structure as per the gNB provided TDD pattern, sidelink bitmap, and subchannels.

#### 7.2.1.4 Retransmissions and new sidelink feedback channel

Differently from LTE C-V2X, which uses fixed MCS and only provides support for blind retransmissions, i.e., the source UE, automatically retransmits without knowing if the initial transmission has been correctly received, NR V2X provides different enhancements to improve reliability of communications, by introducing a feedback channel, the PSFCH. In particular, for unicast and groupcast communications considered by NR V2X, but not by LTE C-V2X, reliability can be improved if the source UE can retransmit the packet once the reception fails at the receiving UE and if the MCS can be adjusted to the actual channel conditions. NR-V2X introduces both blind and feedback-based retransmissions, for unicast and groupcast communications, while for broadcast communications only blind retransmissions are supported. With blind retransmissions, HARQ is implemented only at the receiver for retransmission combining. The transmitting UE chooses the resources within a resource reservation interval to retransmit. In particular the UE

---

[1]The flexible slot is used to provide the necessary guard time for downlink to uplink switching in TDD systems.

retransmits based on the configured value, which can be up to 31. Blind retransmissions are resource inefficient if the initial transmission is successful. On the other hand, feedback-based retransmissions are more resource efficient because the transmitting UE only retransmits if the original transmission is NACKed. In this case, HARQ is implemented both at the transmitter for efficient retransmissions and at the receiver for retransmission combining. In both cases, NR V2X Mode 2 supports a maximum number of PSSCH transmissions of the same MAC Packet Data Unit (PDU), which is pre-configured and whose maximum value is equal to 32. Even if feedback-based retransmissions are more resource efficient, blind retransmissions allow to minimize the latency of the feedback-based retransmissions, as the transmitting UE does not need to wait for a HARQ feedback before sending a retransmission [163]. To enable feedback-based retransmissions, NR V2X introduces the PSFCH.

### 7.2.1.5 Multiplexing of PSCCH, PSSCH, and PSFCH

In LTE C-V2X, PSCCH and PSSCH channels are multiplexed in the frequency domain. The drawback of this approach is that a receiver must buffer the message for the entire sub-frame and can decode the message only at the end of the sub-frame. This may be inefficient in NR V2X due to tight latency constraints of certain messages. To address this problem, different multiplexing options are considered in NR V2X for PSCCH and PSSCH [92]. Among the different options, two out of four consider time multiplexing in NR V2X, i.e., the PSCCH will be transmitted first, followed by the transmission of PSSCH. In time domain, within the symbols available for sidelink in a slot (see description in Section 7.2.1.3), the PSCCH can span over two or three symbols at the beginning of the pre-configured symbols and the PSSCH spans over the remaining number of pre-configured symbols. Finally, both PSCCH and PSSCH are multiplexed in time with the PSFCH. Specifically, every one, two, or four slots available for sidelink, the last two symbols among the pre-configured ones for sidelink, excluding the guard period symbol, are reserved for the PSFCH. In the frequency domain, the PSSCH can occupy up to the maximum number of available subchannels for sidelink, depending on the amount of data to transmit. However, the PSCCH spans over a pre-configured number of consecutive RBs (i.e., $K$ RBs) in the first subchannel in which PSSCH is transmitted, where $K \leq N$ RBs and $N$ is the subchannel size, described in Section 7.2.1.3. NR V2X supports $K = 10$, 12, 15, 20, and 25. Finally, the candidate resources, i.e., RBs for PSFCH are determined as per the PSSCH transmission(s) for which the feedback is generated. For more details, the reader is referred to very comprehensive tutorial of the NR V2X standard in [163] and [178].

The time multiplexing of the PSCCH, PSSCH, and PSFCH in NR V2X is shown in Fig. 7.2, assuming a typical NR slot structure composed of 14 OFDM symbols. As previously mentioned, the number of OFDM symbols used for sidelink is pre-configured. In the example in Fig. 7.2, 14 symbols are available for sidelink, the length of PSCCH is pre-configured to 2 symbols, PSSCH starts at the 3rd symbol with a duration of 8 symbols, $N = 20$ RBs is the subchannel size, and $K = 12$ RBs are used for the PSCCH.

**Figure 7.2:** Slot structure of a slot available for sidelink, with time multiplexing of PSCCH, PSSCH, and PSFCH in NR V2X.

### 7.2.1.6  Resource allocation

NR V2X defines two resource allocation modes for sidelink communications, one centralized and one distributed [92]:

- Mode 1 (A centralized scheduling approach): The NR base station (gNB) schedules sidelink resources to be used by in-coverage UEs for sidelink transmissions.

- Mode 2 (A distributed scheduling approach): The UE autonomously determines sidelink transmission resources within sidelink resources configured by the gNB or pre-configured by the network.

In this thesis, we focus on NR V2X Mode 2 with periodic traffic. Resource reservation for NR V2X Mode 2 under periodic traffic mostly reuses the LTE C-V2X sidelink Mode 4 long-term sensing-based algorithm. It exploits the periodicity and fixed-size assumption of basic safety messages. In addition to the long-term sensing-based resource selection, NR V2X Mode 2 also supports a random resource selection [179]. The difference between sensing-based and random resource selections is that, before selecting the resources from the total available ones, the sensing-based procedure filters those slots which are in use by other UEs, using sensing information. On the other hand, the random selection procedure does not use the sensing information, and directly selects the resources from the total available ones. The random resource selection approach is considered to reduce the complexity of the UE and the power consumption, since the sensing procedure adds complexity and has an energy cost at the transmitting UE that has to continuously sense the channel for resource selection. This random procedure was later approved by 3GPP in NR Release 17 WI [180], as a power saving mechanism. On the other hand, in case of aperiodic traffic, LTE C-V2X Mode 4 resource selection mechanism has been re-engineered for NR V2X Mode 2, since the arrival of future packets cannot be inferred

**Figure 7.3:** Illustration of the SCI split and resource reservation concept in NR V2X.

by sensing previous transmissions from surrounding UEs. For these cases, the use of short-term sensing and dynamic reservation is envisioned in 3GPP [92].

### 7.2.1.7 Sidelink control information

Another key improvement of NR V2X is the split of the SCI. Two SCI formats have been defined [181]: SCI Format 0-1 and SCI Format 0-2, which are sent through different channels, the PSCCH and the PSSCH, respectively. SCI carried on PSCCH is a 1st-stage SCI (SCI Format 0-1), which transports sidelink scheduling information of PSSCH and 2nd-stage-SCI on PSSCH. This sidelink scheduling information includes the priority, time/frequency resource assignment, 2nd-stage-SCI format, MCS, and resource reservation period. The SCI carried on PSSCH is a 2nd-stage-SCI (SCI Format 0-2), which transports information used for the decoding of PSSCH. This includes the HARQ process ID, new data indicator (NDI), redundancy version, source ID, and destination ID.

1st-stage-SCI indicates the reservation of $N_{\max\_reserve}$ (pre-configured) number of sidelink resources within the resource selection window [182]. $N_{\max\_reserve}$ can be 2 or 3 [175]. The resource reservation is indicated in the time resource assignment field of the 1st-stage-SCI. This means that not all the slots in a resource reservation period of a UE carry 1st-stage SCI in the PSCCH; some slots have empty PSCCH and only carry information in the PSSCH, as indicated by a 1st-stage-SCI in a previous slot.

An illustration of the SCI split and the resource reservation mechanism is shown in Fig. 7.3, for the case of $N_{\max\_reserve} = 3$ (i.e., each 1st-stage SCI can indicate up to three sidelink PSSCH resources) and $N_{selected} = 5$. $N_{selected}$ indicates the number of resources that are selected within a selection window, as per [182]. As an example, in NR V2X Mode 2, upon a resource selection trigger, the UE MAC can select various resources for an initial transmission (NDI = 1) and various retransmissions (NDI = 0). According to NR V2X specification, $N_{sci} = \min(N_{\max\_reserve}, N_{selected})$ is the number of resources indicated by a 1st-stage-SCI [173]. In Fig. 7.3, $N_{sci} = 3$ for the first 1st-stage-SCI and $N_{sci} = 2$ for the second 1st-stage SCI (since here, only two remaining slots are left to be indicated after the resources indicated by the first 1st-stage-SCI).

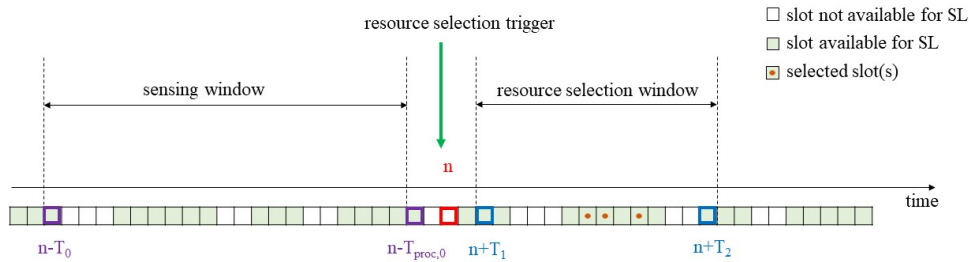**Figure 7.4:** NR V2X Mode 2 resource selection procedure. Example with $T_0 = 20$ slots, $T_{\text{proc},0} = 2$ slots, $T_1 = 2$ slots, and $T_2 = 16$ slots.

## 7.2.2 NR V2X Mode 2

NR V2X Mode 2 considers sensing-based SPS for periodic traffic. This is defined as a distributed scheduling protocol to autonomously select radio resources, in a similar way to what is already considered for LTE C-V2X Mode 4. The sensing procedure takes advantage of the periodic and predictable nature of V2X basic service messages. In particular, sensing-based SPS UEs reserve subchannels in the frequency domain for a random number of consecutive periodic transmissions in time domain. The number of slots for transmission and retransmissions within each periodic resource reservation period depends on the number of blind retransmissions (if any) and the resource selection procedure. The number of reserved subchannels per slot depends on the size of data to be transmitted.

### 7.2.2.1 Resource selection procedure

The sensing-based resource selection procedure is composed of two stages: 1) a sensing procedure and 2) a resource selection procedure [179].

The sensing procedure is in charge of identifying the resources which are candidate for resource selection and is based on the decoding of the 1st-stage-SCI received from the surrounding UEs and on sidelink power measurements in terms of RSRP [173]. The sensing procedure is performed during the so-called *sensing window*, defined by the pre-configured parameter $T_0$ and a UE-specific parameter $T_{\text{proc},0}$ that accounts for the time required to complete SCIs decoding and possibly perform measurements on DMRS for the sensing procedure. Specifically, if at time $n$ the sensing-based resource selection is triggered, the UE will consider the sidelink measurements performed during the interval $[n - T_0, n - T_{\text{proc},0})$. Sidelink RSRP measurements can be computed using the power spectral density of the signal received in the PSCCH or in the PSSCH, for which the UE has successfully decoded the 1st-stage-SCI. PSCCH RSRP and PSSCH RSRP are defined as the linear average over the power contributions (in Watts) of the resource elements that carry DMRS associated with PSCCH and PSSCH [183], respectively.

Based on the information extracted from the sensing, the resource selection procedure determines the resource(s) for sidelink transmissions [179]. For that, another window is defined, the *resource selection window*. The resource selection window is defined by the interval $[n + T_1, n + T_2]$, where $T_1$ and $T_2$ are two parameters that are determined by the UE implementation [173]. $T_2$ depends on the packet delay budget (PDB) and on an
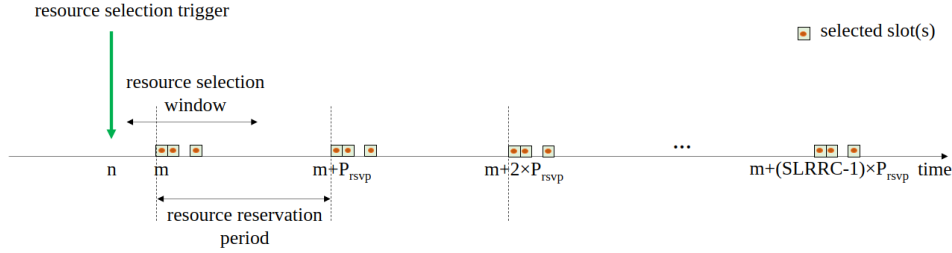
**Figure 7.5:** NR V2X Mode 2 semi-persistent scheduling.

RRC pre-configured parameter called $T_{2,\text{min}}$. In case PDB $> T_{2,\text{min}}$, $T_2$ is determined by the UE implementation and must meet the following condition: $T_{2,\text{min}} <= T_2 <=$ PDB. In case PDB $\leq T_{2,\text{min}}$, $T_2 =$ PDB. $T_1$ is selected so that $T_{\text{proc},1} <= T_1$, where $T_{\text{proc},1}$ is the time required to identify the candidate resources and select a subset of resources for sidelink transmission. The resource selection procedure is composed of two steps. First, the candidate resources within the resource selection window are identified. A resource is indicated as non-candidate if an SCI is received on that slot or the corresponding slot is reserved by a previous SCI, and the associated sidelink RSRP measurement is above a sidelink RSRP threshold [173]. The resulting set of candidate resources within the resource selection window should be at least a $X$ % of the total resources within the resource selection window to proceed with the second step of the resource selection. The value of $X$ is configured by RRC and can be 20 %, 35 % or 50 %. If this condition is not met, the RSRP threshold is increased by 3 dB and the procedure is repeated. Second, the transmitting UE performs the resource selection from the identified candidate resources (which may include initial transmissions and retransmissions). For that, a randomized resource selection from the identified candidate resources in the resource selection window is supported.

To exclude resources from the candidate pool based on sidelink measurements in previous slots, the resource reservation period (which is transmitted by the UEs in the 1st-stage-SCI) is introduced. As only the periodicity of transmissions can be extracted from the SCI, the UE that performs the resource selection uses this periodicity (if included in the decoded SCI) and assumes that the UE(s) that transmitted the SCI will do periodic transmissions with such a periodicity, during $Q$ periods. This allows to identify and exclude the non-candidate resources of the resource selection window. According to [173], $Q = \lceil \frac{T_{\text{scal}}}{P_{\text{rsvp}}} \rceil$, where $P_{\text{rsvp}}$ refers to the resource reservation period decoded from the SCI, and $T_{\text{scal}}$ corresponds to $T_2$ converted to units of ms [173].

As previously mentioned, NR V2X also supports a random resource selection [179]. In this case, the sensing procedure is omitted, and all the resources within the selection window that are part of the resource pool for sidelink are candidates for random selection.

Fig. 7.4 shows the resource selection procedure in NR V2X Mode 2. The figure illustrates the sensing window and resource selection window, with an example that uses $T_0 = 20$ slots, $T_{\text{proc},0} = 2$ slots, $T_1 = 2$ slots, and $T_2 = 16$ slots. Once the resource selection is triggered at time $n$, based on the measurements in the sensing window, the MAC scheduler determines the transmission resources within the resource selection window, which can be used for different MAC PDUs or to perform blind retransmissions.

**Figure 7.6:** NR Sidelink UE control plane.

### 7.2.2.2 Semi-persistent scheduling

Once one or multiple resources are selected, the UE will consider periodic transmissions, using SPS. The transmission interval is defined by the Resource Reservation Period ($P_{\text{rsvp}}$), which is pre-configured by RRC and can take predefined values between 1 ms and 1000 ms [179]. $P_{\text{rsvp}}$ value is included in the 1st-stage-SCI, to allow other UEs to estimate which resources are reserved in the future based on SCI decoding. After using the resource for the number of transmissions equal to the Sidelink Resource Reselection Counter (SLRRC), a resource reselection is triggered. Whether to reselect or not, depends on the configured probability of keeping the current resources, hereafter referred as "probability of resource keep". In particular, once SLRRC reaches zero, the UE either keeps the previous selection or selects new resources based on the pre-configured probability value. The value of SLRRC is randomly selected from the interval $[5, 15]$ for $P_{\text{rsvp}} \geq 100$ ms. For $P_{\text{rsvp}} < 100$ ms, the value of SLRRC is randomly selected from the interval $\lceil 5 \times \frac{100}{\max(20, P_{\text{rsvp}})}, 15 \times \frac{100}{\max(20, P_{\text{rsvp}})} \rceil$ [179]. The standard also defines the maximum number of times that the same resource can be used for SPS through $C_{\text{resel}} = 10 \times \text{SLRRC}$, after which the resource reselection has to be triggered, independently of the probability of resource keep.

An illustration of the SPS procedure for NR V2X Mode 2 is shown in Fig. 7.5. In the example, three resources are selected within the resource selection window ($m$ in the figure is the slot index of the first selected resource), and these allocations are repeated every $P_{\text{rsvp}}$ for SLRRC times. Once the three transmissions in the interval starting at $m + (\text{SLRRC} - 1) \times P_{\text{rsvp}}$ have been carried out, either the same selection is kept or a new resource selection procedure is triggered, based on the probability of resource keep.

## 7.3 NR V2X Simulation Models

This section describes the *ns-3*-based NR V2X simulator that we have built, as an extension of the NR 5G-LENA open source network simulator [184], and the V2X models that we have developed, including the design choices and implementation details.

**Figure 7.7:** NR Sidelink UE data plane.

In our implementation, the data and control plane architectures of 5G-LENA UE nodes [184], shown in Fig. 7.6 and Fig. 7.7, have not been changed. Specifically, the implementation does not add or remove any of the protocol layers in the simulator. However, the extension to support NR V2X functionalities, e.g, sensing, SPS, autonomous resource selection, which are not part of a conventional UE design, involved some modifications at all layers of the protocol stack[2], from the Non-Access Stratum (NAS) and down to the PHY layer. One of the most important changes towards the implementation of NR V2X communications, compared to the typical cellular communications available in 5G-LENA simulator, is the introduction of sidelink, i.e., direct vehicle-to-vehicle communications. For that, the bearer establishment and RRC layer have been fully updated according to NR V2X RRC specification in TS 38.331 [175]. Also, MAC and PHY layers have been redesigned to implement NR V2X Mode 2 procedures using sensing-based SPS, as described in Section 7.2.2, according to TS 38.321 for the MAC layer [179], and TS 38.211 and TS 38.212 for the PHY layer [172,181]. In the simulator, as previously mentioned, we focus on NR V2X Mode 2 for out-of-coverage scenarios with broadcast communications and therefore, for the moment only blind retransmissions are considered.

---

[2]The NR 5G-LENA simulator reuses the upper layers, i.e., RLC and above, of LENA *ns-3* LTE module.

**Table 7.1:** Comparison of D2D, LTE C-V2X, and NR V2X in the standard and ns-3.

| | **D2D standard** | **D2D ns-3** [166] | **LTE C-V2X standard** | **C-V2X ns-3** [167] | **NR V2X standard** | **NR V2X ns-3** [this thesis] |
|---|---|---|---|---|---|---|
| Communication types | groupcast | broadcast | broadcast | broadcast | broadcast, groupcast, unicast | broadcast |
| MCS | QPSK, 16QAM | QPSK, 16QAM | QPSK, 16QAM, 64QAM | QPSK, 16QAM, 64QAM | QPSK, 16QAM, 64QAM, 256QAM | QPSK, 16QAM, 64QAM, 256QAM |
| Waveform | SC-FDMA | SC-FDMA | SC-FDMA | SC-FDMA | OFDMA | OFDMA |
| Frequency range | sub 6 GHz | sub 6 GHz | sub 6 GHz | sub 6 GHz | sub-6 GHz, mmWave | sub-6 GHz, mmWave |
| Subcarrier spacing | 15 kHz | 15 kHz | 15 kHz | 15 kHz | sub-6 GHz: 15, 30, 60 kHz, mmWave: 60, 120 kHz | sub-6 GHz: 15, 30, 60 kHz, mmWave: 60, 120 kHz |
| Duplexing modes | FDD, TDD | FDD | FDD, TDD | FDD | FDD, TDD | TDD |
| Retransmissions | blind | blind | blind | blind | broadcast: blind, groupcast: blind, feedback-based, unicast: blind, feedback-based | broadcast: blind |
| PHY channels | PSCCH, PSSCH, PSDCH, PSBCH | PSCCH, PSSCH, PSDCH, PSBCH | PSCCH, PSSCH, PSBCH | PSCCH, PSSCH, PSBCH | PSCCH, PSSCH, PSBCH, PSFCH | PSCCH, PSSCH |
| Control and data multiplexing | frequency, time | time | frequency | frequency | frequency, time | time |
| Scheduling interval | 1 subframe | 1 subframe | 1 subframe | 1 subframe | 1 slot | 1 slot |
| Sidelink modes | 1 and 2 | 1 and 2 | 3 and 4 | 4 | 1 and 2 | 2 |

Table 7.1 compares the main features of D2D, LTE C-V2X and NR V2X, as defined in the standard. Also, we compare the *ns-3*-based system level simulators available for D2D [166] and C-V2X [167], with the *ns-3* NR V2X simulator presented in this thesis. In Table 7.2, we detail the features and functionalities that are available in the developed NR V2X system-level simulator. The features listed for our simulator are those included in the first release, which allow the evaluation of a full NR V2X system with a subset of NR V2X features, but we plan to further progress with the module's development to support more extensions and functionalities.

<div align="center">**Table 7.2:** NR V2X models.</div>

| | NR V2X |
|---|---|
| Frame structure | TDD NR-compliant frame structure with slots and OFDM symbols of numerology-dependent length [184, 185]: <br> - frame: 10 ms, subframe: 1 ms <br> - each subframe has $2^\mu$ slots (associated to $15 \times 2^\mu$ kHz SCS) <br> - numerologies $\mu$=0, 1, 2, 3, 4 are supported <br> - each slot is composed of 14 OFDM symbols <br> Support for multiple bandwidth parts [184]: <br> - more than one BWP can be configured for sidelink <br> - each BWP can have pre-configured multiple sidelink resource pools, but only one pool can be active at a time |
| Duplexing mode | TDD <br> - the TDD pattern is flexible in length and composition, and can include downlink-only slots, uplink-only slots, or flexible slots (in which downlink and uplink transmissions can occur). An example of TDD pattern is [D F U U U]. |
| Sidelink resource pool | - sidelink transmission is only allowed in uplink-only slots, and whether an uplink slot is available for sidelink or not is specified through the SL bitmap. An example of the SL bitmap is [1 1 1 1 1 1 0 0 0]. <br> - within the uplink slots available for sidelink, the symbols available for sidelink are RRC pre-configured and our default structure is as follows: PSCCH can occupy 1 or 2 starting symbols, and depending on the PSCCH allocation, 2nd to 13th or 3rd to 13th symbols are available for PSSCH, and the 14th symbol is left empty as a guard period <br> - in frequency domain, RRC pre-configures the subchannel size (in number of RBs), and as per this configured size, divides the available bandwidth in number of available subchannels. |
| SL data/control channels | - PSSCH and PSCCH are multiplexed in time <br> - PSSCH and PSCCH are sent and received quasi-omnidirectionally at the UEs |
| Error models | NR PHY abstraction for PSSCH and PSCCH channels [186] including support for MCS Table1 and Table2 [173], MCS LDPC coding and block segmentation [181] |
| Modulation | OFDM |
| Channel Coding | LDPC |
| MCS | QPSK, 16-QAM, 64-QAM, 256-QAM |
| HARQ | NR PHY abstraction for HARQ includes support for HARQ-IR and HARQ-CC |
| Retransmissions | Blind retransmissions, including up to a pre-configured number with retransmission combining |
| Resource allocation | sensing-based and random resource selections are supported |
| Link adaptation | Fixed MCS |
| Antenna models | 3GPP-compliant [187]: <br> - Antenna arrays: 1 uniform planar array per UE, $M \times N$ antenna elements, no polarization <br> - Antenna elements: isotropical and directional radiation are supported |
| Channel models | 3GPP-compliant [188], supporting Urban grid and Highway scenarios, in both sub 6 GHz and mmWave bands |

## 7.3.1 NAS

The establishment and management of sessions occur at the highest layer on the control plane, the NAS. The current functionalities of the NAS layer in *ns-3* involve establishment of Evolved Packet System (EPS) bearers, multiplexing uplink data packets coming from the upper layers into the appropriate EPS bearer by using the Traffic Flow Templates

(TFTs) classifier. A TFT defines the rules for mapping IP packets to the right bearer based on IP addresses, ports, and type of service parameters.

For sidelink, the modifications are similar to the ones introduced in [166]. Specifically, NAS now supports the activation of sidelink bearers, mapping of Internet Protocol (IP) packets to the sidelink bearers based only on the IP destination address of the packets, and the transmission/reception of packets in NAS OFF state to support out-of-coverage scenarios.

## 7.3.2 RRC

The RRC is the control plane protocol in charge of setting important parameters for the session. The modifications in the RRC include the creation of the sidelink bearers upon receiving a notification from NAS, and the pre-configuration of UEs in an out-of-coverage scenario. As mentioned earlier, the model currently focuses on the broadcast communication, therefore, as per the standard, it supports the creation of uni-directional sidelink radio bearers [189].

Regarding the UEs' pre-configuration, the model implements all RRC Information Elements (IEs) needed to configure a UE [175]. This configuration is of key importance to perform sidelink communication when the gNB is absent. These IEs are mainly used for two purposes. The first is to configure the UE's PHY layer parameters, e.g., numerology, symbols per slot, bandwidth, and TDD pattern. The second is to provide the sidelink resource pool(s) information to MAC and PHY layers. It is also worth mentioning that the model allows the configuration of multiple BWPs for sidelink, where for each BWP, more than one resource pool can be configured through RRC. We note that, in spite of supporting multiple resource pools per BWP, only one pool could be active at one time [179]. Moreover, differently from the standard, which uses separate pools for transmission and reception [175], our model uses the same active pool for both.

## 7.3.3 PDCP

The changes introduced in Packet Data Convergence Protocol (PDCP) layer are in line with LTE sidelink [166]. In particular, when it comes to sidelink, it is no longer possible to uniquely identify a logical channel only based on its Logical Channel Identifier (LCID). With sidelink communications, UEs independently assign the LCIDs to logical channels for each destination (i.e., Layer 2 group ID) to which they are transmitting. Thus, it is impossible for UEs to identify the packets if multiple transmitting UEs select the same LCID for the same group. To solve this, two more identifiers, i.e., source Layer-2 ID and destination Layer-2 ID, are included to identify the transmitting UE [185].

## 7.3.4 RLC

The LTE Radio Link Control (RLC) layer in the simulator already supports the so called Unacknowledged Mode (UM), which is the RLC mode used for sidelink broadcast

communications. The only modifications made to the RLC layer are identical to the PDCP layer.

## 7.3.5 MAC

The UE's MAC layer has been extensively modified to transmit and receive sidelink transmissions. In the following, we explain these procedures in detail.

### 7.3.5.1 MAC transmitting procedure

In out-of-coverage scenarios, UEs are required to perform the autonomous resource selection following Mode 2, which could be based on sensing-based or random selection procedures, as explained in Section 7.2.2 and Section 7.2.1.6, respectively. The first significant modification introduced in this respect is the new MAC scheduler interface. This interface allows the implementation of sidelink UE-specific schedulers, which could assign resources following specific strategies, e.g., fixed MCS, adaptive MCS based on CSI, etc. The UE MAC layer is extended to provide all the information needed by a scheduler to perform resource selection. For example, the information related to all the Logical Channels (LCs) of destinations the UE is interested in transmitting to, the total number of available subchannels, $N_{\text{max\_reserve}}$, the maximum number of PSSCH transmissions, and most importantly, the RLC sidelink Buffer Status Reports (BSRs) of each LC that indicate how much sidelink traffic needs to be transmitted. The second important addition is the buffering of the sensing data reported by the UE's PHY layer. This buffer behaves like a sensing window at the time of a sensing-based resource selection. It contains the sensing information for the interval $[n - T_0, n - T_{\text{proc},0})$, where $n$ is the slot at which the resource selection is triggered, and $T_0$ is configured by the RRC while $T_{\text{proc},0}$ is a MAC layer parameter. In what follows we will dive into the details of UE's MAC layer operation to perform sensing-based resource selection.

At slot $n$, when a resource selection is triggered for a destination, the MAC layer draws a random counter (SLRRC) based on the user configured $P_{\text{rsvp}}$ value, which is used to compute $C_{\text{resel}}$. Then, for an active pool configured by the RRC, it computes the candidate resources (i.e., available slots) for sidelink transmission based on the selection window parameters, $T_1$, $T_{2,\text{min}}$, and $T_2$. Since the final resources must be selected based on sensing information, the MAC follows the procedure described in Section 7.2.2.1, to filter out the resources from the total available ones, which could be occupied by the other transmitting UEs. Once the filtered candidate resources' list is ready, the MAC layer forwards it to the scheduler. Our model provides a sample scheduler, which as per the standard [175], first randomly selects a number of slots, i.e., $N_{\text{selected}}$, for sidelink transmissions. The number of $N_{\text{selected}}$ slots depends on the number of slots that are available in the filtered list and the maximum number of configured PSSCH transmissions. If $K$ denotes the total number of available slots, and $N_{\text{PSSCH,maxTx}}$ is the maximum number of PSSCH configured transmissions, then:

$$N_{\text{selected}} = \begin{cases} N_{\text{PSSCH,maxTx}}, & \text{if } K \geq N_{\text{PSSCH,maxTx}} \\ K, & \text{otherwise} \end{cases} \tag{7.1}$$

After randomly selecting the required number of slots, it randomly the selects the required number of contiguous subchannels computed using a fixed MCS strategy.[3] After selecting the $N_{\text{selected}}$ number of slots, the scheduler computes the Transport Block Size (TBS) using the fixed MCS by taking into account the BSR of a LC, and the 5 bytes overhead of 2nd-stage-SCI, which needs to be multiplexed with data. Finally, it prepares a sidelink allocation valid for the first resource reservation period deciding also aspects like which slots from the $N_{\text{selected}}$ have to carry the 1st-stage-SCI, the New Data Indicator (NDI), and the Redundancy Version (RV) number of each slot. The UE MAC layer, upon receiving this allocation, creates the SPS grants based on the configured value of $P_{\text{rsvp}}$ and the already drawn counters, i.e., SLRRC and $C_{\text{resel}}$. After using these grants for a number of transmissions equal to the SLRRC, a resource reselection is triggered. That is, once SLRRC reaches zero, the UE either keeps the previous selection or selects new resources based on the pre-configured probability of resource keep. Finally, if $C_{\text{resel}}$ reaches zero, the resource reselection is triggered, independently of this probability. As already discussed in Section 7.2.1.6, our model also supports a random resource selection. The only difference between the two approaches is that the random resource selection procedure does not filter the slots from the available ones before giving the list to the scheduler, so that the sensing information is not used.

Before forwarding the sidelink packets to the PHY layer, a check is performed at the beginning of each slot to ensure the availability of a valid grant for that slot. If there is, the MAC layer prepares two packet bursts, one for the 1st-stage-SCI and the second for the 2nd-stage-SCI plus data, and assigns a HARQ process ID to the data packet. It also saves this data packet into a HARQ buffer if blind retransmissions are configured. We note that the model allows to configure multiple (no limit for research purposes) sidelink/HARQ processes to allow continuous flow of data. After this, both packet bursts are forwarded to the lower layer. Upon receiving these packet bursts, the PHY, places them in a queue to be transmitted on the configured PSCCH and PSSCH symbols.

### 7.3.5.2   MAC receiving procedure

The UE's MAC layer, upon receiving the PSSCH packet burst from the PHY, first retrieves the 2nd-stage-SCI to read the source Layer-2 ID and the destination Layer-2 ID of the received packet. As mentioned in section 7.3.3, these identifiers are used to map the received packet to its logical channel. If a bearer for the received packet is already established, the data packet is forwarded to the upper layers. Otherwise, the MAC asks the RRC to establish the bearer for the reception. Once this is done, the packet is forwarded to the upper layers.

---

[3]Note that adaptive MCS strategy makes sense in unicast and groupcast communications when the CSI from the receiving UE can be acquired through the PSSCH. As in this implementation we have focused on broadcast communications, without PSFCH, the fixed MCS strategy is the adequate one. However, the implemented scheduler interface is general enough to accommodate more sophisticated schedulers.

## 7.3.6 PHY

Similarly to the MAC layer, the PHY functionality can also be divided into transmitting and receiving procedures. In the following, we describe them in detail.

### 7.3.6.1 PHY transmitting procedure

The 5G-LENA simulator accurately models (as per the standard) the numerology-dependent slot and OFDM symbol granularity. The state-machine of the PHY layer is mainly determined by the definition of the concept of start slot event and variable TTI [62]. When the start slot event is triggered, the processing follows a logical order that involves the MAC and then the scheduler, before returning the control to the PHY. For sidelink, once the control gets back to the PHY, the PHY checks if the MAC has provided an allocation for the current slot. This allocation further consists of variable TTI allocations. The variable TTI means that the number of allocated symbols to physical sidelink channels (i.e., PSCCH and PSSCH) is variable, based on the sidelink configuration. Upon finding the allocation for the slot, the PHY layer transmits PSCCH and PSSCH PDUs on their respective symbols whose duration depend on the configured numerology and Cyclic Prefix (CP).

### 7.3.6.2 PHY receiving procedure

To receive the sidelink transmissions, one of the key enhancements of the PHY is the introduction of the ability to handle collisions/interference, also introduced in [166]. The interference model available in the 5G-LENA simulator was designed for a typical cellular communication. Its design assumes that a UE is interested in transmitting or receiving only from its serving gNB, and assumes no interference from the UEs served by the same gNB. Transmissions from other gNBs/UEs are simply considered as interference. In case of sidelink, especially in broadcast or groupcast, a UE is interested in transmitting to or receiving from multiple surrounding UEs. In this context, UEs in out-of-coverage scenarios or "UE-selected" mode can select the same (or overlapping) resources because the allocation is uncoordinated. Therefore, to determine which packet will be successfully decoded, the new implementation keeps track of the SINR values for each sidelink transmission.

As described earlier, currently our model supports the transmission and reception of timely multiplexed PSCCH and PSSCH. Thus, the PHY first receives signal(s) (i.e., 1st-stage-SCI) transmitted over PSCCH. This signal is used for two purposes: 1) to measure the RSRP required for the sensing-based resource selection, 2) to retrieve the information about the possible PSSCH transmission and retransmissions. The RSRP is computed using the 3 Resource Elements (REs) per RBs, carrying the 1st-stage-SCI, since the simulator does not explicitly include PSCCH DMRS. Moreover, for the sensing-based resource selection, the PHY measures the RSRP of each correctly decoded 1st-stage-SCI, from all the surrounding UEs. On the other hand, after computing the RSRP, if it is from the transmitter of interest, it reads the information encoded in the 1st-stage-SCI to receive the PSSCH transmission and its possible retransmissions.

Concerning the error model used for the reception of PSCCH and PSSCH transmission, we use the existing data plane error model in 5G-LENA [186], since the MCSs defined for PSSCH are the same as the ones defined for PDSCH/PUSCH. Also, we adopt such an error model for the PSCCH, using MCS0.

### 7.3.7 Channel Models

TR 37.885 [188] defines the system-level evaluation methodology for 5G V2X use cases, including the description and modeling of scenarios, deployment, mobility, antenna, traffic, and channel models. For channel modeling, TR 37.885 extends the geometry-based stochastic channel modeling framework introduced in TR 38.901 [190] for typical cellular communications, by adding the possibility to model wireless channel in vehicular environments and sidelink communications in which both the transmitter and the receiver are in motion. Two key scenarios are used for NR V2X evaluation [188]:

- Urban grid, which targets urban environments with a grid of buildings and roads with four lanes (two in each direction) between the buildings, and

- Highway, which targets highway environments with a highway composed of a total of six lanes, considering three lanes in each opposite direction.

For each scenario, TR 37.885 specifies new channel condition models, propagation models, and fast fading parameters capturing the characteristics of each environment.

The developed ns-3 NR V2X module includes the channel and antenna models for both V2X Urban grid and Highway scenarios, as defined in [188].

## 7.4 NR V2X Evaluation Campaigns

This section presents the simulation scenario that we have used to assess NR V2X performance. Then, we present multiple simulation campaigns and discuss the obtained end-to-end results.

### 7.4.1 Scenario and Definition of Neighbor

We consider a V2X Highway scenario, as defined in 3GPP TR 37.885 [188]. The deployment is composed of multiple lanes in a 3.9 km highway road, with an inter-lane distance of 4 m. Within each lane, the inter-vehicle distance is 78 m, which is computed using the formula $max(2, 2 \times$ average speed m/s) defined in [188]. The UE dropping is implemented according to [188] Option A, in which all vehicles (100 %) are of Type 2 (i.e., passenger vehicle with an antenna height of 1.6 m), clustered dropping is not used, and the vehicle speed is set to 140 km/h in all the lanes. We consider 3 lanes with vehicles moving in the same direction, and 50 vehicles per lane. We focus on an out-of-coverage
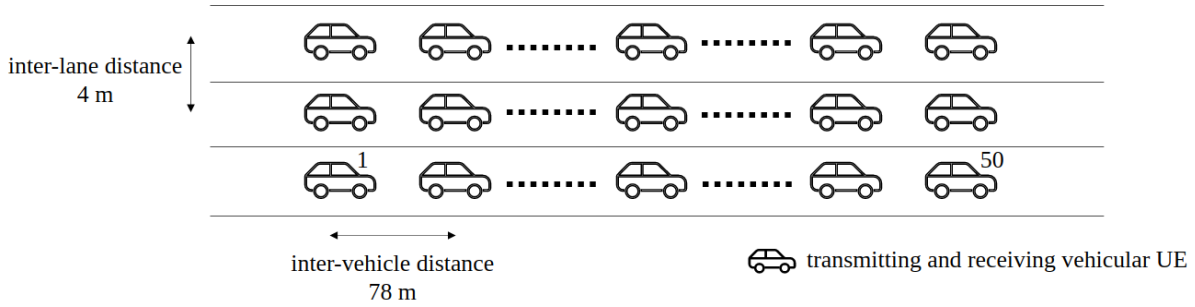
**Figure 7.8:** Highway scenario with 3 lanes and 50 vehicular UEs per lane moving at a speed of 140 km/h and spanning over 3.9 km.

scenario, so that gNBs are disabled in the evaluation [188]. The considered deployment scenario is shown in Fig. 7.8.

We focus on a use case that targets the broadcast of basic service messages and by taking the inspiration from [191] we assume that all vehicular UEs are half duplex[4] transceivers, which have the same packet size, generated at the same rate, and using a fixed MCS. Transmission is done over the 5.9 GHz band, assuming a channel bandwidth of 10 MHz [92]. The traffic model is characterized by periodic packet transmissions, with a packet size of 300 bytes, which are transmitted every 100 ms. This leads to a data rate of 24 kbit/s.

Moreover, in the considered scenario, each vehicular UE is a potential receiver. However, as per 3GPP [188], for the broadcast scenario, the KPIs, e.g., Packet Inter-reception Delay (PIR) and Packet Reception Ratio (PRR), for each UE must be computed by considering only those UEs that are located within a specific range of a certain distance from it, which is known as the "awareness range". We consider an awareness range of 200 m, and we characterize as **neighbors** all those vehicular UEs located within such range from the source UE [191]. We also consider a throughput KPI, which according to its definition in 3GPP standard is computed without considering any awareness range, and is defined in the next subsection [188].

Table 7.3 reports the simulation parameters and functionalities, for NR V2X end-to-end evaluations. Through the simulation campaigns, we study the impact of specific NR V2X parameters, which are listed in Table 7.3 as variations of the baseline configuration (last column).

## 7.4.2 Simulation campaigns

The simulation campaigns are classified into two main blocks. Firstly, in Section 7.4.3 we study the impact of various NR V2X parameters on the performance of the sensing-based resource selection, which by the 3GPP standard, is the default mode of operation for V2X UEs. In particular, we discuss a set of simulation campaigns where we study the impact of the following parameters:

---

[4]In the simulations presented in this chapter, we use sidelink V2X operating band n47, which is a TDD band, therefore, all the UEs in our simulation use half duplex as specified by the 3GPP standard [97].

**Table 7.3:** Main scenario simulation parameters (baseline configuration and its variations).

| Parameter | Value (baseline) | Value (variations) |
|---|---|---|
| **Deployment and propagation parameters**: | | |
| Channel model | 3GPP Highway | |
| Deployment | 3 lanes, 5 vehicles per lane | |
| Carrier frequency | 5.89 GHz | |
| Channel bandwidth | 10 MHz | |
| Noise power spectral density | -174 dBm/Hz | |
| UE antenna height | 1.6 m | |
| UE speed | 140 km/h | |
| **Traffic parameters:** | | |
| Application packet size | 300 Bytes | |
| Inter-packet arrival time | 100 ms | |
| Application load | 24 kbit/s | |
| **Device parameters:** | | |
| UE antennas | uniform planar array 1x2 | |
| UE transmit power | 23 dBm | |
| UE noise figure | 5 dB | |
| **NR V2X parameters and functionalities:** | | |
| Frame structure | $\mu$=0 (SCS=15 kHz) | $\mu$=1 (SCS=30 kHz), $\mu$=2 (SCS=60 kHz) |
| TDD pattern | [D D D F U U U U U U] | |
| Sidelink bitmap | [1 1 1 1 1 1 0 0 0 1 1 1] | |
| Subchannel size ($N$) | 10 RBs | |
| PSCCH symbols | 1 | |
| PSSCH symbols | 12 | |
| Link adaptation | fixed MCS | |
| MCS PSSCH | MCS 14 (MCS Table1) | MCS 4, MCS 7, MCS 20, MCS 28 |
| MCS PSCCH | MCS 0 (MCS Table1) | |
| Error model | NR PHY abstraction based on EESM [186] for PSSCH and PSCCH | |
| Number of PSSCH transmissions ($N_{\mathrm{PSSCH,maxTx}}$) | 5 | 2, 10 |
| HARQ combining method | HARQ incremental redundancy | |
| MAC resource selection | sensing-based resource selection | random resource selection |
| RLC mode | RLC-UM | |
| RLC buffer size | 999999999 Bytes | |
| **NR V2X Mode 2 parameters:** | | |
| Sensing window ($T_0$) | 100 ms | |
| $T_2$ | 33 slots | 17 slots, 65 slots |
| $T_1$ | 2 slots | |
| $T_{\mathrm{proc,0}}$ | 2 slots | |
| Percentage of resources must be selected in a selection window | 20 % | |
| Max num per reserve ($N_{\mathrm{max\_reserve}}$) | 3 | 2 |
| Probability of resource keep | 0 | 0.5, 0.8 |
| Resource reservation period ($P_{\mathrm{rsvp}}$) | 100 ms | |
| RSRP threshold | -128 dBm | |

- NR V2X numerology ($\mu$),

- NR V2X number of PSSCH transmissions of the same MAC PDU (including initial transmission and blind retransmissions) ($N_{\text{PSSCH,maxTx}}$),

- NR V2X Mode 2 selection window length ($T_2$)

- NR V2X maximum number of resources per reservation ($N_{\text{max\_reserve}}$),

- NR V2X Mode 2 probability of resource keep, and

- NR V2X MCS index for PSSCH.

Secondly, in Section 7.4.4, we focus on comparing the performance of the sensing-based and random resource selection procedures (both considered in 3GPP for NR V2X Mode 2). This evaluation will demonstrate the simulator's capability to support both the standardized resource selection procedures and the performance gain that the sensing-based resource selection can provide over the random resource selection. In this case, we consider a concrete system configuration corresponding to the baseline configuration shown in Table 7.3.

For each simulation campaign, 20 random channel realizations are performed, to get statistical significance. A single simulation has the duration of 10 simulated seconds. The constant bit rate applications start randomly within an interval of 100 ms, and run without interruption for 10 seconds.

As output statistics, we focus on the three KPIs defined for V2X evaluations in 3GPP [188], measured at the application layer:

- PIR: interval of time elapsed between two successful packet receptions of packets transmitted by a specific neighbouring UE. We consider the average PIR, averaging over the different successful receptions for each transmit-receive UE pair in the reception range. PIR is a range-based KPI, as per [188].

- PRR: for each packet transmitted by a UE, a ratio of the number of neighboring UEs that successfully receive that packet over the total number of neighboring UEs. We consider the average PRR, averaging over the different packets transmitted by a transmitting UE. PRR is a range-based KPI, as per [188].

- Throughput: total number of correctly received bytes over the simulation time, measured at the application layer, for each transmit-receive UE pair. As per [188], throughput is not range-based, so we consider all throughput values for those UEs that received some data bytes.

For each of the output statistics, we represent the CDF, over the different simulation runs. For each simulation campaign, we show three figures, one for each of the above mentioned output statistics, i.e., (a) PIR, (b) PRR, (c) throughput.

### 7.4.3 Simulation results: Sensing-base resource allocation

#### 7.4.3.1 Impact of numerology

In the first simulation campaign, we evaluate the impact of different NR numerologies. In these tests, we consider three different numerologies: $\mu = 0$ (15 kHz SCS), $\mu = 1$ (30 kHz SCS), and $\mu = 2$ (60 kHz SCS), which are the three numerologies supported in NR standard for sub 6 GHz bands. Our comparison assumes that the same bandwidth is available for all the tested numerologies, which is also the common assumption in 3GPP evaluations. The numerologies are displayed in the legends of the figures as mu-0, mu-1, and mu-2, respectively. Fig. 7.9 shows the CDF statistics of the PIR, PRR, and throughput.

In Fig. 7.9.(a)-(b), we observe that the curves of PIR and PRR for the tested numerologies cross each other in different regions of the plot. The reason is a type of flexibility introduced by each numerology in terms of the number of available subchannels in the frequency domain, and different slot duration in the time domain. In NR, the processing times and the transmission durations are inversely proportional to the SCS, i.e., for a given bandwidth a lower SCS provides higher number of RBs while a higher SCS implies lower timings (i.e., shorter slot duration). For example, in our scenario with 10 MHz bandwidth and a subchannel size of 10 RBs, there are 5, 2, and 1 subchannel(s) available when using $\mu = 0$, $\mu = 1$, and $\mu = 2$, respectively[5]. Therefore, with the considered packet size of 300 bytes, which occupies only one subchannel[6] a lower numerology, e.g., $\mu = 0$ provides the maximum flexibility in the frequency domain, i.e., at a given time 5 UEs can occupy a single slot, which could reduce the collisions in the scenario. On the other hand, increasing the numerology, comes at the cost of a lower number of subchannels, so that we achieve a shorter slot length (as it is inversely proportional to the SCS), which makes that the resulting resource selection window length (set to 32 slots) results in a lower resource selection window in ms. For example, a 32 slots selection window results in 32 ms with $\mu = 0$, 16 ms with $\mu = 1$, and 8 ms with $\mu = 2$. Interestingly, this reduction in terms of ms of the selection window, results in a reduced probability of overlapping of the resource selection windows of different transmitting UEs, because the resource reservation period is defined in ms and so it remains fixed independently of the numerology. This increases the gap in slots/time between the end of the selection window and the beginning of the following reservation period and consequently it helps reduce the probability of collisions between UE's selection windows. However, when increasing the numerology from $\mu = 0$ to $\mu = 1$ we do reduce the number of available subchannels from 5 to 2 but we are also halving the selection window from 32 ms to 16 ms, which results in a very similar performance. In fact, having more subchannels with $\mu = 0$ increases the PRR for some UEs thanks to the fact that UEs can better exploit diversity in frequency domain (see the tail of the PRR in Fig. 7.9.(b)). In this sense, we observe that for the tested scenario reducing the selection window length to half is not enough to achieve an appreciable performance gain. Differently, when the numerology is further increased to $\mu = 2$, we can observe a performance gain achieved as a consequence of a

---

[5]Number of available RBs for a given bandwidth and 12 subcarriers per RB can be computed as: $(BW(Hz) - overhead(Hz))/SCS(Hz) \times 12$. Here, we considered an overhead of 4% of the bandwidth, which is a typical value used in NR.

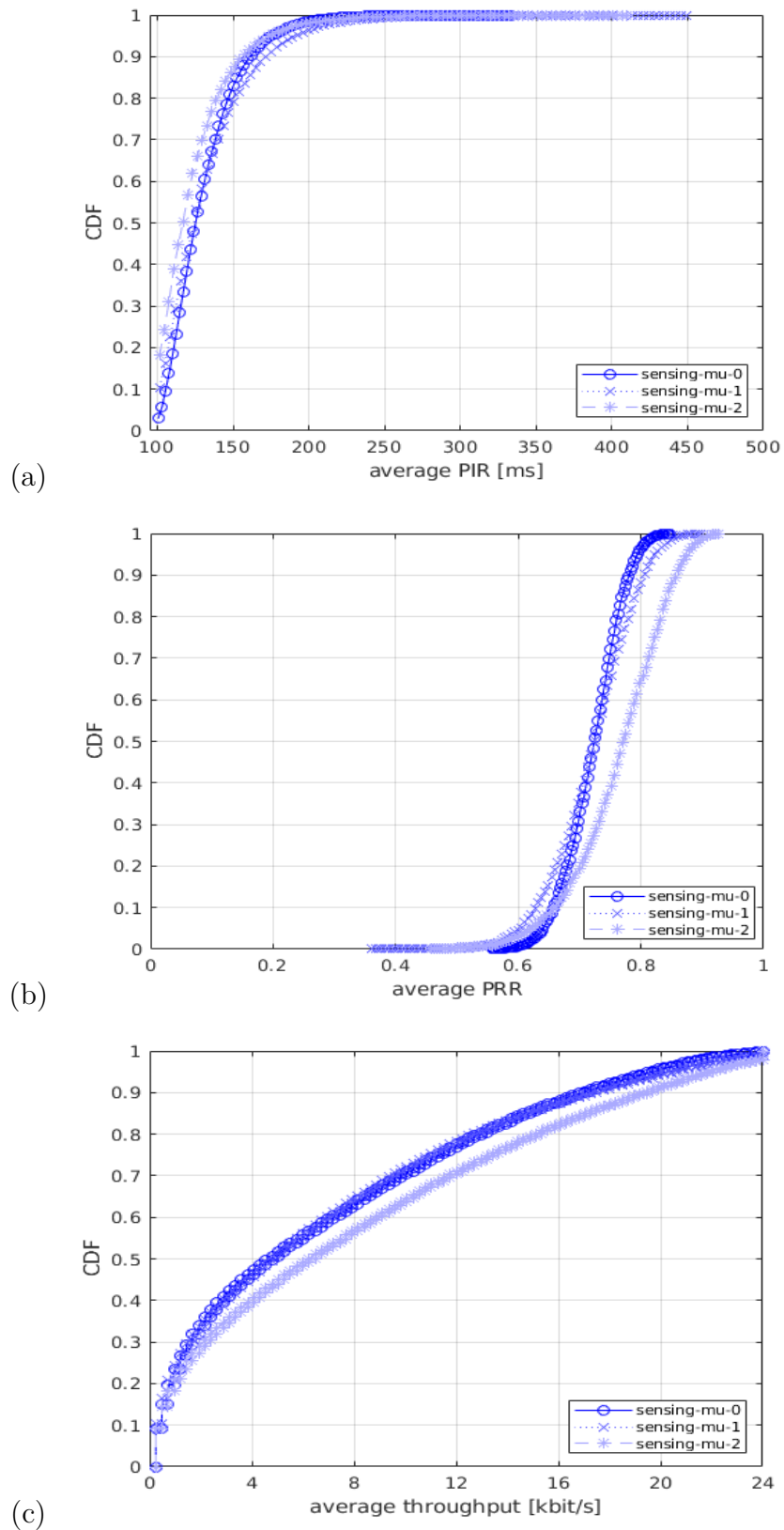[6]MCS 14 of NR MCS Table 1 and 10 RB long subchannel give a TBS of 348 bytes.

(a)

(b)

(c)

**Figure 7.9:** Impact of NR V2X numerology ($\mu$). (a) PIR (ms), (b) PRR, (c) throughput (kbps).

lower selection window of 8 ms (i.e., 4 times lower than $\mu = 0$) as it is clearly shown in Fig. 7.9.(b)-(c). It is important to note that, even if the resource selection is based on sensing, the procedure requires that 20 % of the resources available for SL are candidates to perform a resource selection. If this requirement is not fulfilled, the RSRP threshold is increased by 3 dB until this condition is met. Because of that, collisions can still occur. Ultimately, the number of incorrect PSSCH receptions may imply packet losses, if they can not be recovered by HARQ. This is effectively observed in the PRR and PIR statistics (see Fig. 7.9.(a)-(b)), for which, again, $\mu = 2$ is observed to offer better performance, as it allows reduced packet collisions and a larger number of successfully decoded packets, as a consequence of the shorter slot duration.

All in all, the trade-off of having more subchannels with a lower numerology versus a lower resource selection window length with higher numerology results in comparable performance when using the three numerologies. However, using $\mu = 2$ results in a lower slot duration and helps more UEs to achieve higher PRR, higher throughput, and lower PIR (which appears due to a lower probability of overlapping resource selection windows of different transmitting UEs) compared to $\mu = 0$ and $\mu = 1$.

### 7.4.3.2  Impact of number of PSSCH transmissions

In the second simulation campaign, we assess the impact of using different numbers of PSSCH transmissions of the same MAC PDU. This parameter is also known as $N_{\text{PSSCH,maxTx}}$ and, in case of SPS with blind retransmissions, it corresponds to the maximum possible number of resources that can be selected by the resource selection procedure. In our tests, we use $N_{\text{PSSCH,maxTx}} = 2$, 5, and 10. This includes one initial transmission and $N_{\text{PSSCH,maxTx}} - 1$ blind retransmissions. They are displayed in the legends of the figures as retx-2, retx-5, and retx-10, respectively. Fig. 7.10 shows the CDF statistics of the PIR, PRR, and throughput.

As shown in Fig. 7.10, a lower number of PSSCH transmissions is beneficial in terms of all the performance indicators. The reason is that in the considered scenario, characterized by good propagation conditions, a lower $N_{\text{PSSCH,maxTx}}$ does not saturate the resources, which helps the sensing-based procedure at the transmitting UEs to fully exploit the flexibility in terms of the number of subchannels, i.e., 5 subchannels to choose from at each slot. After properly filtering the resources based on sensing, a lower $N_{\text{PSSCH,maxTx}}$, offers an improvement in PIR, PRR, and throughput compared to higher values of $N_{\text{PSSCH,maxTx}}$, which is due to the fact that the sensing procedure is properly avoiding collisions and a small number of PSSCH transmissions is enough to decode the packets because of the good propagation conditions. This definitely demonstrates the effectiveness of the sensing-based resource selection in vehicular scenarios.

In summary, in a scenario with good propagation conditions where it is likely that UEs are able to sense each other, a low number of PSSCH transmissions shows better performance for all the indicators when using sensing-based resource selection. The result changes when sensing is not activated. A similar campaign has also been conducted to study the impact of $N_{\text{PSSCH,maxTx}}$ for the non sensing case. The results are not shown here for the sake of brevity, but demonstrate that with a random resource selection, more PSSCH
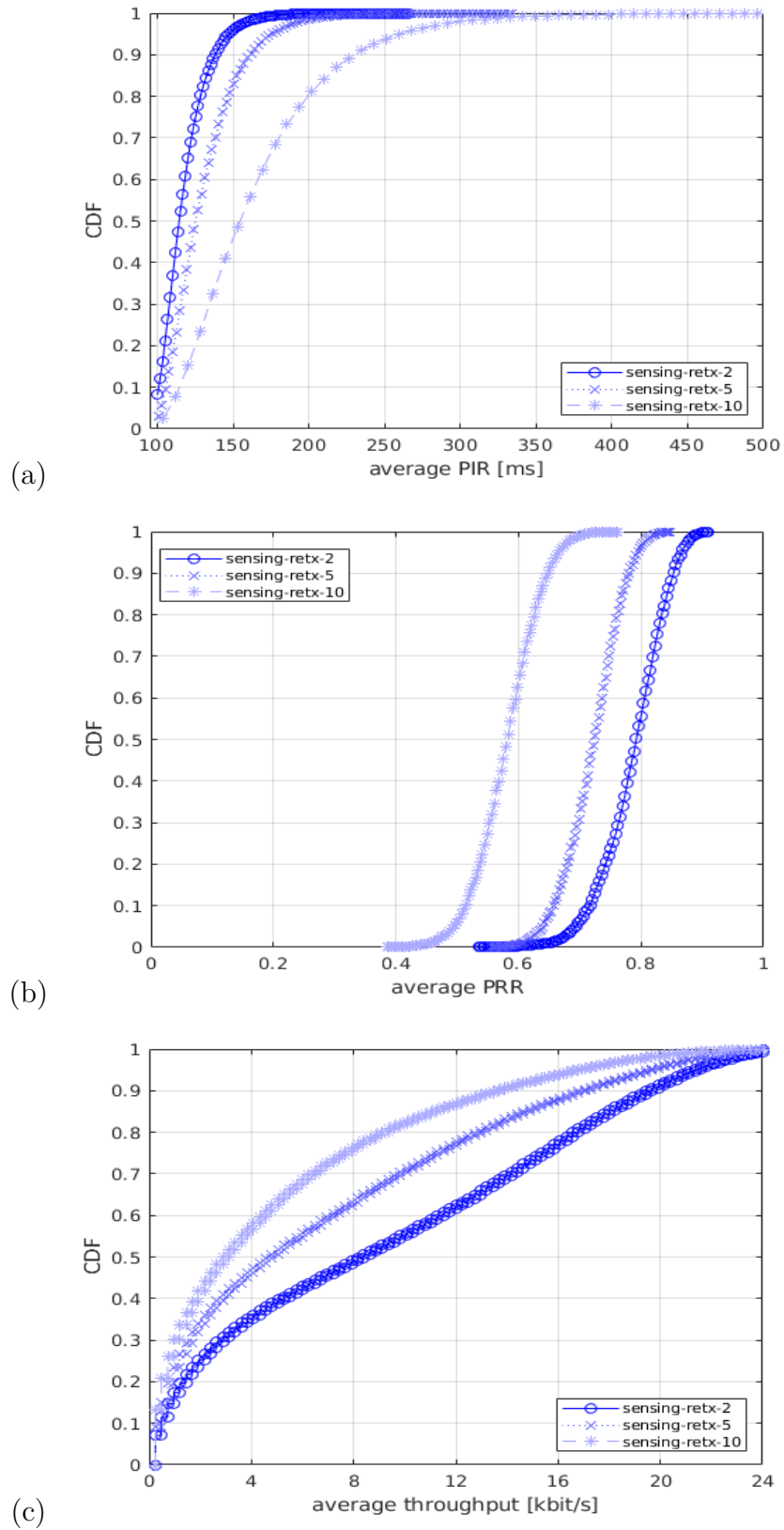
(a)



(b)



(c)

**Figure 7.10:** Impact of NR V2X number of PSSCH transmissions of the same MAC PDU ($N_{\mathrm{PSSCH,maxTx}}$). (a) PIR (ms), (b) PRR, (c) throughput (kbps).

retransmissions of the same MAC PDU are needed to properly decode the packets and improve the throughput, PRR, and PIR performances.

### 7.4.3.3 Impact of resource selection window length

In the third simulation campaign, we evaluate the impact of different selection window lengths. To do so, we select $T_1 = 2$ slots and vary $T_2$ values. These two parameters, determine the start and the end point of the resource selection window. For $T_2$, we consider three different values: $T_2 = 17$ slots, $T_2 = 33$ slots, and $T_2 = 65$ slots, which result into a selection window length $(T_2 - T_1 + 1)$ of 16 slots, 32 slots, and 64 slots, respectively. Fig. 7.11 shows the results in terms of the PIR, PRR, and throughput.

In terms of PIR and PRR, the impact of different selection window length is not linear and different effects can be highlighted from the analysis of the results. On the one hand, a lower $T_2$ causes PRR to decrease for almost 50% of the UEs, as observed in Fig. 7.11.(b), which is also reflected by higher PIR values in some cases in Fig. 7.11.(a). The reason is that a lower $T_2$ generates more collisions due to the reduced number of slots in the resource selection procedure to select from. On the other hand, a higher $T_2$ improves the performance, both in terms of PIR and PRR because a larger selection window provides more freedom to better randomize the resources among the various UEs. However, in some cases, the performance is very similar to what is observed with a lower $T_2$. This is because by increasing the resource selection window length in ms, we also increase the probability of overlapping of the resource selection windows of different transmitting UEs, which ultimately results in more collisions. Due to the combination of these divergent effects, we observe that the curves of different $T_2$ cross showing that in some cases a lower $T_2$ can be beneficial, while in others it is not. In particular, we observe that high-PRR UEs benefit from $T_2 = 17$ slots, while low-PRR UEs from $T_2 = 65$ slots.

In general, a larger resource selection window shows benefits in terms of throughput, PIR, and PRR metrics, at the cost of a slight increase of the PIR of some UEs. As a result, it may be better to have more resources to select, even if for some specific cases, more collisions may appear because of the longer resource selection window length.

### 7.4.3.4 Impact of maximum number of resources per reservation

In the fourth simulation campaign, we vary the maximum number of resources per reservation ($N_{\text{max\_reserve}}$). We consider the two values permitted in the NR V2X standard: $N_{\text{max\_reserve}} = 2$ and $N_{\text{max\_reserve}} = 3$. They are displayed in the legends of the figures as maxReserve-2 and maxReserve-3, respectively. Fig. 7.12 shows the results in terms of the PIR, PRR, and throughput.

The end-to-end performance obtained using the two values are very similar, but when considering $N_{\text{max\_reserve}} = 3$ the UEs experiences slightly higher PIR and slightly lower throughput, compared to $N_{\text{max\_reserve}} = 2$, as shown in Fig. 7.12.(a) and Fig. 7.12.(c), respectively. Also, the achieved PRR is higher with $N_{\text{max\_reserve}} = 2$ than with $N_{\text{max\_reserve}} = 3$ (see Fig. 7.12.(b)). The reason is that when there is a loss in the PSCCH channel, the UE has more opportunities to correctly receive and decode the
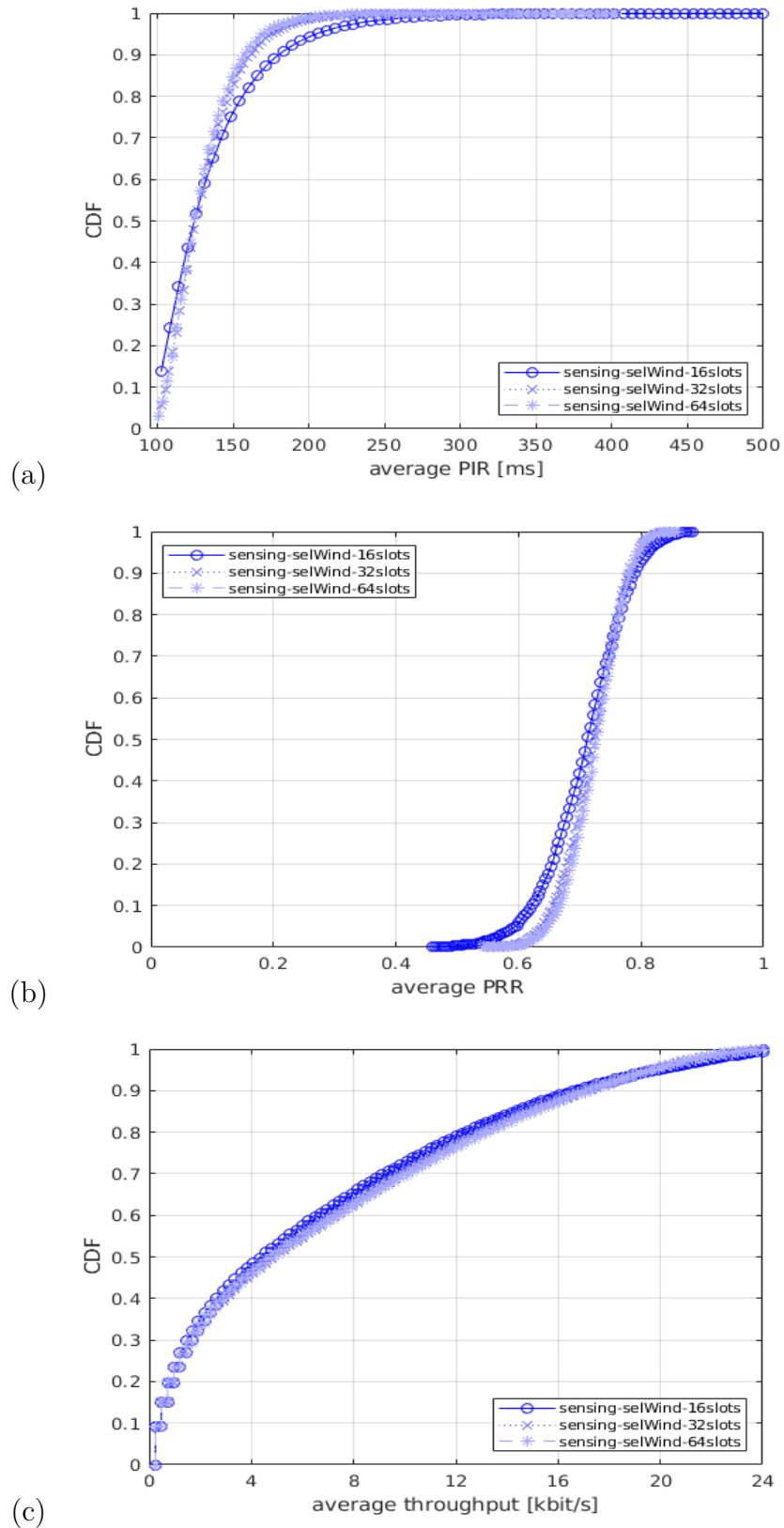
(a)

(b)

(c)

**Figure 7.11:** Impact of NR V2X selection window $(T_2 - T_1 + 1)$. (a) PIR (ms), (b) PRR, (c) throughput (kbps).
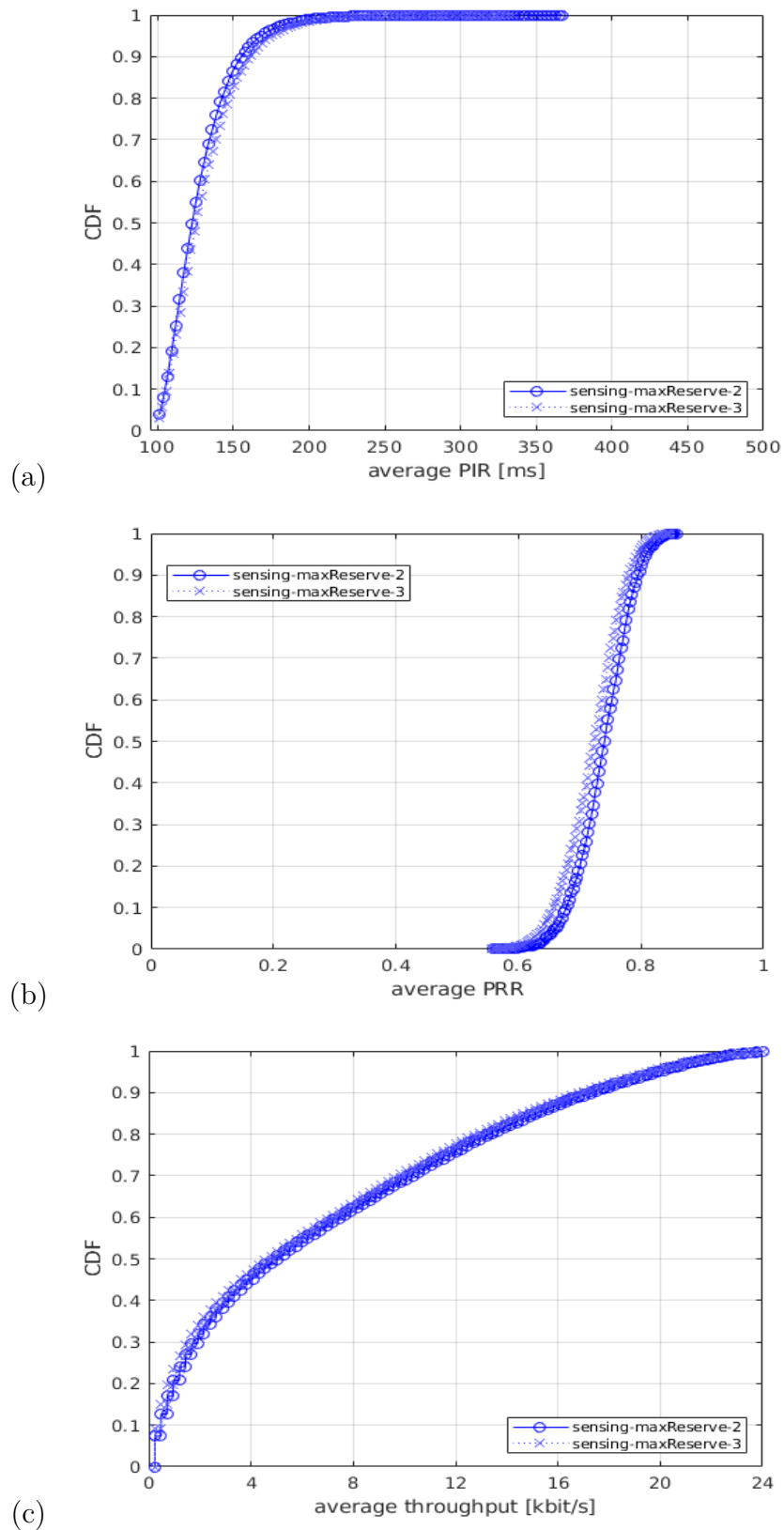
(a)

(b)

(c)

**Figure 7.12:** Impact of NR V2X maximum number of resources per reserve ($N_{\mathrm{max\_reserve}}$). (a) PIR (ms), (b) PRR, (c) throughput (kbps).

1st-stage-SCI using a retransmission with $N_{\text{max\_reserve}} = 2$, compared to the case of $N_{\text{max\_reserve}} = 3$, because with a lower $N_{\text{max\_reserve}}$ there are more slots that carry 1st-stage-SCI messages in a resource reservation period. For example, in a simulation with $N_{\text{PSSCH,maxTx}} = 5$ and $N_{\text{max\_reserve}} = 3$, there are only two 1st-stage-SCI messages transmitted in a resource reservation period (see Fig. 7.4), while with $N_{\text{max\_reserve}} = 2$, there are three 1st-stage-SCI messages. Specifically, if a UE fails to decode two 1st-stage-SCIs, it has no more opportunities to decode the data with $N_{\text{max\_reserve}} = 2$. However, in the same situation, with $N_{\text{max\_reserve}} = 2$, the receiving UE has another opportunity to decode the third 1st-stage-SCI and, potentially, the associated sidelink data. In addition, a larger $N_{\text{max\_reserve}}$ value may cause the loss of sensing information because of the inherent characteristic of the sensing-based procedure, and, hence, more collisions. This is why, we observe a slightly better performance in terms of PIR, PRR, and throughput, when using $N_{\text{max\_reserve}} = 2$. However, let us note that said situation does not happen very often in our scenario, and that is why the difference between the performance is not much significant.

Consequently, the sensing-based resource selection is shown to be slightly more efficient when using a lower number of $N_{\text{max\_reserve}}$, from all the metrics, because it allows for a more accurate sensing procedure. Particularly, it allows to detect more 1st-stage-SCIs from neighbor UEs, and so perform a better resource selection by excluding a larger set of slots that are being occupied by neighbor UEs. Also, for the receiving UEs, it provides more opportunities to decode the data.

### 7.4.3.5  Impact of probability of resource keep

In the fifth simulation campaign, we study the impact of the probability of keeping the resources during reselection, by testing different values. As explained in Section 7.2.2.2, once SLRRC reaches zero, the UE either keeps the previous selection or selects new resources based on the pre-configured probability value. We consider three values, 0, 0.5, and 0.8. We note that, 0 and 0.8 are the standard minimum and the maximum values for the probability of keeping the resources [175]. They are identified in the legends of the figures as ProbResKeep-0, ProbResKeep-0.5, and ProbResKeep-0.8, respectively. Fig. 7.13 shows the results in terms of the PIR, PRR, and throughput. Let us note that for a transmitting UE, ProbResKeep-0 triggers the resource reselection procedure more often, as compared to ProbResKeep-0.5 and ProbResKeep-0.8.

Interestingly, we can see a trade-off in the obtained PIR and PRR performance when using different probabilities of keeping the resources. On the one hand, UEs that have correctly selected the resources (meaning their selection is not colliding with other UEs and its obtained KPIs are good enough) are benefited by not switching the selection, i.e., higher values of probability of resource keep are better. On the other hand, the UEs that did a wrong resource selection at the beginning (i.e., selected resources that may collide with other UEs' selections) can benefit by reselecting the resources more often, i.e., by using lower values of probability of resource keep. For this reason, we see the crossing of the curves for PIR, PRR, and throughput curves.
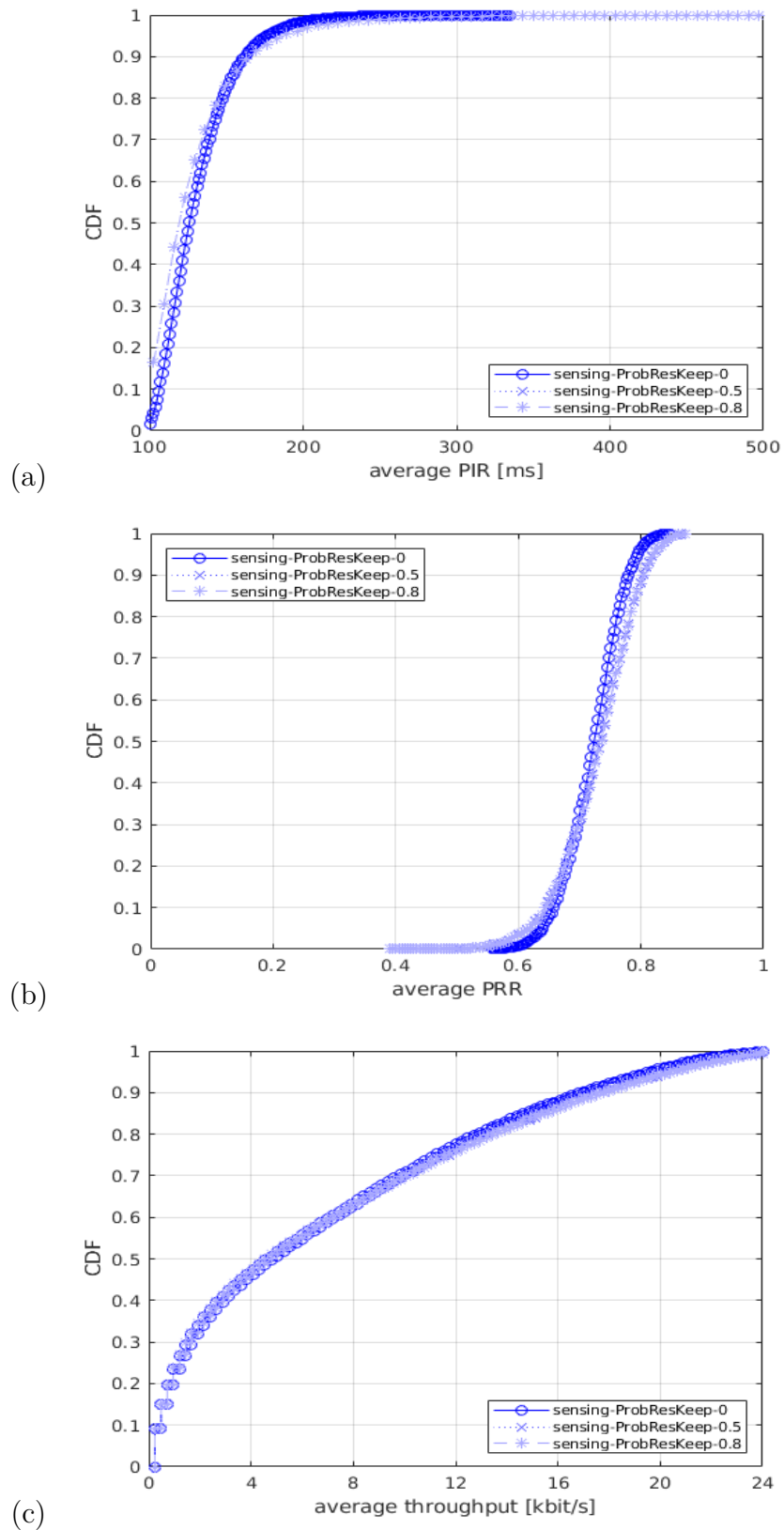
(a)

(b)

(c)

**Figure 7.13:** Impact of NR V2X probability of resource keep. (a) PIR (ms), (b) PRR, (c) throughput (kbps).
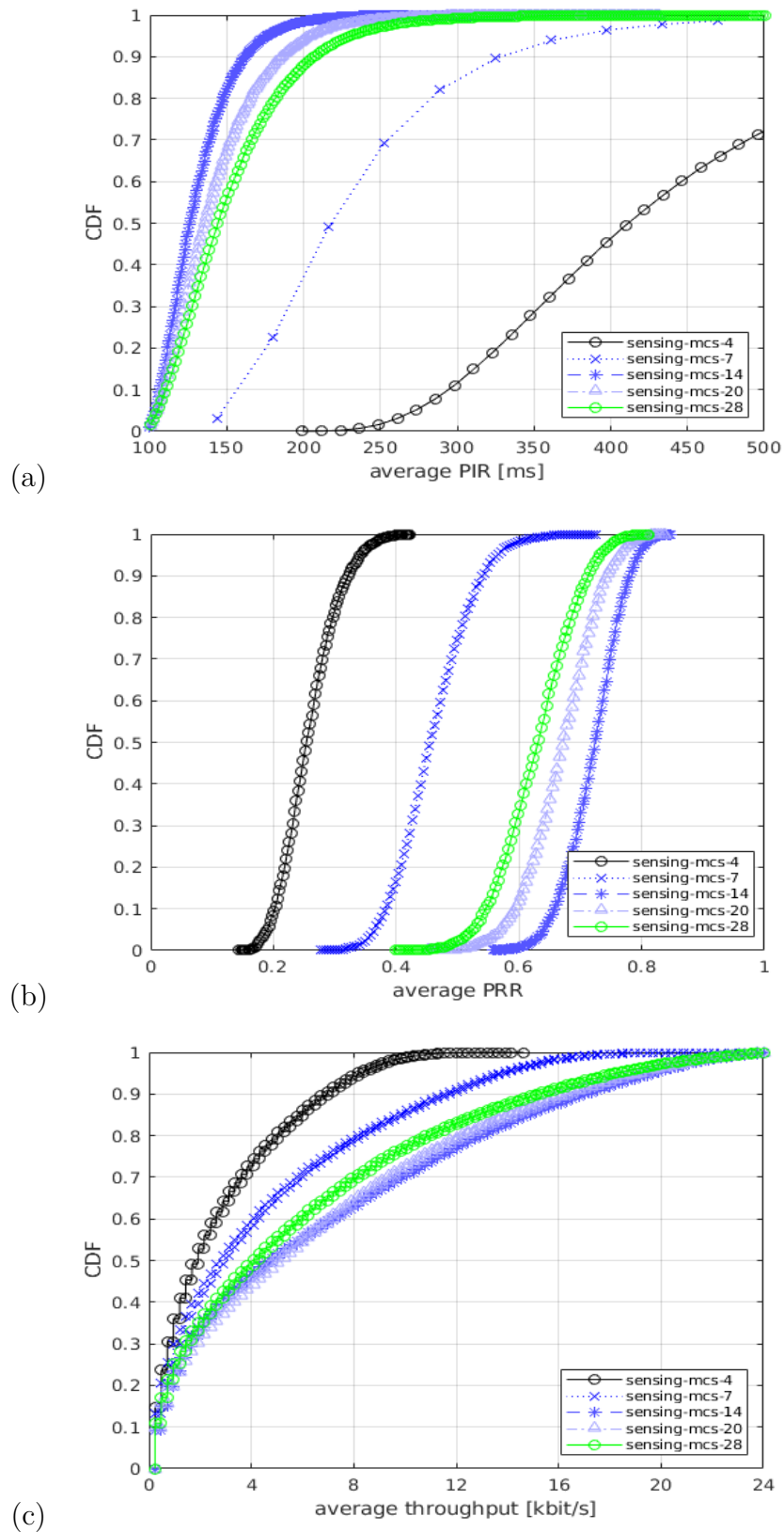
(a)

(b)

(c)

**Figure 7.14:** Impact of NR V2X MCS. (a) PIR (ms), (b) PRR, (c) throughput (kbps).

### 7.4.3.6   Impact of MCS

In the sixth simulation campaign, we evaluate the impact of different MCSs. We have used three MCS indexes adequate for V2X scenarios: MCS4, MCS7, and MCS14. Also, we include two more MCSs, just for comparison purposes: MCS20 and MCS28, even though, these two MCSs may be high for the considered reception range of 200 m. All the considered MCSs are displayed in the legends of the figures as mcs-4, mcs-7, mcs-14, mcs-20, and mcs-28, respectively. Fig. 7.14 shows the statistics of the PIR, PRR, and throughput.

Simulation results confirm that MCS14 is adequate for the considered simulation scenario with 200 m reception range. As expected, if a higher MCS is used (e.g., MCS20 and MCS28), the end-to-end performance in terms of PIR, PRR and throughput is degraded as compared to MCS14 (see Fig. 7.14.(a)-(c)). This is because a higher MCS is not suitable for large distances and leads to incorrect PHY receptions and packet losses that cannot be recovered even with HARQ combining. On the other hand, if an MCS lower than MCS14 is used, we also observe a performance degradation in all the considered KPIs (PIR, PRR, throughput), as shown in Fig. 7.14 for MCS4 and MCS7. In this case, even though lower MCS are more robust to packet losses and large distances, the reason of degradation is that the amount of data that fits in one subchannel gets reduced with a lower MCS, because of the reduced modulation order and the lower effective code rate. In particular, for the considered configuration and traffic pattern, with MCS14, one data packet can fit in one subchannel, as explained in Sec. 7.4.3.1. On the other hand, two subchannels are needed with MCS7 and three subchannels are required for MCS4. As more subchannels are needed with lower MCSs, the frequency diversity gain gets reduced. That is, less UEs can be multiplexed in frequency domain within the same slot, without interference. Accordingly, lower MCSs experience also PIR, PRR and throughput degradation as compared to MCS14.

All in all, an intermediate MCS (i.e., MCS14 in our case) is shown to be beneficial for V2X scenarios with a reception range of 200 m and the broadcast use cases, because it allows exploiting the frequency multiplexing gain and overcoming propagation losses, simultaneously.

## 7.4.4   Simulation results:   Sensing-based vs Random resource selection

In the last campaign, we consider the baseline configuration (i.e., numerology 0, number of PSSCH transmissions = 5, $T_2 = 32$ slots, and $N_{\mathrm{max\_reserve}} = 3$) and we focus on comparing sensing-based and random resource selection procedures for NR V2X. Notice that the sensing-based resource selection is defined by 3GPP Release 16, but Release 17 is also considering the random resource selection for power saving purposes, as previously discussed in Section 7.2.2.1. The two techniques are labeled in the legends of the figures as sensing and random, respectively. Fig. 7.15 shows the CDF of the PIR, PRR, and throughput.
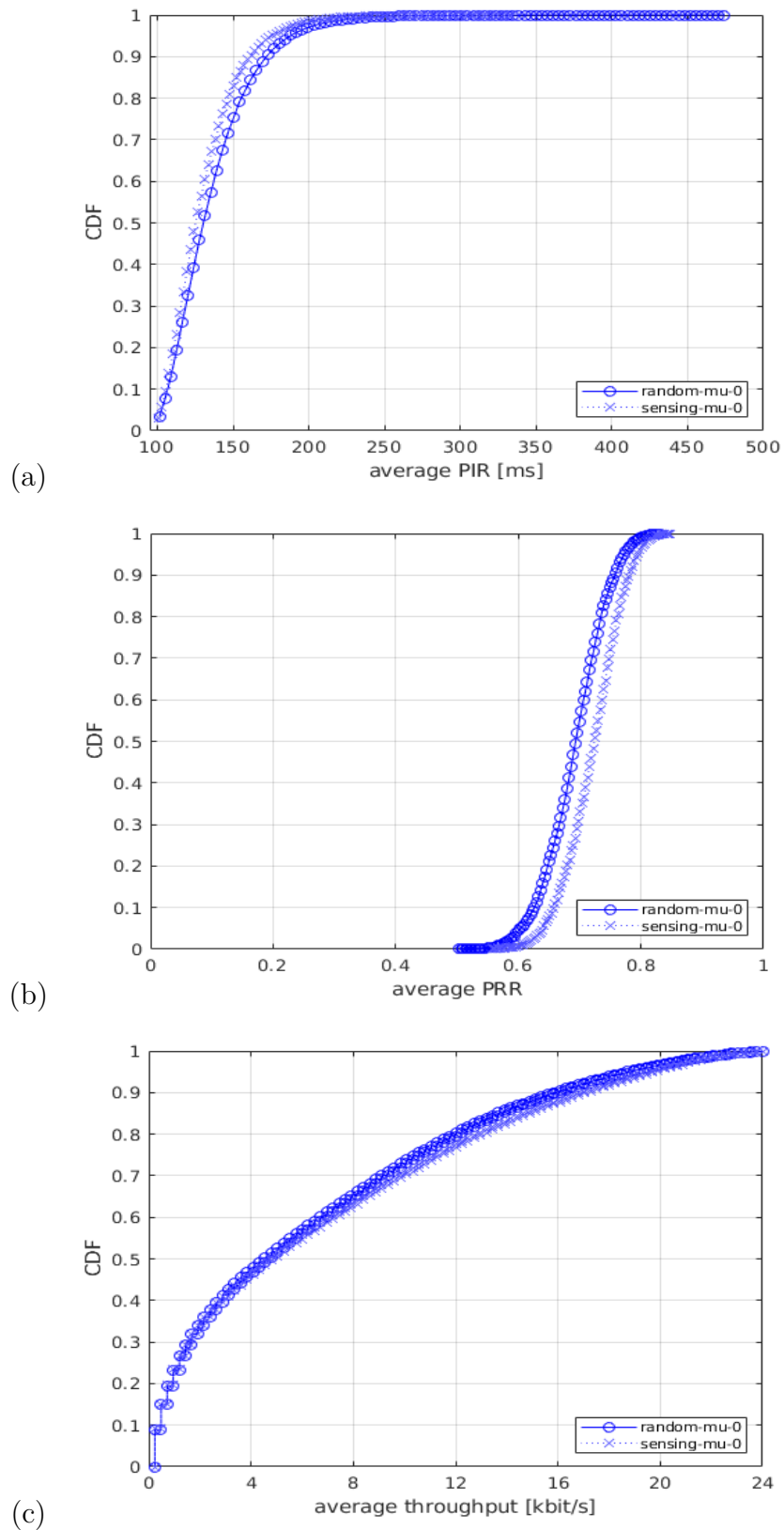
(a)

(b)

(c)

**Figure 7.15:** Impact of NR V2X resource selection procedure: sensing vs random. (a) PIR (ms), (b) PRR, (c) throughput (kbps).

Simulation results confirm that the sensing procedure ends up in a reduced PIR (see Fig. 7.15.(a)), an increased PRR (see Fig. 7.15.(b)), and a larger throughput (see Fig. 7.15.(c)). This is because sensing-based resource selection allows reducing the number of simultaneous PSSCH transmissions and incorrect PSSCH receptions in the reception range, as compared to the random resource selection procedure. As a consequence, due to the effectiveness of the sensing, we observe the PIR, PRR, and throughput improvement of sensing over the random selection procedure, in all the percentiles of the output statistics. So, these results confirm the expectations about the improvements given by sensing and show its performance gains in an end-to-end system-level simulator. The question that remains open though is whether the improvements provided by the sensing procedure are considerable enough to compensate for the increased complexity and power consumption that sensing involves. The answer to this question is however out of the scope of this study and is left for a future work.
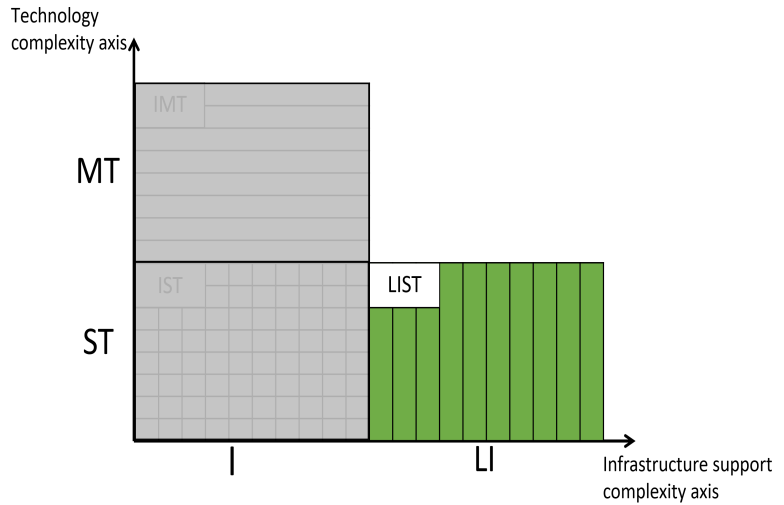
## 7.5 Conclusions

In this chapter, we have presented an open source, full-stack, end-to-end, standard-compliant network simulator for NR V2X, based on an extension of the already available *ns-3* NR 5G-LENA simulator. We have started by reviewing the history of the different radio access technologies developed by 3GPP for sidelink communications and we have provided an exhaustive overview of NR V2X technology, currently under development in 3GPP, with special emphasis on NR V2X Mode 2 for autonomous resource selection, which is the main focus of this work. Successively, we have described with details the design and implementation of the developed simulator, which provides a useful and readable introduction to the module for a prospective and interested user. We have focused our work on broadcast communications for out-of-coverage scenarios, following the specifications of NR V2X Mode 2. For that, we have described the RRC pre-configuration and the NR V2X-compliant procedures at PHY and MAC layers, using UE autonomous resource selection based on sensing and semi-persistent scheduling.

Finally, we have presented a complete set of simulation campaigns, including the impact assessment of key NR V2X parameters, such as the numerology, the resource selection window, the number of retransmissions, the maximum number of resources per reservation, the probability of resource keep, and the MCS, as well as a comparison of sensing and random based resource selection procedures. The seven simulation campaigns that we have conducted have highlighted the following:

1. Only $\mu = 2$ exhibits clear benefit to improve PIR, PRR, and throughput performance metrics.

2. A low number of blind retransmissions is already optimal in scenarios with good propagation conditions, because of the effectiveness of the sensing-based resource selection.

3. A dependency between the resource selection window length and the number of collisions is observed.

4. The maximum number of resources per reservation does not have a noticeable impact on the end-to-end performance.

5. A trade-off is observed between the probability of resource keep and the performance KPIs for NR V2X with periodic traffic.

6. An intermediate MCS is shown to be optimal for V2X scenarios and broadcast use cases.

7. Appreciable but not significant gains are obtained by using sensing-based resource selection, in comparison to the random resource selection strategy, which is considered in 3GPP Release 17.

With these campaigns, we have only touched the tip of the iceberg, and many more studies can be conducted by the research community, considering the proposed open-source platform as a basis for analysis. In this line, in the next chapter, we leverage this simulator to explore the energy-performance trade-off offered by different resource allocation procedures in MODE 2. Based on this, we propose a solution that dynamically balances the mentioned trade-off.

# Chapter 8

# Exploiting Energy-Performance trade-off in 3GPP NR V2X using Fuzzy Logic

The 3GPP's NR V2X technology is based on sidelink communication. It enables a UE to communicate directly with other UE(s) without sending packets first to its base station (gNB in NR). As explained in the previous chapter, based on the resource allocation method, sidelink communication is classified into two categories, i.e., Mode 1 and Mode 2 [92]. In Mode 1, the resource allocation is managed by the base station (i.e., centralized). It thus applies to scenarios where the various UEs are inside the coverage of the base station (i.e., in-coverage scenarios). On the other hand, Mode 2 is a distributed scheduling approach in which the UEs themselves carry out the resource allocation, with no need to be in the coverage area, i.e., it supports out-of-coverage communications.

Resource reservation for NR V2X Mode 2 under periodic traffic uses a long-term sensing-based algorithm, which exploits the periodicity and fixed-size assumption of basic safety messages. Thanks to the sensing, a more accurate and collision-preventing resource selection mechanism can be achieved. However, it comes at the cost of increased energy consumption due to continuous sensing. To reduce the energy consumption at V2X devices, random and partial sensing mechanisms are standardized by 3GPP for NR V2X Mode 2 but the performance is hampered due to collisions resulting from less accurate resource selection [173, 180]. Therefore, it creates an energy-performance trade-off that leads us to the problem: how much sensing should be performed by each vehicle in the network so that energy-performance is balanced depending on each vehicle's environment

(e.g., the density of surrounding transmitting vehicles) and situation (e.g., battery level or energy consumption)?

In this chapter, we present a fuzzy logic-based partial sensing duty cycle mechanism for NR V2X Mode 2 that exploits the energy-performance trade-off in vehicular scenarios. In particular, we use fuzzy logic to harness its capabilities of handling uncertainty, vagueness, and dissimilar inputs to make decisions. Uncertainty and vagueness are the two characteristics of typical heterogeneous scenarios, such as vehicular scenarios, where exact traffic situations cannot be predicted with certainty. In such conditions, we might have to make a decision about the amount of sensing to be performed based on dissimilar inputs, which are not directly comparable but can be processed in a homogeneous manner using fuzzy sets. This study considers one set of such inputs and explains them in detail in Sec 8.3. Furthermore, we take advantage of fuzzy logic interpretability and explainability characteristics by proposing different fuzzy rules based on different objectives, which are particularly valuable in integrating user or operator preference about sensing and performance in complex decision-making scenarios. In this line, to validate the models, we have interfaced Matlab (for the fuzzy logic system) with an extension of the open-source, end-to-end, *ns-3* 5G-LENA simulator [184], developed to support NR V2X capabilities [192]. Specifically, we used the Semi-Online-Evaluation approach, explained in Chapter 1 to evaluate the performance of the proposed models.

The chapter is structured as follows. In Section. 8.1 we discuss the related work. Sec. 8.2 reviews NR V2X Mode 2 resource allocation, introduces the energy-performance trade-off and derives the energy consumption for different resource selection mechanisms, including random, sensing and partial sensing. Sec. 8.3 describes the proposed fuzzy logic-based partial sensing duty cycle mechanism. Sec.8.4 presents the simulation results using *ns-3*. Finally, Sec. 8.5 concludes the chapter.

## 8.1 Related work

The authors in [193] proposed a frequency-selective partial sensing mechanism that uses a reduced number of subchannels for sensing to decrease the energy consumption of V2X UEs. However, the question about how to select those reduced number of subchannels still needs to be answered. Authors in [194] also acknowledge a trade-off between energy consumption and the performance of V2X UEs. To address this, the authors formulate an energy efficiency maximization problem analytically, considering the latency and reliability constraints. However, it is suboptimal if the number of UEs is high. Alternatively, they proposed a heuristic algorithm to select the sensing or random resource selection procedure and the corresponding dedicated pool of resources to serve periodic and aperiodic traffic. Specifically, it is assumed that a central control or roadside unit provides traffic-related information to the UE. If the density of periodic traffic flow is greater than 70%, UE uses the sensing-based procedure to select resources from a dedicated resource pool for periodic traffic. Otherwise, the UE performs random resource selection on the corresponding pool of resources. The result shows that splitting the resource pools to serve the dedicated traffic and using an appropriate resource selection procedure for the UEs with aperiodic traffic improves the overall energy consumption by maintaining the same PRR achieved by employing sensing. However, the performance

gains achieved by the algorithm are dependent on the availability of enough bandwidth to have a dedicated resource pool for each traffic type. Therefore, it could be challenging to maintain the same performance when such a split is not possible, e.g., under low bandwidth situations. Moreover, as per the algorithm, the UE either uses sensing or does not sense, i.e., random resource selection. This indicates that it might not be able to identify medium-density scenarios in which employing partial sensing could benefit energy efficiency, given the performance constraints. To tackle these dynamic situations, the use of machine learning in V2X communication has shown promising results. For example, the authors in [195] proposed a fuzzy logic-based resource allocation algorithm for LTE V2X Mode 3, which is similar to the NR V2X Mode 1. The algorithm maximizes the resource reusability and satisfies the V2X service requirements, i.e., the latency, by self-adopting to the changes such as interference level in the network. In this thesis, different from the solutions presented in [193, 194], we propose a solution to regulate the sensing period dynamically instead of operating at extremes, i.e., sensing or no-sensing, and without modifying the standard sensing procedure. Furthermore, compared to the fuzzy logic-based solution in [195], which is designed to maximize the spectrum efficiency in MODE 1, we focus on achieving a balanced trade-off between the energy consumption and the performance of a V2X UE in MODE 2.

## 8.2 NR V2X Mode 2 and the Energy-Performance Trade-off

Towards the end of release 17, besides sensing and random resource selection, the 3GPP has also standardized a partial sensing mechanism [173]. Both mechanisms, random and partial sensing, are considered as power saving mechanisms, especially relevant for the use cases of public safety, pedestrian UEs in V2X scenarios, and electric vehicles where UEs have battery limited capacity and must operate efficiently [194]. In the partial sensing, the sensing information considers decoding only of a part of the entire data. Thus, when partial sensing is used, the power consumption is reduced as much as the decoding time of information data is reduced.

Accordingly, sensing-based resource reservation is the mechanism that can achieve a more accurate resource selection (and so higher throughput and lower probability of collision) at the cost of a higher power consumption. The partial sensing reduces the power consumption but also gets a reduced performance, because of the less accurate resource selection. Finally, the random resource selection is the mechanism entailing lower power consumption, but, at the same time, higher performance degradation. Therefore, there is a clear energy-performance trade-off when using different resource allocation methods for NR V2X. Let us note that the trade-off is more pronounced in scenarios where there are nodes that act only as transmitters, which need to do additional decoding and processing in case of using sensing or partial sensing-based resource selections.

In what follows, we derive the energy consumption of NR V2X under different resource selection methods. In general, the time to transmit/receive PSCCH and PSSCH transmissions are proportional to the number of PSCCH/PSSCH transmissions/receptions and their time duration. In this work, based on the NR UE power model presented in [196],

we model the power and the energy consumption of a NR V2X UE when it acts as only a transmitter, using random, sensing, and partial sensing-based resource selection procedures. We also note that, all the powers and the energies mentioned in the following are in the units of milliwatt (mW) and Joule (J), respectively.

**Random resource selection:** When using random based resource selection, a UE will not receive any transmission from other surrounding UE(s) but only transmits PSCCH (i.e., 1st stage SCI) and PSSCH (i.e., 2nd stage SCI + Data). Since we are only transmitting, no power/energy would be consumed for reception and a UE would consume its power only when performing transmission(s). In this case, the transmit energy consumption results:

$$E^{\text{rnd}} = P_{\text{tx,pscch}} \times T_{\text{tx,pscch}} + P_{\text{tx,pssch}} \times T_{\text{tx,pssch}} \tag{8.1}$$

where $T_{\text{tx,pscch}}$ and $T_{\text{tx,pssch}}$ are the total time spent to transmit PSCCH and PSSCH, respectively, and $P_{\text{tx,pscch}}$ and $P_{\text{tx,pssch}}$ are the power consumed to transmit PSCCH and PSSCH, respectively. As per [197], for UE total available transmit power of 23 dBm, its PSCCH/PSSCH transmit power consumption is: $P_{\text{tx,pscch}} = P_{\text{tx,pssch}} = 700$ mW.

**Sensing resource selection:** In this case, in addition to transmitting PSCCH and PSSCH, a UE would receive PSCCH (i.e., 1st stage SCI) messages and buffer them as sensing information to perform the resource selection. Therefore, the total energy consumption is:

$$E^{\text{sens}} = E^{\text{rnd}} + P_{\text{rx,pscch}} \times T_{\text{rx,pscch}} \tag{8.2}$$

where $T_{\text{rx,pscch}}$ is the total time spent to receive PSCCH and $P_{\text{rx,pscch}}$ is the power consumed to receive PSCCH. Thus, the term $P_{\text{rx,pscch}} \times T_{\text{rx,pscch}}$ reflects the energy spent to decode the SCI messages. Table 1 in [196] details a power consumption of $P_{\text{rx,pscch}} = 100$ mW to monitor and process PSCCH for the case of 4 receive RF chains and 100 MHz bandwidth. These power levels must be scaled appropriately, depending on the number of used RF chains and bandwidth at the UEs. For example, for 1 RF chain and 40 MHz bandwidth, by using the scaling factors provided in [196], the PSCCH receive power results: $P_{\text{rx,pscch}} = 100 \times 1.4 \times 0.325 = 45.5$ mW.

The UEs that perform sensing-based resource selection can achieve a better selection and a better performance as compared to random resource selection, at the cost of an additional energy consumption, as shown in (8.2). For that reason, to address the trade-off, partial sensing mechanisms are being envisioned by 3GPP.

**Partial sensing resource selection:** In partial sensing, the UEs switch between sensing and random resource selections based on a partial sensing window. In this case, the total energy consumption can be modelled as:

$$E^{\text{part}} = \alpha E^{\text{sens}} + (1 - \alpha) E^{\text{rnd}} \tag{8.3}$$

where $\alpha \in [0, 1]$ is the relation of the partial sensing window over the sensing window used for typical sensing. Note that the extreme cases result into the sensing-based resource selection (for $\alpha = 1$) and the random resource selection (for $\alpha = 0$).
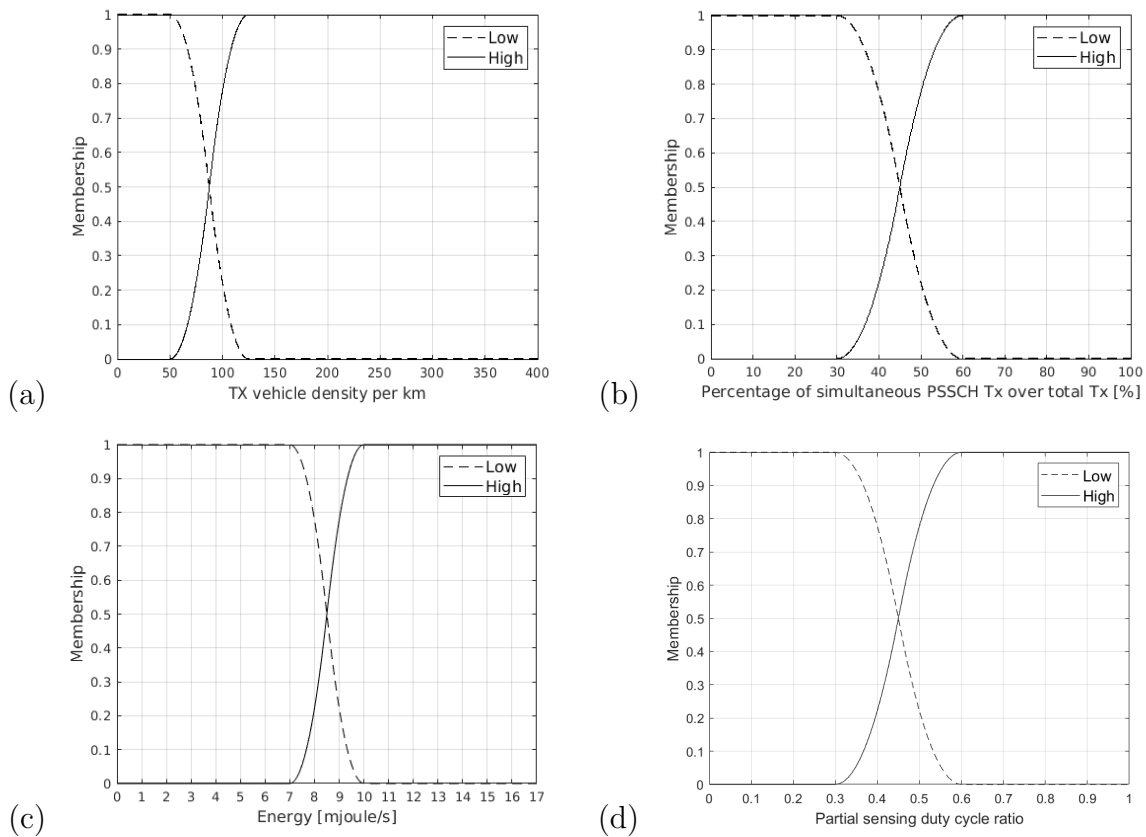
**Figure 8.1:** Input and output fuzzy membership functions. (a) The fuzzy set for vehicle density per km: "Low" and "High": "Low" and "High", (b) The fuzzy set for the percentage of simultaneous PSSCH Tx over total TX: "Low" and "High", (c) The fuzzy set for the energy consumption, (d) The fuzzy set for partial sensing duty cycle ratio during 1 second period: "Low" and "High".

## 8.3 Partial Sensing Mechanism using FL

In our implementation of the partial sensing mechanism, UEs are configured with a sensing duty cycle. This duty cycle has two parameters that define an on-off behaviour for sensing, the duty-cycle period and the duty cycle ratio. Specifically, the period defines how often the on-off pattern repeats (defined in time, e.g., every 1 sec), and the ratio is the fraction of the period a UE performs the sensing. To exploit the energy performance trade-off under specific network objectives, e.g., improving the performance or reducing the energy consumption, the duty-cycle ratio could be configured statically between 0%, i.e., no-sensing, and 100%, i.e., full-sensing. However, the dynamicity of the NR V2X scenarios can make it challenging to select an appropriate duty cycle. Therefore, it should be configured based on some key metrics that may change dynamically and reflect the change in a network, specific locations, and its performance. This section proposes a distributed partial sensing mechanism using Fuzzy Logic (FL) that resides inside a UE and can adaptively configure the duty-cycle ratio over a fixed period.

In this line, three fuzzy input variables are identified: the energy consumption of a UE, the percentage of simultaneous Physical Sidelink Shared Channel (PSSCH) transmissions over total transmissions by all the UEs in a coverage area, and the density of the

transmitting (TX) vehicles per km inside a coverage area. Following the NR V2X energy model derived in Sec. 8.2, the energy consumption to transmit and receive sidelink transmissions can be computed locally by a UE. For the other two fuzzy input variables, we assume that those are broadcast by a RSU since it has a global view of the local network that can improve the accuracy of these variables. Alternatively, a UE can also compute them locally at the cost of accuracy. For example, the vehicle density in the coverage area can be discovered by a UE performing full-sensing during which it decodes the first stage SCI and can count unique layer-2 source identifiers. However, its accuracy may be affected when some of the first stage SCI messages are not decoded, either due to the half-duplex nature of the NR V2X UE or under partial sensing where a reduced number of slots are decoded. Nevertheless, the objective of this work is to manifest an idea of dynamically adjusting the duty cycle ratio using these parameters. Therefore, obtaining them locally or globally is implementation-specific and is out of the scope of this study.

Two input fuzzy sets, "High" and "Low" are defined for each of the fuzzy input variables. As shown in Fig. 8.1, their membership functions are complimentary in nature, i.e., if $f(x)$ and $g(x)$ denote the membership function of "High" and "Low", respectively, then, $f(x) + g(x) = 1 \ \forall \ x$. We note that, given the novelty of the problem discussed in this chapter, as per the knowledge of the authors, currently there is no experimental or simulation-based data available that could be used to decide the thresholds for "High" or "Low" for all these fuzzy variables. Therefore, these thresholds were selected based on our expert knowledge by simulating several vehicle densities under sensing and random resource selection procedures, which are the two extreme cases to understand the energy and performance trade-off. Fig. 8.1(a) shows the membership function of fuzzy input variable "the density of the transmitting vehicles per km inside a lane". Here, the transmitter density per km is computed as [198]:

$$D = \frac{N_{\text{tx}} \times 1000}{L} \tag{8.4}$$

where $D$ is the density of the transmitting UEs per km, $N_{\text{tx}}$ is the number of transmitting UEs occupying a length of the lane, and $L$ is the length of the lane occupied by the vehicles (in m). It is assumed that if $D$ is above or equal to 125, its membership in "High" is unity, and if it is lower than or equal to 50, its membership in "Low" is unity. The threshold for the "Low" fuzzy set for this membership function is based on the assumption that using any of the resource selection procedures the percentage of simultaneous PSSCH transmissions should be lower or equal to 30 %, which in our experimental campaigns resulted from 50 transmitting vehicles per km. On the other hand, if this percentage is above or equal to 60 %, the vehicle density is considered as high, i.e., 125 transmitting vehicles per km, in this case. Based on this assumption, Fig. 8.1(b) shows the membership function for the second fuzzy input variable "the percentage of simultaneous PSSCH transmissions over total transmissions by all the UEs in a coverage area", its membership function in "High" is unity if it's above or equal to 60 %. On the other hand, if it is less than or equal to 30 % its membership in "Low" is unity. Finally, Fig. 8.1(c) shows the membership function of "the energy consumption of a UE", where it is assumed that if the energy consumption is above or equal to 10 mJoules/s, its membership in "High" is unity, and if it is less than or equal to 7 mJoules/s, its membership in "Low" is unity. In our case, using random resource selection, a UE, irrespective of the transmitter vehicle

**Table 8.1:** Proposed Fuzzy Rules

| | Rule prerequisite | Objective 1 (energy + performance) | Objective 2 (energy) | Objective 3 (performance) |
|---|---|---|---|---|
| 1 | **IF** the energy consumed is *High* and the number of simultaneous transmissions are *High* and the density of the vehicles is *High*, | **THEN** duty cycle is *High* | **THEN** duty cycle is *Low* | **THEN** duty cycle is *High* |
| 2 | **IF** the energy consumed is *Low* and the number of simultaneous transmissions are *High* and the density of the vehicles is *High*, | **THEN** duty cycle is *High* | **THEN** duty cycle is *High* | **THEN** duty cycle is *High* |
| 3 | **IF** the energy is *High* and the number of simultaneous transmissions are *Low* and the density of the vehicles is *High*, | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* |
| 4 | **IF** the energy consumed is *Low* and the number of simultaneous transmissions are *Low* and the density of the vehicles is *High*, | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* | **THEN** duty cycle is *High* |
| 5 | **IF** the energy consumed is *High* and the number of simultaneous transmissions are *High* and the density of the vehicles is *Low*, | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* | **THEN** duty cycle is *High* |
| 6 | **IF** the energy consumed is *Low* and the number of simultaneous transmissions are *High* and the density of the vehicles is *Low*, | **THEN** duty cycle is *High* | **THEN** duty cycle is *High* | **THEN** duty cycle is *High* |
| 7 | **IF** the energy consumed is *High* and the number of simultaneous transmissions are *Low* and the density of the vehicles is *Low*, | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* | **THEN** duty cycle is *Low* |
| 8 | **IF** the energy consumed is *Low* and the number of simultaneous transmissions are *Low* and the density of the vehicles is *Low*, | **THEN** duty cycle is *Low* | **THEN** duty cycle is *High* | **THEN** duty cycle is *High* |

density, at maximum consumes 7 mJoules/s; thus, any value of the energy consumption equal to or below this value is considered as low. Now, to select the threshold for "High", we observed that with the lowest transmitter vehicle density, i.e., 50, a UE using sensing-based resource selection may consume, at maximum, 10 mJoules/s. Therefore, the consumption higher than this value is considered as high. At last, along with the above three fuzzy input variables, one fuzzy output variable is identified: the partial sensing duty cycle ratio. Fig. 8.1(d) shows the membership functions of the two fuzzy sets, "High" and "Low", assigned to this variable. For the "High" fuzzy set, if the duty
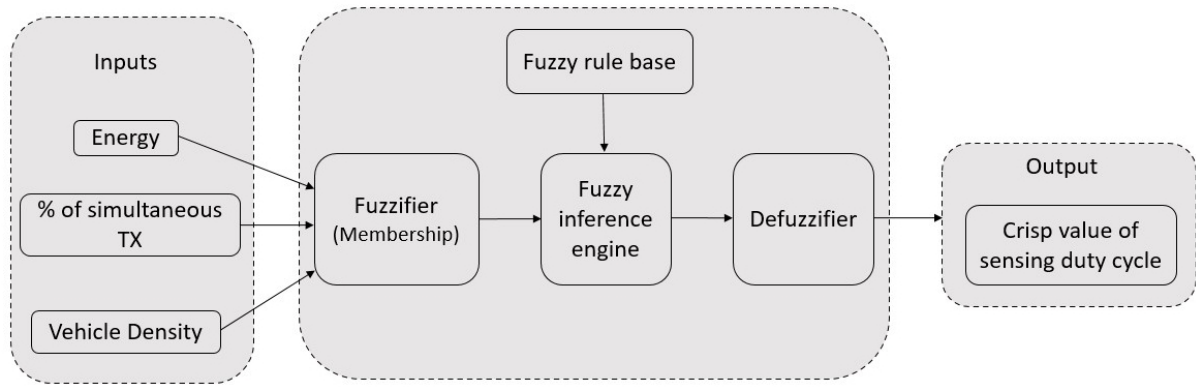
**Figure 8.2:** Block diagram of fuzzy inference system for NR V2X partial sensing

cycle ratio is above or equal to 0.6, its membership is unity. On the other hand, if it is lower or equal to 0.3, its membership in "Low" is unity.

Based on the values of the fuzzy input variables, we define a set of rules to achieve the desired behaviour. In particular, we designed three different set of rules, depending on the objective. The first objective focuses on improving both the energy and performance (PRR or throughput) simultaneously. The second objective is designed to prioritize energy, i.e., reducing energy consumption is prioritized over improving the PRR or throughput. The third objective prioritizes performance (PRR or throughput) over energy. Table 8.1 shows the proposed fuzzy rules, depending on the fuzzy input variables' condition, for the three different objectives. When energy is prioritized (Objective 2), the duty cycle is in more cases "Low" (i.e., tending to no sensing) to reduce the energy consumption, except for the cases in which the energy consumed is already "Low" and sensing can help improving the performance (i.e., rule prerequisites 2, 6 and 8). On the other hand, when performance is prioritized (Objective 3), the duty cycle is "High" in major part of the conditions or rule prerequisites. There are two exceptions (prerequisites 3 and 7), in which the number of collisions is "Low" and energy consumed is "High", and thus we can try to balance energy because the number of collisions in the system is already controlled, and sensing may not be so much needed. Finally, Objective 1 follows intermediate rule conditions between Objectives 2 and 3, trying to balance between energy and performance under each condition.

The block diagram of the proposed Fuzzy Inference System (FIS) for adaptive selection of the duty-cycle ratio is shown in Fig 8.2. The proposed FIS is built using the fuzzy logic toolbox of MATLAB that uses the Mamdani inference system [199]. At first, for each input to the model, the "Fuzzifier" block finds the degree of membership (also known as membership value) with which a particular input belongs to the fuzzy set, i.e., "High" and "Low" using their respective membership functions in Fig 8.1 (a) to (c). Then, the "Fuzzy inference engine" generates a single output fuzzy set using those input membership values and the "Fuzzy rule base (i.e., 8 rules in Table 8.1)". To do so, the engine uses the well-known "And" based fuzzy implication method and combines the output of each rule using the "Sum" aggregation procedure. After that, the final aggregated fuzzy set is defuzzified using the "Defuzzifier" block. Specifically, it uses the MATLAB built-in "centroid" method to generate a single value of the sensing duty cycle ratio [200].

# 8.4 Simulation Results

For the evaluation, ns-3 5G-LENA system-level simulator extension for NR V2X is used. The implementation details of the ns-3 NR V2X module are summarized in [62].

## 8.4.1 Scenario

We consider a typical V2X highway scenario, as defined in 3GPP TR 37.885 [188]. The deployment is composed of three lanes with an inter-lane and inter-vehicle distance of 4 m, and we deploy 50 vehicles per lane. In the considered scenario, each vehicular UE is a potential receiver, but only a subset of them do transmit data. The reason of selecting the above scenario parameters is because the evaluation methodology proposed by the 3GPP standard, mainly considers range-based KPIs. For example, the Packet Reception Ratio (PRR) is defined as the ratio of vehicles that successfully receive a given packet, considering only those UEs that are located within a specific range, called, the "awareness range". Therefore, in this work we limit the awareness range to 200 m (hence 4 m inter-vehicle distance), and study the energy-performance trade-off by varying the number of transmitting vehicles within this range. In particular, we consider two different vehicle TX densities; either 10 or 25 TX vehicles per lane.

The UE dropping is implemented according to [188] Option A, in which all vehicles (100 %) are of Type 2 (i.e., passenger vehicle with an antenna height of 1.6 m), clustered dropping is not used, and the vehicle speed is set to 140 km/h in all the lanes. We focus on an out-of-coverage scenario, so that gNBs are disabled in the evaluation [188].

As key performance indicators, we consider the PRR and the energy consumption. For the energy consumption, the model presented in Sec. 8.2 is used. The rest of deployment and configuration parameters are detailed in Table 8.2.

## 8.4.2 Results

Fig. 8.3 and Fig. 8.4 show the CDF statistics of the PRR and energy consumption, respectively, for the case of (a) 10 TX vehicles per lane and (b) 25 TX vehicles per lane. In each figure, we compare the three fuzzy logic objectives proposed in Sec. 8.3 for partial sensing duty cycling (obj-1: energy+performance, obj-2: energy, and obj-3: performance). Also, we consider three baseline strategies: sensing-based resource selection (fixedDc-1), random resource selection (fixedDc-0), and a partial sensing using a fixed duty cycle of 50% (fixedDc-0.5).

We can observe that the proposed fuzzy logic approach allows to trade-off in terms of energy and PRR among the baselines sensing and random resource selection strategies. Among the various proposed objectives and rules for the fuzzy inference system, we can see that, in both scenarios, obj-2 prioritizes the energy, while obj-3 prioritizes PRR performance at the cost of a higher energy consumption. For the low TX density scenario (see Fig. 8.3.(a) and Fig. 8.4.(a)), obj-1 is similar to obj-2, since energy is primed because of the low collision probability. Instead, for the higher TX density scenario (see

**Table 8.2:** Simulation Parameters

| Parameter | Value |
|---|---|
| Number of vehicles | 150 (50 per lane) |
| Number of TX vehicles | 30 or 75 (10 or 25 TX per lane) |
| Propagation scenario | 3GPP V2X Highway |
| UE antenna height | 1.6 m |
| UE antenna | 1x2 antenna array |
| UE transmit power | 23 dBm |
| UE speed | 140 km/h |
| UE noise figure | 5 dB |
| Carrier frequency | 5.9 GHz |
| Bandwidth | 10 MHz |
| Numerology | 0 (15 kHz subcarrier spacing) |
| RB overhead | 0.04 |
| Duplexing mode | TDD |
| TDD pattern | [D D D F U U U U U U] |
| Sidelink bitmap | [1 1 1 1 1 1 0 0 0 1 1 1] |
| PSSCH and PSCCH multiplexing | Time based<br>PSCCH symbols: 1<br>PSSCH symbols: 12 |
| MCS PSSCH | MCS14 |
| MCS PSCCH | MCS0 |
| Subchannel size | 10 RBs |
| Sensing window | 100 ms |
| $T_2$ | 33 slots |
| $T_1$ | 2 slots |
| $T_{proc,0}$ | 2 slots |
| Probability of resource keep | 0 |
| Resource reservation window | 100 ms |
| RLC | Unacknowledged Mode |
| Transport protocol | UDP |
| Traffic | Periodic packet transmissions, with a packet size of 300 bytes, transmitted every 100 ms, leading to a data rate of 24 kbits/s. |
| Simulation duration | 10 s |

Fig. 8.3.(b) and Fig. 8.4.(b)), obj-1 nearly matches obj-3 KPIs, because PRR performance is prioritized owing to the higher interference conditions.

All in all, the proposed fuzzy logic-based partial sensing duty cycle allows handling the energy-performance trade-off in an independent manner at each vehicle and, depending on the pre-defined fuzzy logic rules, different KPIs can be prioritized while balancing all of them simultaneously. The benefit as compared to a fixed partial sensing duty cycle is the adaptability to each vehicles' observations. This has shown the applicability and suitability of AI/ML to choose a balanced resource allocation strategy in complex and dense vehicular scenarios.
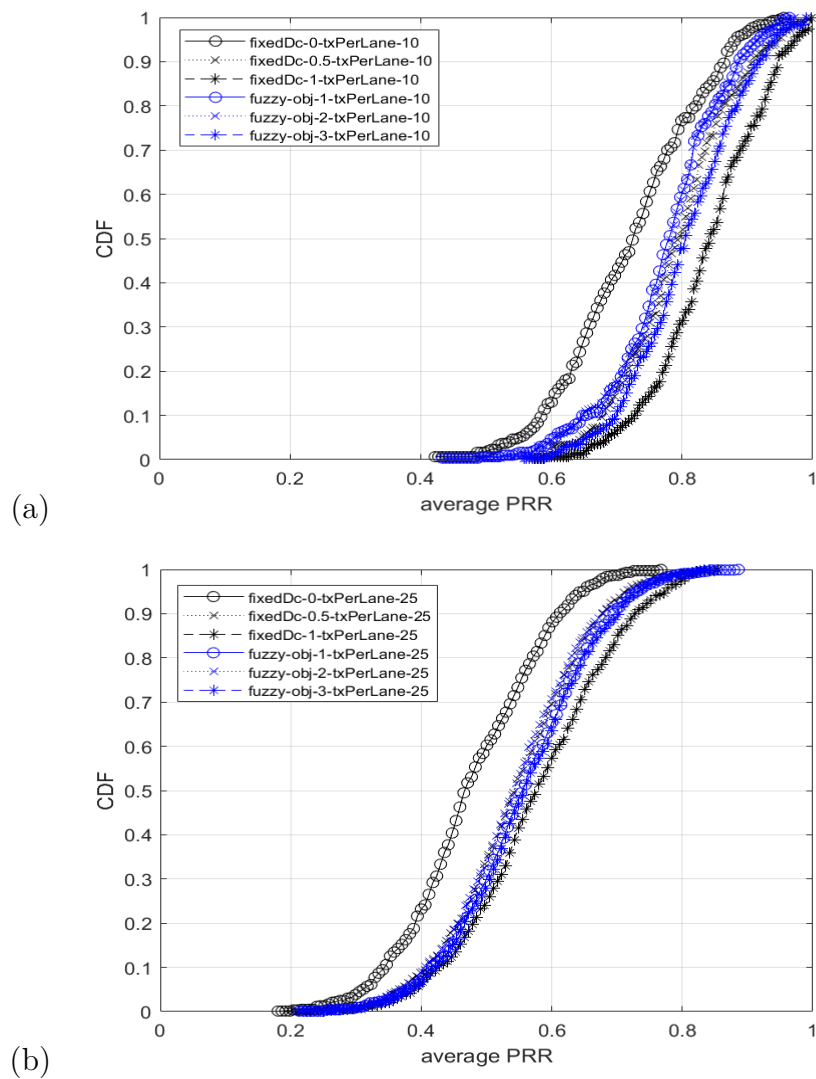
**Figure 8.3:** Fixed partial sensing duty cycle vs Fuzzy logic-based partial sensing duty cycle. (a) PRR with 10 TX per lane, (b) PRR with 25 TX per lane.
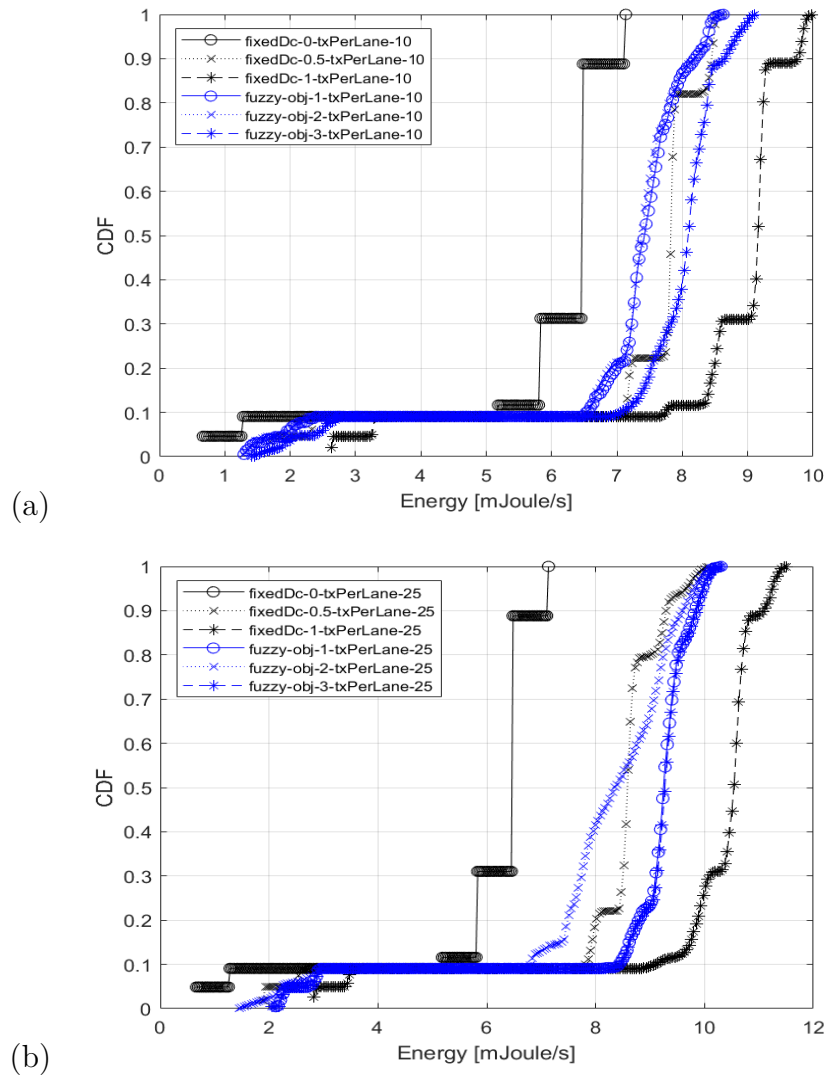
(a)



(b)

**Figure 8.4:** Fixed partial sensing duty cycle vs Fuzzy logic-based partial sensing duty cycle. (a) energy consumption with 10 TX per lane, (b) energy consumption with 25 TX per lane.

## 8.5 Conclusions

In this chapter, we address the complexity that encompasses the limited infrastructure-based single-technology axes. In particular, we investigate the energy-performance trade-off offered by the sensing and random resource selection methods of the 3GPP NR V2X Mode 2. Specifically, to address this energy-performance trade-off, we propose a fuzzy logic-based partial sensing resource selection procedure. This solution answers our fifth research question (RQ5): "How to achieve a balance between the energy consumption and the performance of a NR V2X UE in limited infrastructure-based scenarios using AI/ML?". The proposed scheme dynamically adjusts the partial sensing duty cycle ratio over a fixed sensing window by using a fuzzy inference system. The fuzzy inference system takes as inputs the energy consumption at the UE, the number of simultaneous transmissions, and the vehicle's density, and considers various target objectives (or fuzzy rules). The results show that the proposed scheme achieves a good balance between the energy consumption and the performance of a UE under three different objectives, i.e., simultaneously improving both the energy and PRR or prioritizing the energy consumption or improving the PRR. Through a simple but effective fuzzy logic-based model, our work has motivated the use of an AI/ML-based resource selection method that not only focuses on improving the performance but also the energy consumption, which has always been one of the key objectives in mobile networks.

Specifically, to address this energy-performance trade-off, we propose a fuzzy logic-based partial sensing resource selection procedure. Our scheme dynamically adjusts the partial sensing duty cycle ratio over a fixed sensing window by using a fuzzy inference system. The fuzzy inference system takes as inputs the energy consumption at the UE, the number of simultaneous transmissions and the vehicle's density, and considers various target objectives (or fuzzy rules). The results show that the proposed scheme achieves a good balance between the energy consumption and the performance of a UE under three different objectives, i.e., simultaneously improving both the energy and PRR or prioritizing the energy consumption or improving the PRR. Through a simple but effective fuzzy logic-based model, our work has motivated the use of AI/ML-based resource selection method that not only focuses on improving the performance but also the energy consumption, which has always been one of the key objectives in mobile networks.

# Chapter 9

# Conclusions and future work

New generation networks, such as 5G and 6G, are envisioned to improve network performance and user experience by supporting innovative technologies and applications. However, as technology evolves and demands increase, the network operation and management become more complex, especially the RAN being the most challenging to deploy, manage and operate. In this respect, AI/ML-based techniques have attracted much attention from academia and industry to automate RAN management. In literature, various solutions have been proposed in the context of SON and AI/ML. However, as RAN evolves to meet the ever-growing needs of future mobile users, various new challenges arise that must be addressed to prevent potential issues related to network performance and user experience. In this line, this thesis has deeply investigated some of the most relevant challenges and provided novel solutions to prove the applicability of AI/ML in tackling those issues. This chapter first presents the conclusion of this thesis while highlighting its main contributions along the way. Finally, the chapter ends by giving pointers for potential research directions for future work.

## 9.1   Concluding Remarks

In Chapter 1, we provided the motivation, the problem statement, and the research approach we adopted to conduct the research in this thesis. Then, in Chapter 2, for the reader´s convenience, we provided the fundamental knowledge of all the relevant functionalities and mobile technologies we discussed throughout the thesis. Based on the

motivation and problem statement presented in Chapter 1, in there, we formulated the following high-level (HQ) question:

*HQ: How can AI/ML be used to automate the increasing RAN management complexity along two axes: 1. infrastructure- and limited infrastructure-based RAN scenarios and 2. single- and multi-access technology RAN scenarios?*

The two complexity axes mentioned above reflect the level of heterogeneity in which RAN is evolving. The first axis captured the complexities that arise due to the infrastructure of a network, i.e., the RAN operating with or without BSs. We labeled them infrastructure- and limited (because of the presence of a roadside unit) infrastructure-based RAN. Whereas the second axis covered the challenges due to the number of technologies that need to coexist, i.e., single-access technology RAN using only LTE or NR technology, and the multi-access technology RAN that also contains other technologies, such as WiFi. The main contributions that helped us to answer the above high-level question based on these axes are distributed from chapters 3 to 8. In the following, we present the most important conclusions and contributions that resulted from each of these chapters.

In Chapter 3, we focused on the complexity of the HO management use case under the infrastructure-based single-technology scenarios. We targeted this use case because seamless mobility is one of the most important features that all generations of mobile networks have promised to deliver. In this chapter, using a simplified simulation scenario comprised of three BS and three UEs, we showed that the QoE resulting from a HO becomes crucial when there exists an obstacle partially blocking the coverage of the target cell. Therefore, relying on the simple target cell selection, such as the quality of the received signal of the serving and target cell, becomes too shortsighted and can affect users' QoE in the long run.

To tackle the above issue, as the main contribution of this chapter, we proposed an AI/ML solution based on FFNN using supervised learning. This solution can select the next target cell that is expected to provide an improved QoE by leveraging the knowledge acquired from past HO decisions. The results extracted from an Offline-Evaluation showed that the proposed scheme improved 1) the number of completed downloads and 2) the users´ QoE based on the time to download a file, outperforming the benchmark HO scheme. Based on the results, we can confidently conclude about the effectiveness and useability of AI/ML solutions such as FFNN to handle the complexity of HO management use case in unpredictable radio conditions such as the one presented in this thesis. Additionally, in this chapter, we presented the open-source *ns-3* models to simulate obstacles and the deterministic HO algorithm that can be used to create a database for supervised learning in such studies.

In Chapter 4, we continued with the infrastructure-based single-technology scenarios. We addressed two RAN use cases, i.e., the HO management and initial the MCS selection targeting a realistic and complex multicell simulation scenario. Here, we chose the initial MCS selection use case for two main reasons. First, the MCS selection directly impacts spectrum resource utilization and perceived quality of service, which become critical with the growing number of devices that the future mobile networks have to serve. Second, to advocate using MTL approaches to handle multiple RAN use cases that can function at the same or different layers of the mobile protocol stack. We argued that RAN is

more rapidly evolving, which indicates the emergence of a vast amount of use cases that the RAN might have to handle. Therefore, having AI/ML models for each use case will exhaust the network resources, e.g., the computation cost to train these models, which increases linearly with the number of use cases.

In this respect, the main contributions of Chapter 4 are:

- we proposed a shared database of measurements extracted from *ns-3* standard-compliant BS and UE protocol stacks. This database, which we believe can be used to target multiple use cases, was used to train and evaluate the single-task and multitask AI/ML models based on supervised learning for the two targeted use cases.

- we presented two single-task learning models based on LSTM for HO and initial MCS use cases. The offline performance evaluation of the HO solution showed that the proposed model outperforms the benchmark HO algorithm by increasing the number of UEs able to download a file and decreasing the time to download the file. Similarly, the LSTM model for the MCS use case learned the dynamics of the radio environment using the proposed database. Then, based on this knowledge, it assigned an appropriate initial MCS that is not limited to MCS 0. This is an improvement over the benchmark approach that always assigns MCS 0 due to the lack of channel status report from a newly connected UE. The results obtained showed a significant increase in the achieved initial throughput for UEs that established a connection with a new BS.

- lastly, we proposed the MTL models based on LSTM AE and an MLP. The thesis demonstrated the use of this architecture using parallel and incremental MTL paradigms to address the HO management at layer 3 and initial MCS use cases at layer 2, jointly. The results of the extensive Offline-Evaluation showed that the MTL models achieved similar performance as those achieved using dedicated single-task models. This proved the efficiency of the AE in compressing the inputs without losing important information. It enabled us to use this compressed representation of the input feature space in a shared manner to address HO and the initial MCS use case jointly. Thus, lifting the burden of the network to train a model from scratch when a new use case is added to the setup.

After tackling the issues characterized by the infrastructure-based single-access technology scenarios, we moved into a higher level of complexity by switching the second axis to multi-access technology. In particular, Chapters 5 and 6 dealt with the fairness issues due to the coexistence of LAA, and LTE-U with WiFi, since, without guaranteeing their fairness, these mobile technologies cannot function in the unlicensed spectrum. In Chapter 5, after deeply studying LAA´s channel access mechanism, we identified some significant differences between the LAA and WiFi CW adaptation procedures. These differences pose challenges in achieving fairness for LAA in the unlicensed spectrum.

The main contribution of Chapter 5 is the AI/ML solution that improved the CW adaptation procedure of LAA by inferring the number of NACKs that can be received in a TxOP. The proposed solution achieved a good trade-off between WiFi fairness and

LAA performance, overcoming the schemes that sacrificed LAA performance in terms of throughput and latency for better coexistence with WiFi. Connected to this, the contribution in Chapter 6 aimed to quantify the fairness in coexistence scenarios instead of making qualitative comparisons of throughput and latency CDFs. In particular, we proposed a framework based on the KS-test that statistically quantified the fairness of LAA and LTE-U when they coexist with WiFi. The evaluation study based on this framework revealed that LAA provides better fairness, while LTE-U introduces more collisions. Such framework and its possible extension can benefit the fairness evaluations where CDF curves intertwine, making it challenging to conclude fairness.

Finally, in Chapters 7 and 8, our attention shifted to address the complexity that exists in the limited infrastructure and single-access technology scenarios. We learned that in such scenarios, due to the absence of BSs, the mobile devices operate autonomously, making decisions based on their perspective of the channel. This introduces challenges for mobile devices to operate without BS support. This thesis identified one such challenge of energy-performance trade-off due to continuous sensing in NR-V2X technology.

In this respect, the main contribution in Chapter 7 is the first *ns-3* open-source and standard-compliant simulator for NR-V2X. This chapter included the results of the deep simulation study that, among other findings, revealed the energy-performance trade-off in NR-V2X. This trade-off is then studied to propose a novel fuzzy logic-based partial sensing resource selection mechanism in Chapter 8. The proposed scheme automatically balanced the energy consumption and the UE performance, allowing improvements in both energy efficiency and PRR. The main contribution of our scheme is that it can automatically adjust the sensing duty cycle of the NR-V2X UE based on the dynamics of the V2X scenarios. Hence, it exploits the energy-performance trade-off automatically, which is impossible by employing a static configuration.

In a nutshell, the main contributions of this Ph.D. thesis are twofold. On the one hand, using the two identified complexity axes, it has presented an in-depth research on some of the most relevant use cases in the RAN and their inherent complexities impacting the RAN management. Building on that, the thesis proposed novel AI/ML solutions and demonstrated their effectiveness in automating the increasing complexity of the RAN management over the benchmark schemes in realistic simulation scenarios. On the other hand, the thesis also contributed to the open-source community by extending and implementing new open-source simulation models in *ns-3* and 5G-LENA simulators. We hope that our contribution can help the research community to perform wide-scaled studies, such as ours, by coupling these models with different AI/ML frameworks.

## 9.2 Future work

The AI/ML technologies are expected to play a vital role in achieving zero-touch 6G networks. By leveraging these technologies, mobile networks can benefit in terms of cost, energy consumption, reliability, and operational efficiency. In 6G, AI/ML will augment networks to react quickly and efficiently to unpredictable situations and traffic needs through predictive orchestration mechanisms. To ensure effective network orchestration, decisions should be based on a holistic end-to-end perspective of the network and enforce

actions from devices to RAN disaggregated functions, edge computing, core elements, and cloud components [201].

In this context, the work that we have proposed can be extended in future works considering: 1) not only the RAN but an end-to-end perspective where network orchestration is achieved through AI/ML; 2) more use cases can be analyzed taking into account this end-to-end vision; 3) AI/ML training costs, which can become huge, pose an issue that has to be addressed considering the computational complexity and the energy efficiency during the design process. In this line, the framework we have proposed based on multi-task methodologies can also be extended beyond the RAN; 4) other AI/ML training approaches can also be employed. We have focused on AI/ML in a centralized manner, and we believe this approach should be continued; however, we also consider that distributed AI/ML allowing distributed training among nodes in the network should be considered to explore new horizons and evaluate benefits and advantages compared to centralized solutions; 5) by also taking into account the recent development in the *ns-3* simulator to perform AI/ML studies. When conducting our research, the *ns-3* simulator did not support any specialized module to link *ns-3* with external AI/ML frameworks. However, Hao Yin et al. recently contributed the "ns3-ai" module to connect Python-based AI/ML frameworks with *ns-3* using the shared memory concept [202]. This module can facilitate future studies extending our proposed AI/ML solutions.

# Bibliography

[1] "GSMA," https://www.gsma.com/, accessed: 2021-08-01.

[2] J. Wang, H. Roy, and C. Kelly, "OpenRAN: The next generation of radio access networks,," https://telecominfraproject.com/openran/, Nov 2019.

[3] 3GPP TS 32.500 V0.5.1, *Self-Organising Networks (SON): Concepts and requirements*, Release 8.

[4] Huawei, "White Paper: Next Generation SON for 5G," https://www.huawei.com/en/huaweitech/industry-insights/outlook/mobile-broadband/insights-reports/next-generation-son-for-5g, accessed: 2023-06-18.

[5] J. Moysen and L. Giupponi, "From 4G to 5G: Self-organized network management meets machine learning," *Computer Communications*, vol. 129, pp. 248–268, Sept 2018.

[6] N. Baldo, L. Giupponi, and J. Mangues, "Big Data Empowered Self Organized Networks," in *Proc. IEEE 20th European Wireless Conference*, Barcelona, Spain, May 2014, pp. 1–8.

[7] M. Gheisari, G. Wang, and M. Z. A. Bhuiyan, "A survey on deep learning in big data," in *IEEE International Conference on Computational Science and Engineering (CSE) and IEEE International Conference on Embedded and Ubiquitous Computing (EUC)*, vol. 2, 2017, pp. 173–180.

[8] C. Zhang, P. Patras, and H. Haddadi, "Deep learning in mobile and wireless networking: A survey," *IEEE Communications Surveys Tutorials*, vol. 21, no. 3, pp. 2224–2287, 2019.

[9] Gino Masini, Yin Gao, Sasha Sirotkin, "Artificial Intelligence and Machine Learning in NG-RAN: New Study in RAN3," https://www.3gpp.org/ftp/Information/Highlights/2021_Issue02/mobile/index.html#p=7/.

[10] GSMA, "The Mobile Industry and AI," https://www.gsma.com/betterfuture/resources/the-mobile-industry-and-ai, accessed: 2023-03-17.

[11] O-RAN Alliance, "O-RAN Use Cases and Deployment Scenarios," https://www.o-ran.org/resources.

[12] 3GPP TS 29.520, *Network Data Analytics Services; Stage 3, Rel. 15*, July 2018.

[13] 3GPP TR 37.817, *Study on enhancement for Data Collection for NR and EN-DC, Rel. 17*, April 2022.

[14] 3GPP TR 38.843, *Study on Artificial Intelligence (AI)/Machine Learning (ML) for NR air interface (Rel. 18)*, June. 2022.

[15] Ericsson, "Mobile radio access networks: What policy makers need to know," https://www.ericsson.com/en/blog/2020/9/ran-what-policy-makers-need-to-know, accessed: 2023-03-17.

[16] O-RAN.WG2, *O-RAN Working Group 2, Non-RT RIC & A1 Interface: Use Cases and Requirements, v06.00*, April 2021.

[17] O-RAN.WG3, *O-RAN Working Group 3, Use Cases and Requirements, v01.00*, August 2021.

[18] C. Zhang, P. Patras, and H. Haddadi, "Deep learning in mobile and wireless networking: A survey," *IEEE Communications Surveys Tutorials*, vol. 21, no. 3, pp. 2224–2287, Mar 2019.

[19] Y. Yao, H. Zhou, and M. Erol-Kantarci, "Deep reinforcement learning-based radio resource allocation and beam management under location uncertainty in 5g mm wave networks," in *2022 IEEE Symposium on Computers and Communications (ISCC)*, 2022, pp. 1–6.

[20] S. Vittal and A. F. A, "Self optimizing network slicing in 5g for slice isolation and high availability," in *2021 17th International Conference on Network and Service Management (CNSM)*, 2021, pp. 125–131.

[21] S. S. Mwanje and A. Mitschele-Thiel, "Distributed cooperative Q-learning for mobility-sensitive handover optimization in LTE SON," in *Proc. IEEE Symposium on Computers and Communications (ISCC)*, Funchal, Portugal, Jun 2014, pp. 1–6.

[22] Z. Becvar and P. Mach, "Adaptive Hysteresis Margin for Handover in Femtocell Networks," in *Proc. 6th International Conference on Wireless and Mobile Communications*, Valencia, Spain, Nov 2010, pp. 256–261.

[23] J. Wu, J. Liu, Z. Huang, and S. Zheng, "Dynamic fuzzy Q-learning for handover parameters optimization in 5G multi-tier networks," in *Proc. International Conference on Wireless Communications Signal Processing (WCSP)*, Nanjing, China, Dec 2015, pp. 1–5.

[24] Z. Wang, L. Li, Y. Xu, H. Tian, and S. Cui, "Handover Optimization via Asynchronous Multi-User Deep Reinforcement Learning," in *Proc. IEEE International Conference on Communications (ICC)*, Kansas City, MO, USA, Jul 2018, pp. 1–6.

[25] C. Lee, H. Cho, S. Song, and J. Chung, "Prediction-Based Conditional Handover for 5G mm-Wave Networks: A Deep-Learning Approach," *IEEE Vehicular Technology Magazine*, vol. 15, no. 1, pp. 54–62, 2020.

[26] Y. Xu, W. Xu, Z. Wang, J. Lin, and S. Cui, "Load balancing for ultradense networks: A deep reinforcement learning-based approach," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 9399–9412, 2019.

[27] P. E. Iturria-Rivera and M. Erol-Kantarci, "Qos-aware load balancing in wireless networks using clipped double q-learning," in *2021 IEEE 18th International Conference on Mobile Ad Hoc and Smart Systems (MASS)*, 2021, pp. 10–16.

[28] A. Giannopoulos, S. Spantideas, N. Kapsalis, P. Gkonis, L. Sarakis, C. Capsalis, M. Vecchio, and P. Trakadas, "Supporting intelligence in disaggregated open radio access networks: Architectural principles, ai/ml workflow, and use cases," *IEEE Access*, vol. 10, pp. 39 580–39 595, 2022.

[29] H. Zhang, H. Zhou, and M. Erol-Kantarci, "Federated deep reinforcement learning for resource allocation in o-ran slicing," in *GLOBECOM 2022 - 2022 IEEE Global Communications Conference*, 2022, pp. 958–963.

[30] B. Agarwal, M. A. Togou, M. Ruffini, and G.-M. Muntean, "Qoe-driven optimization in 5g o-ran-enabled hetnets for enhanced video service quality," *IEEE Communications Magazine*, vol. 61, no. 1, pp. 56–62, 2023.

[31] A. Lacava, M. Polese, R. Sivaraj, R. Soundrarajan, B. S. Bhati, T. Singh, T. Zugno, F. Cuomo, and T. Melodia, "Programmable and customized intelligence for traffic steering in 5g networks using open ran architectures," *IEEE Transactions on Mobile Computing*, pp. 1–16, 2023.

[32] L. Bonati, S. D'Oro, M. Polese, S. Basagni, and T. Melodia, "Intelligence and learning in o-ran for data-driven nextg cellular networks," *IEEE Communications Magazine*, vol. 59, no. 10, pp. 21–27, 2021.

[33] Sandra Lagén Morancho, *Coordination strategies for interference management in MIMO dense cellular networks (PhD Thesis)*. Universitat Politècnica de Catalunya, 2017.

[34] Koutlia Aikaterini, *Radio Resource Management Strategies for Interference Mitigation in 4G Heterogeneous Wireless Networks (PhD Thesis)*. Universitat Politècnica de Catalunya, 2016.

[35] Marco Miozzo, *Energy sustainability of next generation cellular networks through learning techniques (PhD Thesis)*. Universitat Politècnica de Catalunya, 2018.

[36] 3GPP TS 22.261, *Service requirements for the 5G system; Stage 1, Rel. 19*, March 2023.

[37] Huawei, *6G: The Next Horizon (From Connected People and Things to Connected Intelligence)*, https://www-file.huawei.com/-/media/corp2020/pdf/tech-insights/1/6g-white-paper-en.pdf?la=en, accessed: 2023-06-18.

[38] S. Sesia, I. Toufik, and M. Baker, *LTE - The UMTS Long Term Evolution - from theory to practice*. Wiley, 2011.

[39] M. Chui, J. Manyika, M. Miremadi, N. Henke, R. Chung, P. Nel, and S. Malhotra, "Notes from the AI frontier insights from hundereds of use cases," *McKinsey Global Institute*, 2018.

[40] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.

[41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[42] Trinh Hoang Duy, *Data analytics for mobile traffic in 5G networks using machine learning techniques (PhD Thesis).* Universitat Politècnica de Catalunya, 2020.

[43] Y. Huang, S. Liu, C. Zhang, X. You, and H. Wu, "True-data testbed for 5g/b5g intelligent network," *Intelligent and Converged Networks*, vol. 2, no. 2, pp. 133–149, 2021.

[44] O. A. Fernando, H. Xiao, and J. Spring, "Developing a testbed with p4 to generate datasets for the analysis of 5g-mec security," in *2022 IEEE Wireless Communications and Networking Conference (WCNC)*, 2022, pp. 2256–2261.

[45] Biljana Bojovic, *Cellular and Wi-Fi technologies evolution: from complementarity to competition (PhD Thesis).* Universitat Politècnica de Catalunya, 2022.

[46] Mathworks, "Physical Layer Abstraction for System-Level Simulation," https://se.mathworks.com/help/wlan/ug/physical-layer-abstraction-for-system-level-simulation.html, accessed: 2023-04-10.

[47] K. Wehrle, M. Günes, and J. Gross, *Modeling and Tools for Network Simulation.* Springer Berlin Heidelberg, 2010.

[48] ns-3, "ns-3 Manual," https://www.nsnam.org/docs/release/3.38/manual/singlehtml/index.html, accessed: 2023-04-12.

[49] Syed Muhammad Asif Qamar and Zoraze Ali, *3GPP LTE vs. IEEE 802.11p/WAVE: Competition or Coexistence? (Master Thesis).* Department of Electrical and Information Technology Faculty of Engineering, Lund University, Sweden, 2013.

[50] ns-3, "ns-3 Tutorial," https://www.nsnam.org/docs/release/3.38/tutorial/singlehtml/index.html, accessed: 2023-04-12.

[51] ns-3, "ns-3 Model Library," https://www.nsnam.org/docs/release/3.38/models/singlehtml/index.html, accessed: 2023-04-12.

[52] ns-3, "Network Simulator," http://code.nsnam.org/ns-3-dev, accessed: 2023-04-12.

[53] S. Sakai, G. Gambugge, R. Takaki, J. Seki, J. Bazzo, and J. P. Miranda, "Performance comparison of a custom emulation-based test environment against a real-world lte testbed," in *Proceedings of the 2015 Workshop on ns-3*, 2015, pp. 106–111.

[54] A. Marinescu, I. Macaluso, and L. A. DaSilva, "System level evaluation and validation of the ns-3 lte module in 3gpp reference scenarios," in *Proceedings of the 13th ACM Symposium on QoS and Security for Wireless and Mobile Networks*, 2017, pp. 59–64.

[55] G. . V9.0.0, *Further advancements for E-UTRA physical layer aspects*, Mar. 2010.

[56] K. Koutlia, B. Bojovic, Z. Ali, and S. Lagén, "Calibration of the 5g-lena system level simulator in 3gpp reference scenarios," *Simulation Modelling Practice and Theory*, vol. 119, p. 102580, 2022.

[57] Zoraze Ali, "ns-3 extension to simulate coverage holes and deterministic handover in LTE," https://github.com/ZorazeAli/ns-3-dev-obstacle, accessed: 2023-06-18.

[58] LENA, "LAA and LTE-U implementation," https://bitbucket.org/cttc-lena/ns-3-lena-dev-lte-u/src/, accessed: 2023-06-18.

[59] Fraunhofer, "ELENA extension for ns-3 LTE module," https://gitlab.cc-asp.fraunhofer.de/elena-ns3-lte/elena, accessed: 2023-06-18.

[60] ns-3, "Public Safety Communications, Models to support applications and scenarios for first responders," https://apps.nsnam.org/app/publicsafetylte/, accessed: 2023-06-18.

[61] ns-3, "mmWave Cellular Network Simulator," https://apps.nsnam.org/app/mmwave/, accessed: 2023-06-18.

[62] CTTC, "5G LENA Project," https://5g-lena.cttc.es/download/.

[63] Jessica Moysen, *Self Organisation for 4G/5G Networks(PhD Thesis)*. Universitat Politècnica de Catalunya, 2016.

[64] N. Baldo, "Cognitive radios and networks," Ph.D. dissertation, UNIVERSIT**A**' DI PADOVA, Jan 2009, phD thesis.

[65] cplusplus.com, "The system() function of C++," https://cplusplus.com/reference/cstdlib/system/, accessed: 2023-06-18.

[66] C. M.Bishop, *Pattern Recognition and Machine Learning*. NY,USA: Springer, 2006.

[67] B. D. Ripley, *Pattern Recognition and Neural Networks*. Cambridge,UK: Cambridge University Press, 1996.

[68] Ericsson, "Reducing mobility interruption time in 5G networks," https://www.ericsson.com/en/blog/2020/4/reducing-mobility-interruption-time-5g-networks/, accessed: 2022-10-03.

[69] G. T. . V10.4.0, *Radio measurement collection for Minimization of Drive Tests (MDT); Overall description*, Release 10.

[70] D. López-Pérez, M. Ding, H. Claussen, and A. H. Jafari, "Towards 1 Gbps/UE in Cellular Systems: Understanding Ultra-Dense Small Cell Deployments," *IEEE Communications Surveys Tutorials*, vol. 17, no. 4, pp. 2078–2101, 2015.

[71] 3GPP TS 23.234 V6.10.0, *3GPP system to Wireless Local Area Network (WLAN) interworking; System description*, September 2006.

[72] 3GPP TS 24.327 V8.6.0, *Mobility between 3GPP Wireless Local Area Network (WLAN) interworking (I-WLAN) and 3GPP systems; General Packet Radio System (GPRS) and 3GPP I-WLAN aspects; Stage 3*, June 2010.

[73] 3GPP TS 36.300 V12.10.0, *Overall description; Stage 2*, June 2016.

[74] 3GPP TS 36.300 V13.14.0, *Overall description; Stage 2*, March 2020.

[75] 3GPP TS 36.300 V14.12.0, *Overall description; Stage 2*, March 2020.

[76] Qualcomm Technologies Inc., "LTE-U coexistence mechanism," May 2015.

[77] G. . V16.0.0, *TSG RAN; NR; Study on NR-based access to unlicensed spectrum*, Dec. 2018.

[78] B. Chen, J. Chen, Y. Gao, and J. Zhang, "Coexistence of lte-laa and wi-fi on 5 ghz with corresponding deployment scenarios: A survey," *IEEE Comm Surveys Tutorials*, vol. 19, pp. 7–32, 2017.

[79] 3GPP TR 36.889 V13.0.0, *Study on Licensed-Assisted Access to Unlicensed Spectrum*, Release 13.

[80] 3GPP TS 36.213 V14.1.0, *Physical layer procedures*, Release 14.

[81] LTE-U Forum, "LTE-U SDL Coexistence Specifications, V1.3," October 2018.

[82] LTE-U Forum, "LTE-U CSAT Procedure TS, V1.3," October 2015.

[83] Qualcomm Technologies Inc., "LTE-U Technology and Coexistence," San Diego, USA, May 2015.

[84] U.S. Department of Transportation, *Vehicle-to-vehicle Communication Technology*, https://www.nhtsa.gov/sites/nhtsa.dot.gov/files/documents/v2v_fact_sheet_101414_v2a.pdf, July 2017.

[85] 3GPP TR 22.886 V16.2.0, *Study on Enhancement of 3GPP Support for 5G V2X Services*, Dec. 2018.

[86] 5G Automotive Association, *5G Automotive Vision, White Paper*, Oct. 2015.

[87] D. Jiang and L. Delgrossi, "IEEE 802.11p: Towards an international standard for wireless access in vehicular environments," in *VTC Spring 2008 - IEEE Vehicular Technology Conference*, 2008, pp. 2036–2040.

[88] G. Naik, B. Choudhury, and J. Park, "IEEE 802.11bd and 5G NR V2X: Evolution of radio access technologies for V2X communications," *IEEE Access*, vol. 7, pp. 70 169–70 184, 2019.

[89] R. Molina-Masegosa and J. Gozalvez, "LTE-V for sidelink 5G V2X vehicular communications: A new 5G technology for short-range vehicle-to-everything communications," *IEEE Vehicular Technology Magazine*, vol. 12, no. 4, pp. 30–39, 2017.

[90] 3GPP TS 22.185 V16.0.0, *Service Requirements for V2X Services*, Jun. 2020.

[91] 5GAA, *An Assessment of LTE-V2X (PC5) and 802.11p Direct Communications Technologies for Improved Road Safety in the EU*, Dec. 2017.

[92] 3GPP TR 38.885 V16.0.0, *Study on NR Vehicular to Everything (V2X)*, Mar. 2019.

[93] 3GPP RP-190766, *5G V2X with NR Sidelink (Release 16), Work Item description*, Dec. 2018.

[94] S. Parkvall, E. Dahlman, A. Furuskar, and M. Frenne, "NR: The New 5G Radio Access Technology," *IEEE Communications Standards Magazine*, vol. 1, no. 4, pp. 24–30, Dec 2017.

[95] 3GPP TS 23.303 V20.0.0, *Proximity-based services (ProSe), Stage 2*, June 2020.

[96] 3GPP TR 36.843 , V12.0.1, *Study on LTE Device to Device Proximity Services; Radio Aspects*, Mar. 2014.

[97] 3GPP TS 38.101, *User Equipment (UE) radio transmission and reception; Part 1: Range 1 Standalone, Rel. 17*, Mar. 2021.

[98] 3GPP TR 38.886, *V2X Services based on NR; User Equipment (UE) radio transmission and reception. Rel. 16*, Mar. 2021.

[99] J. Wang, C. Jiang, H. Zhang, Y. Ren, K. C. Chen, and L. Hanzo, "Thirty Years of Machine Learning: The Road to Pareto-Optimal Wireless Networks," *IEEE Communications Surveys Tutorials*, vol. 22, no. 3, pp. 1472–1514, 2020.

[100] P. V. Klaine, M. A. Imran, O. Onireti, and R. D. Souza, "A Survey of Machine Learning Techniques Applied to Self-Organizing Cellular Networks," *IEEE Communications Surveys Tutorials*, vol. 19, no. 4, pp. 2392–2431, 2017.

[101] I. Pappalardo, A. Zanella, and M. Zorzi, "Upper Bound Analysis of the Handover Performance in HetNets," *IEEE Communications Letters*, 2017.

[102] P. Tseng, K. Feng, and C. Huang, "POMDP-Based Cell Selection Schemes for Wireless Networks," *IEEE Communications Letters*, vol. 18, no. 5, pp. 797–800, 2014.

[103] F. Guidolin, I. Pappalardo, A. Zanella, and M. Zorzi, "Context-Aware Handover Policies in HetNets," *IEEE Transactions on Wireless Communications*, vol. 15, no. 3, pp. 1895–1906, 2016.

[104] ns-3, "Buildings Module Design Documentation," http://www.nsnam.org/docs/release/3.22/models/html/buildings-design.html, accessed: 2023-06-18.

[105] Christer Ericson, *Real-Time Collision Detection.* San Francisco,CA: Morgan Kaufmann, 2005.

[106] CTTC, "LTE-EPC Network Simulator (LENA)," https://vimeo.com/353023881, accessed: 2023-06-18.

[107] M. Nogaard and O. Ravn and N. K. Poulsen and L. K. Hansen, *Neural networks for modelling and control of dynamic systems.* London,UK: Springer-Verlag, 2000.

[108] 3GPP TR 32.862 V14.0.0, *Study on Key Quality Indicators (KQIs) for service experience.*

[109] ITU-T Recommendation, *Estimating end-to-end performance in IP networks for data applications*, G.1030, Series G.

[110] M. Khan and U. Toseef, "User utility function as quality of experience (QoE)," in *Proc. of the Tenth International Conference on Networks*, The Netherlands, Jan 2011, pp. 99–104.

[111] J. Mendoza, I. de-la Bandera, D. Palacios, and R. Barco, "Qoe optimization in a live cellular network through rlc parameter tuning," *Sensors*, vol. 21, no. 16, 2021.

[112] B. Ripley and W. Venables, "Feed-Forward Neural Networks and Multinomial Log-Linear Models," http://cran.r-project.org/web/packages/nnet/nnet.pdf, accessed: 2023-06-18.

[113] J. Sola and J. Sevilla, "Importance of input data normalization for the application of neural networks to complex industrial problems," *Nuclear Science, IEEE Transactions on*, vol. 44, no. 3, pp. 1464–1468, Jun 1997.

[114] Y. Bengio, A. Courville, and P. Vincent, "Representation Learning: A Review and New Perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1464–1468, Jun 2013.

[115] Y. Ju, J. Guo, and S. Lui, "A Deep Learning Method Combined Sparse Autoencoder with SVM," in *Proc. International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery*, Xi'an, China, Sept 2015, pp. 257–260.

[116] M. Ghifary, W. B. Kleijn, M. Zhang, and D. Balduzzi, "Domain Generalization for Object Recognition with Multi-task Autoencoders," in *Proc. International Conference on Computer Vision*, Santiago, Chile, 2015, pp. 2551–2559.

[117] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. early access, pp. 1–1, Feb 2021.

[118] M. Kanakis, D. Bruggemann, S. Saha, S. Georgoulis, A. Obukhov, and L. Van Gool, "Reparameterizing Convolutions for Incremental Multi-Task Learning Without Task Interference," in *Proc. European Conference on Computer Vision (ECCV)*, Glasgow, UK, 2020, pp. 689–707.

[119] J. P. Leite, P. H. P. de Carvalho, and R. D. Vieira, "A flexible framework based on reinforcement learning for adaptive modulation and coding in OFDM wireless systems," in *Proc. IEEE Wireless Communications and Networking Conference (WCNC)*, Paris, France, Jun 2012, pp. 809–814.

[120] R. Bruno, A. Masaracchia, and A. Passarella, "Robust Adaptive Modulation and Coding (AMC) Selection in LTE Systems Using Reinforcement Learning," in *Proc. IEEE 80th Vehicular Technology Conference (VTC2014-Fall)*, Vancouver, BC, Canada, Dec 2014, pp. 1–6.

[121] L. Zhang, J. Tan, Y. Liang, G. Feng, and D. Niyato, "Deep Reinforcement Learning-Based Modulation and Coding Scheme Selection in Cognitive Heterogeneous Networks," *IEEE Transactions on Wireless Communications*, vol. 18, no. 6, pp. 3281–3294, Apr 2019.

[122] A. Shahmansoori, "Sparse Bayesian Multi-Task Learning of Time-Varying Massive MIMO Channels With Dynamic Filtering," *IEEE Wireless Communications Letters*, vol. 9, no. 6, pp. 871–874, Feb 2020.

[123] A. Rago, G. Piro, G. Boggia, and P. Dini, "Multi-Task Learning at the Mobile Edge: An Effective Way to Combine Traffic Classification and Prediction," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 9, pp. 10 362–10 374, Jun 2020.

[124] N. Ye, X. Li, H. Yu, L. Zhao, W. Liu, and X. Hou, "DeepNOMA: A Unified Framework for NOMA Using Deep Multi-Task Learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 4, pp. 2208–2225, Jan 2020.

[125] Z. Ali, L. Giupponi, M. Miozzo, and P. Dini, "Multi-Task Learning for Efficient Management of Beyond 5G Radio Access Network Architectures," *IEEE Access*, vol. 9, pp. 158 892–158 907, 2021.

[126] Z. Ali, M. Miozzo, L. Giupponi, P. Dini, S. Denic, and S. Vassaki, "Recurrent Neural Networks for Handover Management in Next-Generation Self-Organized Networks," in *Proc. IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*, London, UK, Sept 2020, pp. 1–6.

[127] S. Hochreiter and J. A. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov 1997.

[128] J. Wang, J. Tang, Z. Xu, Y. Wang, G. Xue, X. Zhang, and D. Yang, "Spatiotemporal modeling and prediction in cellular networks: A big data enabled deep learning approach," in *Proc. IEEE INFOCOM - IEEE Conference on Computer Communications*, May 2017, pp. 1–9.

[129] H. D. Trinh, L. Giupponi, and P. Dini, "Mobile traffic prediction from raw data using LSTM networks," in *Proc. IEEE 29th Annual International Symposium on Personal, Indoor, and Mobile Radio Communication (PIMRC)*, Bologna, Italy, Sept 2018, pp. 1827–1832.

[130] F. A. Gers, J. A. Schmidhuber, and F. A. Cummins, "Learning to forget: Continual prediction with lstm," *Neural Comput.*, vol. 12, no. 10, pp. 2451–2471, Oct 2000.

[131] M. Crawshaw, "Multi-task learning with deep neural networks: A survey," *ArXiv*, vol. abs/2009.09796, Sept 2020.

[132] M. Miozzo, Z. Ali, L. Giupponi, and P. Dini, "Distributed and Multi-Task Learning at the Edge for Energy Efficient Radio Access Networks," *IEEE Access*, vol. 9, pp. 12 491–12 505, 2021.

[133] A. M. Dai and Q. V. Le, "Semi-supervised sequence learning," in *Proc. Advances in Neural Information Processing Systems 28*, 2015, pp. 3079–3087.

[134] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, "Stacked convolutional auto-encoders for hierarchical feature extraction," in *Proc. 21st International Conference on Artificial Neural Networks*, Espoo, Finland, Jun 2011, pp. 52–59.

[135] R. M. French, "Catastrophic forgetting in connectionist networks," *Trends in Cognitive Sciences*, vol. 3, no. 4, pp. 128–135, Apr 1999.

[136] F. Chollet *et al.*, "Keras," https://keras.io, 2015, accessed: 2023-06-18.

[137] Small Cell Forum, "Combining the benefits of licensed and unlicensed technologies (Release 7.0)," https://scf.io/en/documents/094_-_Combining_the_benefits_of_licensed_and_unlicensed_technologies.php, accessed: 2023-06-18.

[138] Y. Li, G. Drory, Y. Cohen, B. Yang, M. Fischer, S. Levy, S. Verma, and S. Adhikar, *WIFI-coordinated LAA-LTE*, March 31, 2016, U.S. Patent 20 160 095 110.

[139] Y. Song, K. W. Sung, and Y. Han, "Coexistence of Wi-Fi and Cellular With Listen-Before-Talk in Unlicensed Spectrum," *IEEE Communications Letters*, vol. 20, pp. 161–164, 2016.

[140] X. Yan, H. Tian, and C. Qin, "A Markov-Based Modelling with Dynamic Contention Window Adaptation for LAA and WiFi Coexistence," in *2017 IEEE 85th Vehicular Technology Conference (VTC Spring)*, 2017, pp. 1–6.

[141] 3GPP TS 36.211, *IEEE Standard 802.11, Part 11:Wireless LAN Medium Access Control (MAC) and Physical Layer(PHY) Specifications*, 2007.

[142] T. Tao, F. Han, and Y. Liu, "Enhanced LBT algorithm for LTE-LAA in unlicensed band," in *PIMRC*, Hong Kong, China, 2015, pp. 1907–1911.

[143] F. Hao, C. Yongyu, H. Li, J. Zhang, and W. Quan, "Contention window size adaptation algorithm for LAA-LTE in unlicensed band," in *ISWCS*, Poznan, Poland, 2016, pp. 476–480.

[144] IEEE, "IEEE standard for information technology-telecommunications and information exchange between systems local and metropolitan area networks," *IEEE Std*, 2012.

[145] E. Perahia and R. Stacey, *Next Generation Wireless LANs 802.11n and 802.11ac*. Cambridge University Press, 2013.

[146] M. S. Gast, *802.11 ac: A survival guide*. O'Reilly Media, Inc., 2013.

[147] M. Gast, *802.11 n: A survival guide.* O'Reilly Media, Inc., 2012.

[148] ns-3, "ns-3 Tutorial," https://www.nsnam.org/wiki/LAA-WiFi-Coexistence, accessed: 2023-06-18.

[149] Qualcomm et al., *WF on CW adjustment based on eNB sensing 3GPP TSG RAN WG1 #82*, August 2015.

[150] S. Dama, A. Kumar, and K. Kuchi, "Performance evaluation of laa-lbt based lte and wlan's co-existence in unlicensed spectrum," in *2015 IEEE Globecom Workshops (GC Wkshps)*, Dec 2015, pp. 1–6.

[151] B. Jia and M. Tao, "A channel sensing based design for lte in unlicensed bands," in *Communication Workshop (ICCW), 2015 IEEE International Conference on*, June 2015, pp. 2332–2337.

[152] S. Choi and S. Park, "Performance analysis of various co-existence methods with wi-fi in unlicensed bands," in *Information and Communication Technology Convergence (ICTC), 2015 International Conference on*, Oct 2015, pp. 729–732.

[153] D. Vose, *Risk Analysis - A Quantitative Guide*, 3rd ed. Wiley, 2008.

[154] J. D. G. a. John W. Pratt, *Concepts of Nonparametric Theory*, 1st ed., ser. Springer Series in Statistics. Springer-Verlag New York, 1981.

[155] I. Banri and T. Ayumu, *Open Innovation, Productivity, and Export: Evidence from Japanese firms*, Feb 2013, discussion papers.

[156] R. Jain, *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling.* Wiley, 1990.

[157] "R statistical functions," https://stat.ethz.ch/R-manual/R-devel/library/stats/html/stats-package.html.

[158] M. Gonzalez-Martín, M. Sepulcre, R. Molina-Masegosa, and J. Gozalvez, "Analytical Models of the Performance of C-V2X Mode 4 Vehicular Communications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 2, pp. 1155–1166, 2019.

[159] R. Molina-Masegosa, J. Gozalvez, and M. Sepulcre, "Comparison of IEEE 802.11p and LTE-V2X: An evaluation with periodic and aperiodic messages of constant and variable size," *IEEE Access*, vol. 8, pp. 121 526–121 548, 2020.

[160] B. Toghi, M. Saifuddin, H. N. Mahjoub, M. O. Mughal, Y. P. Fallah, J. Rao, and S. Das, "Multiple Access in Cellular V2X: Performance Analysis in Highly Congested Vehicular Networks," in *2018 IEEE Vehicular Networking Conference (VNC)*, 2018, pp. 1–8.

[161] T. Zugno, M. Drago, M. Giordani, M. Polese, and M. Zorzi, "Toward Standardization of Millimeter-Wave Vehicle-to-Vehicle Networks: Open Challenges and Performance Evaluation," *IEEE Communications Magazine*, vol. 58, no. 9, pp. 79–85, 2020.

[162] K. Ganesan, J. Lohr, P. B. Mallick, A. Kunz, and R. Kuchibhotla, "NR Sidelink Design Overview for Advanced V2X Service," *IEEE Internet of Things Magazine*, vol. 3, no. 1, pp. 26–30, 2020.

[163] S. Lien, D. Deng, C. Lin, H. Tsai, T. Chen, C. Guo, and S. Cheng, "3GPP NR Sidelink Transmissions Toward 5G V2X," *IEEE Access*, vol. 8, pp. 35 368–35 382, 2020.

[164] M. H. C. Garcia, A. Molina-Galan, M. Boban, J. Gozalvez, B. Coll-Perales, T. Şahin, and A. Kousaridas, "A Tutorial on 5G NR V2X Communications," *IEEE Communications Surveys Tutorials*, pp. 1–1, 2021.

[165] C. Campolo, A. Molinaro, F. Romeo, A. Bazzi, and A. O. Berthet, "5G NR V2X: On the Impact of a Flexible Numerology on the Autonomous Sidelink Mode," in *2019 IEEE 2nd 5G World Forum (5GWF)*, 2019, pp. 102–107.

[166] R. Rouil, F. J. Cintrón, A. B. Mosbah, and S. Gamboa, "Implementation and Validation of an LTE D2D Model for ns-3," in *Workshop on Ns-3*, 2017.

[167] F. Eckermann, M. Kahlert, and C. Wietfeld, "Performance Analysis of C-V2X Mode 4 Communication Introducing an Open-Source C-V2X Simulator," in *2019 IEEE 90th Vehicular Technology Conference (VTC2019-Fall)*, 2019, pp. 1–5.

[168] M. Malinverno, F. Raviglione, C. Casetti, C.-F. Chiasserini, J. Mangues-Bafalluy, and M. Requena-Esteso, "A multi-stack simulation framework for vehicular applications testing," in *Proceedings of the 10th ACM Symposium on Design and Analysis of Intelligent Vehicular Networks and Applications*, ser. DIVANet '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 17–24. [Online]. Available: https://doi.org/10.1145/3416014.3424603

[169] G. Cecchini, A. Bazzi, B. M. Masini, and A. Zanella, "LTEV2Vsim: An LTE-V2V Simulator for the Investigation of Resource Allocation for Cooperative Awareness," in *2017 5th IEEE International Conference on Models and Technologies for Intelligent Transportation Systems (MT-ITS)*, 2017, pp. 80–85.

[170] M. Drago, T. Zugno, M. Polese, M. Giordani, and M. Zorzi, "MilliCar: An ns-3 Module for mmWave NR V2X Networks," in *Workshop on Ns-3*, 2020.

[171] G. . V16.0.0, *TSG RAN; NR; Physical layer; General Description*, Jan. 2020.

[172] *TSG RAN; NR; Physical Channels and Modulation*, 3GPP TS 38.211 V16.4.0, Jan. 2021.

[173] G. . V16.4.0, *TSG RAN; NR; Physical Layer Procedures for Data*, Jan. 2021.

[174] G. 38.202, *TSG RAN; NR; Services Provided by the Physical Layer*, Release 15, v16.2.0, Sep. 2020.

[175] G. . V16.2.0, *TSG RAN; NR; Radio Resource Control (RRC); Protocol specification*, Sep. 2020.

[176] 3GPP TS 36.331, *Radio Resource Control (RRC); Protocol specification, Rel. 16*, Jul. 2020.

[177] *Final Report of 3GPP TSG RAN WG1 100bis-e v1.0.0*, 3GPP TSG RAN WG1 Meeting 101-e, Apr. 2020.

[178] M. H. C. Garcia, A. Molina-Galan, M. Boban, J. Gozalvez, B. Coll-Perales, T. Şahin, and A. Kousaridas, "A Tutorial on 5G NR V2X Communications," *IEEE Communications Surveys Tutorials*, pp. 1–1, 2021.

[179] G. . V16.3.0, *TSG RAN; NR; Medium Access Control (MAC) Protocol Specification*, Jan. 2021.

[180] G. . M. RP-193231, *New WID on NR Sidelink Enhancement*, Dec. 2019.

[181] G. . V16.4.0, *TSG RAN; NR; Multiplexing and Channel Coding*, Jan. 2021.

[182] G. . V16.4.0, *TSG RAN; NR; Physical Layer Procedures for Control*, Jan. 2021.

[183] G. . V16.4.0, *TSG RAN; NR; Physical Layer Measurements*, Jan. 2021.

[184] N. Patriciello, S. Lagen, B. Bojovic, and L. Giupponi, "An E2E simulator for 5G NR networks," *Simulation Modelling Practice and Theory*, vol. 96, p. 101933, 2019. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1569190X19300589

[185] G. . V16.3.0, *TSG RAN; NR; Overall description; Stage 2*, Sep. 2020.

[186] S. Lagen, K. Wanuga, H. Elkotby, S. Goyal, N. Patriciello, and L. Giupponi, "New Radio Physical Layer Abstraction for System-Level Simulations of 5G Networks," in *IEEE International Conference on Communications*, June 2020.

[187] 3GPP TR 38.802 V14.2.0, *Study on New Radio (NR) Access Technology; Physical Layer Aspects*, Sep. 2017.

[188] 3GPP TR 37.885 V15.3.0, *Study on Evaluation Methodology of New Vehicle-to-Everything (V2X) Use Cases for LTE and NR (Rel. 15)*, Jun. 2019.

[189] G. . v16.2.0, *TSG RAN; NR; Radio Link Control (RLC) protocol specification*, Jan. 2021.

[190] 3GPP TR 38.901 V15.0.0, *Study on Channel Model for Frequencies from 0.5 to 100 GHz*, Jun. 2019.

[191] A. Bazzi, G. Cecchini, A. Zanella, and B. M. Masini, "Study of the impact of phy and mac parameters in 3gpp c-v2v mode 4," *IEEE Access*, vol. 6, pp. 71 685–71 698, 2018.

[192] Z. Ali, S. Lagén, L. Giupponi, and R. Rouil, "3gpp nr v2x mode 2: Overview, models and system-level evaluation," *IEEE Access*, vol. 9, pp. 89 554–89 579, 2021.

[193] T. Kim, G. Noh, J. Kim, H. Chung, and I. Kim, "Enhanced resource allocation method for 5g v2x communications," in *2021 International Conference on Information and Communication Technology Convergence (ICTC)*, 2021, pp. 621–623.

[194] D. M. Soleymani, L. Ravichandran, M. R. Gholami, G. Del Galdo, and M. Harounabadi, "Energy-efficient autonomous resource selection for power-saving users in nr v2x," in *2021 IEEE 32nd Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, 2021, pp. 972–978.

[195] M. Zhang, Y. Dou, P. H. J. Chong, H. C. B. Chan, and B.-C. Seet, "Fuzzy logic-based resource allocation algorithm for v2x communications in 5g cellular networks," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2501–2513, 2021.

[196] M. Lauridsen, D. Laselva, F. Frederiksen, and J. Kaikkonen, "5G New Radio User Equipment Power Modeling and Potential Energy Savings," in *2019 IEEE 90th Vehicular Technology Conference (VTC2019-Fall)*, 2019, pp. 1–6.

[197] 3GPP TR 38.840 V16.0.0, *Study on User Equipment Power Saving in NR*, June. 2019.

[198] "The Engineering ToolBox," //https://www.engineeringtoolbox.com/vehicle-flow-density-highway-design-d_1831.html/, accessed: 2021-12-16.

[199] "Mamdani fuzzy inference system," https://se.mathworks.com/help/fuzzy/mamfis.html/, accessed: 2022-06-27.

[200] "Defuzzification Methods," https://se.mathworks.com/help/fuzzy/defuzzification-methods.html/, accessed: 2022-06-27.

[201] Ericsson, "Hexa-X: 6G technology and its evolution so far," https://www.ericsson.com/en/blog/2021/7/hexa-x-6g-technology-6g-use-cases, accessed: 2023-06-10.

[202] H. Yin, P. Liu, K. Liu, L. Cao, L. Zhang, Y. Gao, and X. Hei, "Ns3-ai: Fostering artificial intelligence algorithms for networking research," in *Proceedings of the 2020 Workshop on Ns-3*, ser. WNS3 2020, 2020, p. 57–64.