





Project: H2020-ICT-2014-2 5G NORMA

Project Name:

5G Novel Radio Multiservice adaptive network Architecture (5G NORMA)

Deliverable D4.2 RAN architecture components – final report

 Date of delivery:
 30/06/2017

 Start date of Project:
 01/07/2015

Version: 1.0 Duration: 30 months **Document properties:**

Document Number:	H2020-ICT-2014-2 5G NORMA/D4.2
Document Title:	RAN architecture components – final report
Editor(s):	Mark Doll
Authors:	Hajo Bakker, Mark Doll, Diomidis Michalopoulos, Vinh Van Phan, Peter Rost, Peter Schneider (Nokia); Vincenzo Sciancalepore (NEC); Jorge Rivas (Atos); Alessandro Colazzo (Azcom); Sina Khatibi, Kunjan Shah (Nomor); Vasilis Friderikos, Oliver Holland (King's College London); Bin Han, Shreya Tayade (TU Kaiserslautern); Albert Banchs, Maria Cristina Márquez Colás (Universidad Carlos III de Madrid)
Contractual Date of Delivery:	30/06/2017
Dissemination level:	Public
Status:	Final
Version:	1.0
File Name:	5G NORMA D4.2

Revision History

Revision	Date	Issued by	Description	
1.0	21.07.2017	5G NORMA WP4	Final version	

Abstract

The main goal of 5G NORMA is to propose a multi-tenant multi-service mobile network architecture that adapts the use of the mobile network resources to the service requirements, the variations of the traffic demands over time and location, and the network topology, individually and concurrently for multiple tenants that share the infrastructure. This is the final deliverable of WP4 containing the findings of WP4 after the third 5G NORMA design iteration. Three options for RAN slicing, namely slice-specific RAN, splice-specific radio bearer and slice-aware RAN, are presented. Standardization relevance and potential of 5G NORMA innovations are discussed and a comparison with respective aspects of the 3GPP next generation architecture is made. Then, we show how the functionally decomposed c/d-layer can be adapted to specific services through suitable function selection and placement as well as how to use the functionally decomposed c/d-layer to realize a multi-service radio access. The c/d-layer function blocks are characterized and categorised into data layer, distributed and centralized control, respectively. Besides these architectural aspects, we introduce specific solutions for multi-tenancy RRM and admission control, for multi-connectivity support to increase reliability, integrate mm-wave and realize virtual cells, for signalling optimizations of mMTC and finally discuss the usefulness of geolocation DBs for mobile networks. Security considerations for both architecture design and specific solutions complete the 5G NORMA view on flexible RAN design.

Keywords

5G, mobile radio network architecture, RAN architecture components, functional decomposition, multi-tenancy, multi-service, multi-connectivity, NFV, network slicing, RAN slicing, physical network functions, SDMC

Table of Contents

1	Introduction	. 14
1.1	Scientific highlights	. 14
1.2	Outline and relation to other work packages	. 15
2	Flevible network design	17
21	RAN Slicing	17
$\frac{2.1}{2.1}$	1 Slice-specific RAN (Option 1)	19
2.1	 Slice-specific radio bearer (Option 2) 	20
2.1.	3 Slice-aware shared RAN (Option 2)	. 20
2.1.	4 Exemplary use of RAN Slicing	23
2.1.	Standardization relevance and potential	. 23
2.2	Use of functional decomposition for service adoptation	20
3 3 1	Eulering selection and placement	· 49
2.1	Multi compies radio access	. 27
3.2	Socurity considerations	. 54
5.5	Security considerations	. 34
4	Control and data layer	. 36
4.1	Centralized control	. 36
4.2	Distributed control	. 40
4.3	Data layer	. 41
4.4	Security considerations	. 43
4.4.	1 Securing inter-domain interfaces	. 43
4.4.	2 Securing intra-domain interfaces	. 45
5	Multi-tenancy	. 47
5.1	Multi-tenancy aspects of RRM	.47
5.1.	1 Multi-tenancy radio-resource management	.47
5.1.	2 Reinforcement learning for network slice resource management	.49
5.2	Multi-tenancy in multi-RAT environments	.54
5.2	1 Admission control	. 54
5.2	2 Dynamic resource sharing	. 56
5.3	Security considerations	. 58
5.3.	1 Isolation between multiple tenants	. 58
5.3.	2 Trust relationships	. 59
6	Multi-technology architecture in HetNets	60
61	Multi-connectivity algorithm	60
61	1 The Inter-RAT Link Controller	60
6.1	 Slice specific constraints 	60
6.1	3 Proposed algorithm	62
6.2	Clustering of mm-wave access points controlled by a 5G low band coverage layer	63
63	Virtual cells and multi-cell coordination	65
63	1 Novel virtual cell algorithm	65
63	 Provide virtual cell using O-learning 	69
63	3 Implementation of virtual cells in Demo 1	71
64	User-centric connection area	.71
65	Mobile edge computing	75
6.6	Massive machine-type communication RAN congestion control	.,, 77
6.7	Geolocation database	. 81
	Conducions	07
/	COLICIUSIONS	. 03
8	Annex A: Why no virtualized connectivity on the air interface?	. 85

9	Annex B: List of function blocks	86
9.1	Data layer	86
9.2	Distributed control	86
9.3	Centralized control	87
10	References	89
11	Multi-tenancy in multi-RAT environments	94
11	Admission control	
11.2	2 Dvnamic resource sharing	95
10		00
12	Inter-slice resource sharing	 98
12.	Poleted work	98
12.2	A mathematical programming formulation approach	99
12	1 Evaluation	101
12.	5 Concluding remarks	101
12.0	6 References	103
10	N - 14'- 1	105
13 12	Multiple connectivity at the different layers	105
13.	Δrebitecture and process	105
15.2	Architecture and process	105
14	Data-layer and control-layer design for multi-connectivity	108
14.	Motivation and problem statement	108
14.2	2 Related work	109
14.	Architectural approaches	110
14.4	Finoughput evaluation of multi-connectivity	111
14.		112
15	Architectural approaches for multi-connectivity of mm-wave APs and 5G low band	112
15	Motivation and problem statement	113
15.	Major results	113
15.3	Architectural approaches for provisioning of 5G-mmAPs	113
15.4	Evaluation	116
15.5	5 Security considerations	118
15.0	6 References	118
16	Virtual calls and multi-call coordination	110
10 16	Virtual Cens and Indid-Cen Cool dination	119
16.	Numeric results	119
16.2	The implementation of the Virtual Cell in Demo 1	121
16.4	4 References	123
17	Elemitele 5C commiss flow (CE) with in SE Och differentiation and multi compositivity	
1/	Flexible 5G service-flow (SF) with in-SF QoS differentiation and multi-connectivity	124
17.	PDCP multiplexing of SFs from different network slices	147
1 / • ·		124
17.2	2 Coordinated dynamic scheduling and semi-persistent scheduling based allocations	124
17.2	2 Coordinated dynamic scheduling and semi-persistent scheduling based allocations	124 129
17.2 17.3	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References. 	124 129 132
17.2 17.3 18	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References. User-centric connection area 	124 129 132 133
17.2 17.3 18 18.1	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References	124 129 132 133 133
17.2 17.3 18 18.1 18.2	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References User-centric connection area Motivation and problem statement Major results 	124 129 132 133 133
17.2 17.3 18 18.2 18.2	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References. User-centric connection area Motivation and problem statement Major results Related work 	124 129 132 133 133 133 134
17.2 17.3 18 18.2 18.2 18.2	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References User-centric connection area Motivation and problem statement Major results Related work UCA concept 	124 129 132 133 133 133 134 134
17.2 17.3 18.1 18.2 18.2 18.2 18.4	 Coordinated dynamic scheduling and semi-persistent scheduling based allocations References	124 129 132 133 133 133 134 134 135

18.7	Security considerations	140
18.8	References	140
19 Mas	ssive machine-type communication RAN congestion control	142
19.1	Impact analysis of D2D link exception	142
19.2	Enhanced grouping processes	142
19.3	Enhanced transmission frame structure	145
19.4	Geolocation database integration	146
19.5	Security considerations	147
19.6	Summary	148
17.0	Summary	1 10
20 Geo	location databases, use of geolocation information and associated opportunit	ties
20 Geo	location databases, use of geolocation information and associated opportunit	ties
20 Geo 20.1	location databases, use of geolocation information and associated opportunit Motivation and problem statement	ties 149 149
20 Geo 20.1 20.2	location databases, use of geolocation information and associated opportunit Motivation and problem statement	ties 149 149 150
20 Geo 20.1 20.2 20.3	location databases, use of geolocation information and associated opportunit Motivation and problem statement Major results Related work	ties 149 150 152
20 Geo 20.1 20.2 20.3 20.4	location databases, use of geolocation information and associated opportunit Motivation and problem statement Major results Related work Signalling procedures	ties 149 150 152 153
20 Geo 20.1 20.2 20.3 20.4 20.5	location databases, use of geolocation information and associated opportunit Motivation and problem statement Major results Related work Signalling procedures Security considerations	ties 149 149 150 152 153 157

List of Figures

Figure 2-1: Three RAN slicing options
Figure 2-2: Functional control and data layer architecture, RAN slicing Option 1 (slice-specific RAN)
Figure 2-3: Functional control and data layer architecture, RAN slicing Option 2 (slice-specific radio bearer) (function types)
Figure 2-4: Integration of RAN slicing Option 2 and multi-connectivity based on MAC and PDCP layer, respectively
Figure 2-5: Functional control and data layer architecture, RAN slicing Option 3 (shared RAN)
Figure 2-6: Exemplary deployment using RAN slicing in the context of an industrial campus deployment
Figure 2-7: Overall architecture of 3GPP 5G network (Figure 4.1-1 of TS 38.300)
Figure 2-8: Overview of functional split between 3GPP NG-RAN and 5GC (Figure 4.2-1 of TS 38.300)
Figure 2-9: The gNB architecture with CU and DUs (Figure 11.1.3.8-1 of TR 38.801)26
Figure 2-10: Functional split options between CU and DU (Figure 11.1.1-1 of TR 38.801) 26
Figure 3-1: Function selection and placement (RAN slicing Option 2)
Figure 3-2: Multi-service vs. multi-tenancy and inter-tenant multi-connectivity vs (single tenant) multi-connectivity
Figure 3-3: Example 5G New Radio multi-service RAT for RAN slicing Option 2
Figure 4-1: 5G NORMA SDMC interfaces
Figure 4-2: Abstraction layer turning non-SDN devices into SDN-controllable devices
Figure 5-1: Reinforcement learning for slice admission control
Figure 5-2: Network slice admission control
Figure 5-3: Pseudocode of the admission control algorithm
Figure 5-4: Revenue vs ρ_i/ρ_e
Figure 5-5: Utility gains for different approaches as a function of network size
Figure 6-1: Slice specific requirements and UE mapping
Figure 6-2: Data duplication mode (left) and data split mode (right)
Figure 6-3: Flow diagram of the proposed algorithm
Figure 6-4: Deployment of a UE specific multi-connectivity cluster of mmAPs on a 5G-LB coverage layer
Figure 6-5: Activity Management within MC-Cluster, triggering of data transmission with flags
Figure 6-6: An example of the virtual cell concept (based on [CSS+16])
Figure 6-7: The network throughput increases as a function of the cell edge threshold (Scenario A)
Figure 6-8: The total network throughput and the three phases of the Q-learning

Figure 6-9: The key elements of Demo 1 of WP671
Figure 6-10: Signalling messages towards the core based on 4G LTE72
Figure 6-11: Reduction of signalling messages towards the core with the NORMA UCA concept
Figure 6-12: Gain over LTE for different RRC timer values and different UE speeds
Figure 6-13: Gain over LTE for different paging are sizes
Figure 6-14: Mobile Edge Computing platform supporting Slice QoS Monitoring
Figure 6-15: Message sequence chart of MEC slice QoS monitoring/enforcement
Figure 6-16: Different congestions and their sources in LTE-A networks
Figure 6-17: Comparing the performances of different RAN congestion controlling methods 78
Figure 6-18: The RAN topology in D2D-based grouped RA
Figure 6-19: Collision densities with respect to RACH resource dedication in a 2-DC-case (DC1+DC2), numerical results obtained from 500 iterations of Monte-Carlo test
Figure 6-20: Performances of different allocation methods, when DC1 is required to have an average collision rate of 0.02
Figure 11-1. The distribution of the revenues obtained by random smart policies compared to the proposed algorithms
Figure 11-2. Revenue in perturbed scenario, $\rho_i / \rho_e = 5$
Figure 11-3: Normalised utility gain G _w as a function of the maximum allowed number of handovers m
Figure 11-4: Computational complexity of GLLG and SoA algorithms
Figure 11-5: Capacity saving96
Figure 11-6: Improvement on the user throughput97
Figure 12-1: Comparison between loose and tight coupling using a toy example scenario 99
Figure 12-2: Aggregate rate vs. number of user
Figure 12-3: CDF of throughput per user
Figure 13-1: C/d-layer architecture for multi-connectivity support in a multi RAT environment
Figure 13-2: Message sequence chart for the inter-RAT link selection process 107
Figure 14-1 The LTE RAN Architecture
Figure 14-2 The 5G RAN Architecture considered in [5GN-D41] vs the existing LTE RAN Architecture
Figure 14-3 Transport connectivity options to small cells [NGMN-SBH] 110
Figure 14-4: The small-cell aggregation and centralized scenario architectural approach 110
Figure 14-5: Throughput performance (at application layer) of approach 1 112
Figure 14-6: Throughput performance (at application layer) of approach 2 112
Figure 15-1: PDCP based MC with two variants for location of RRC protocol stack
Figure 15-2: Architecture approach -2: split in PDCP_H and PDCP_L and related protocol stack
Figure 15-3: Architecture approach-3: MAC based MC and related protocol stack 115

Figure 15-4: Signalling messages flow between the node and the UE for approach-1a 116
Figure 15-5: Comparison of number of messages between different approaches
Figure 16-1: The network throughput increases as a function of the cell edge threshold (Scenario B)
Figure 16-2: The network throughput increases as a function of the cell edge threshold (Scenario C)
Figure 16-3: The throughput of the cells in download direction
Figure 16-4: The throughput of the cells in download direction
Figure 16-5: Network reconfiguration into edge-cloud
Figure 16-6: Software demo with Edge-Cloud Configuration
Figure 16-7: The beam coordination method
Figure 17-1: PDCP structure for PDCP multiplexing of SFs from different NSIs into one RB125
Figure 17-2: PDCP multiplexing of SFs from different NSIs into one RB 125
Figure 17-3: Adding a NSI-SF into a multiplexed RB
Figure 17-4: Illustration of SN mapping inquiry procedure of PDCP 127
Figure 17-5: Possible extension to cover NSI specific c-layer RRC, referred to as RRC_Slice128
Figure 17-6: TB formed based on the coordination between SPS and dynamic scheduling 130
Figure 17-7: Coordinated dynamic scheduling and SPS resource allocations
Figure 18-1: (a) Assignment and (b) update of UCA for a UE traversing the path shown by the dashed line
Figure 18-2: Saving of signalling messages towards the core with UCA compared to 4G LTE
Figure 18-3: Idle/active modelling for simulation study based on SDP arrivals and UE mobility
Figure 18-4: Definition of a UCA based on applying a window size to RSRP based UE measurements
Figure 18-5: UCA size and UCA updates with respect to UCA window size
Figure 18-6: Gain over LTE for different UE speeds with linear mobility
Figure 18-7: Gain over LTE for different UE speeds and different mobility schemes
Figure 18-8: Gain over LTE for different packet inter-arrival times
Figure 18-9: Gain over LTE for different RRC timer values and different UE speeds
Figure 18-10: Gain over LTE for different paging are sizes
Figure 19-1: The simulated impact of D2D link failure reports through extra RA processes 142
Figure 19-2: The global group updating process
Figure 19-3: The group leaving process, triggered by detaching events
Figure 19-4: The group leaving process, triggered by a GM losing D2D link but holding macro- cell link
Figure 19-5: The group leaving process, triggered by a GM losing D2D and macro-cell links 144
Figure 19-6: The group leaving process, triggered by handover events

Figure 19-7: Frame structure of the grouped mMTC transmission based on D2D; the terms R/R, UDT, A/C and DDT denote request/report, uplink data, acknowledgment/command, respectively
Figure 19-8: GDB signalling (generic case)
Figure 20-1: Virtualized edge network slices achieving a more direct path compared with (fixed) network elements in a Tactile Internet remote surgical operation example
Figure 20-2: Map of channel availability for a Class 3 white space device, for a large area of England, for a 30m height above ground level transmitter, \geq 30 dBm allowed power
Figure 20-3: CCDFs of channel availability for different classes of white space device, for a large area of England, for a 30 m height above ground level transmitter,≥30 dBm allowed power scenario
Figure 20-4: Example of latency CDFs for optimised (using the GDB) and non-optimised Network Function placement
Figure 20-5: Signalling—resource management (generalisation)154
Figure 20-6: Signalling—rendezvous
Figure 20-7: Signalling—indirect access
Figure 20-8: Signalling—Link to mMTC congestion control, MAC structuring
Figure 20-9: Signalling—Link to mMTC congestion control reflected in GDB signalling (partial)

List of Tables

Table 2-1: 5G NORMA innovations vs. 3GPP 5G, grouped according to the 5 key innovations 27
Table 4-1: Subset of parameters obtained for eMBMS Control
Table 5-1: Traffic Class Requirements (similar QCI from [23.203]) 49
Table 6-1: Virtual cell scenario configurations 67
Table 6-2: Process of MEC triggering MAC Scheduling to adapt scheduling
Table 6-3: Simulation specification
Table 6-4: Results in simulations under different DC combinations
Table 12-1: Simulation parameters used for numerical investigations 101
Table 18-1: Simulation parameter for UCA performance assessment
Table 20-1: Statistical results on number of allowed channels 151
Table 20-2: Scenarios (MBD: Mobile Broadband Downlink; IBP/MBU: Internal Broadband Provisioning/Mobile Broadband Uplink)

List of Acronyms and Abbreviations

3GPP	3rd Generation Partnership Project	ETSI	European Telecommunication Standard
5GC	5G Core		Institute
A/D	Analog-to-Digital Converter	EU	European Union
ABSF	Almost-Blank Sub-Frame	E-UTRA	Evolved Universal Terrestrial Radio Access
АСК	Acknowledgement	E-UTRAN	Evolved Universal Terrestrial Radio Access
AI	Air Interface		Network
AMCC	Activity Management within a Multi-	FCC	Federal Communications Commission
	Connectivity cluster	FDD	Frequency Division Duplexing
AMF	Access and Mobility Management Function	FE	Function Element
AMMC	Activity Management within a Multi-	FEC	Forward Error Correction
	Connectivity cluster	FFT	Fast Fourier Transform
ANDSF	Access Network Discovery and Selection	ForCES	Forwarding and Control Element Separation
	Function	FPTAS	Fully Polynomial Time Approximation
AP	Access Point		Schemes
API	Application Programming Interface	GBR	Guaranteed Bit Rate
ARP	Allocation and Retention Priority	GC	Group Coordinator
ARQ	Automatic Repeat Request	GDB	Geolocation Database
AS	Access Stratum	GLLG	Greedy Local Largest Gain
	Autonomous Timing Advance	GIVI	Group Member
BCCH	Broadcast Control Channel	ginb ctr	Next Generation Node B
	Carrier Aggregation	GIP	GPRS Tunnelling Protocol
CAIVIADIVI	Context-Aware Multiple Attribute Decision	GWCN	Gateway Core Network
CADEV	Making Consisted Funder distance	HZUZU	Horizon 2020
	Capital Experioritore	HARQ	Hybrid Automatic Repeat Request
CDRS	Citizens Brodubanu Radio Service	Heinei	
CCIVICC	Connectivity Cluster		High Speed Dacket Access
CD	Confidential Degree		High Speed Packet Access
	Channel		Hardware
	Context Identification		Inter Arrival Time
	Command Line Interface		Inter-Cell Interference Cancelation
CN	Core Network		Information and Communication
CoMP	Coordinated Multipoint Transmission and		Technologies
com	Recention	IFTF	Internet Engineering Task Force
CP	Control Plane	iFFT	Inverse East Fourier Transformation
CPC	Cognitive Pilot Channel	IoT	Internet of Things
CPRI	Common Public Radio Interface	IP	Internet Protocol
CQI	Channel Quality Indication	IT	Information Technology
CRC	Cyclic Redundancy Check	JT/JP	(Downlink) Joint Transmission/(Uplink) Joint
CRE	Cell Range Expansion		Processing
CRS	Common Reference Signal	LC	Logical Channel
CRS	Cell-specific Reference Signal	LMMC	Link Management within MC-Cluster
CSI	Channel State Information	LMMC	Link Management within a Multi-Connectivity
CSI-RS	Channel State Information Reference Signal		cluster
D/A	Digital-to-Analog Converter	LNA	Low-Noise Amplifier
D2D	Device-to-Device	LSA	Licensed Shared Access
DC	Device Class	LTE	Long-Term Evolution
DCCH	Dedicated Control Channel	LTE-A	Long Term Evolution Advanced
DCN	Data Communication Network	MAC	Medium Access Control
DG	Distributed Greedy	MBMS	Multimedia Broadcast Multicast Services
DL	Downlink	MBR	Maximum Bit Rate
DM-RS	Demodulation Reference Signal	MBSFN	Multimedia Broadcast Single Frequency
DPF	Direct Provisioning Function		network
DRB	Data Radio Bearer	MC	Multi-Connectivity
DRX	Discontinuous Reception	MCCH	Multicast Control Channel
DS	Dynamic Scheduling	MCH	Multicast Channel
E2E	End-to-End	MCN	Multi-hop Cellular Network
EC	European Commission	MCPTT	Mission Critical Push to Talk
eDECOR	Enhancements of Dedicated Core Networks	MCS	Modulation and Coding Scheme
eIMTA	Enhanced Interference Mitigation and Traffic	MCS	Modulation and Coding Scheme
	Adaptation	MEC	Mobile Edge Computing
eMBMS	Evolved Multimedia Broadcast Multicast	MeNB	Master Evolved Node B
	Services	MIB	Master Information Block
eMBSFN	Evolved Multimedia Broadcast Single	MIMO	Multiple-Input Multiple-Output
500	Frequency network	MMC	Massive Machine Communication
EPC	Evolved Packet Core	MIME	Nobility Management Entity
EPS	Evolved Packet System	MMSE	IVIINIMUM Mean Square Error

mMTC	Massive Machine-Type Communication	RRH	Remote Radio Head
MN	Moving Network	RRM	Radio Resource Management
MNO	Mobile Network Operator	RRU	Radio Resource Utilisation
MOCN	Multi-Operator Core Network	RSRP	Reference Signal Received Power
MORA	Multi-Operator Resource Sharing	RTT	Round Trip Time
MPL	Mean Path Length	SAE	System Architecture Evolution
MPTCP	Multi-Path Transmission Control Protocol	SBI	Southbound Interface
MSP	Mobile Service Provider	SCTP	Stream Control Transmission Protocol
MTA	Moving Tracking Area	SDMC	Software-Defined Mobile Network Control
MTC	Machine Type Communication	SDM-C	Software-Defined Mobile Network Controller
MTCD	Machine-Type Communication Device	SDMN	Software-Defined Mobile Network
МТСН	Multicast Traffic Channel	SDM-O	Software-Defined Mobile Network
MUX	Multiplexer		Orchestrator
MVNO	Mobile Virtual Network Operators	SDM-X	Software-Defined Mobile Network
NAS	Non-Access Stratum		Coordinator
NBI	Northbound Interface	SDP	Short Data Packet
NETCONF	Network Configuration Protocol	SDU	Service Data Unit
NFV	Network Function Virtualisation	SeNB	Secondary Evolved Node B
NG	Next Generation	SF	Service Flow
NG-RAN	Next Generation Radio Access Network	SFN	System Frame Number
NR-RAN	New Radio Radio Access Network	S-GW	Serving Gateway
NSI	Network Slice Instance	SIB	System Information Block
NSI-SF	Network Slice Instance Service Flow	SINR	Signal-to-Interference-plus-Noise Ratio
NVS	Network Virtualisation Substrate	SLA	Service Level Agreement
OFDM	Orthogonal Frequency-Division Multiplexing	SON	Self-Organised Network
OFDMA	Orthogonal Frequency-Division Multiple	SPS	Semi-Persistent Scheduling
0.115	Access	SRB	Signalling Radio Bearer
ONF	Open Networking Foundation	SSH	Secure Snell
OPEX	Operational Expenditure	SSL	Secure Sockets Layer
			Soltware
	Power Ampiller		Transmission Control Protocol
	Paging Control Channel		Time Division Dupley
	Protocol Data Unit		Time Division Duplex
P-GW	Packet Data Network Gateway		Transport Laver Security
PHV	Physical Laver	тр	Transmission Point
РНУ ТР	Physical Laver Transmission Point	 тті	Transmission Time Interval
PRACH	Physical Random Access Channel	TV	Television
PRB	Physical Resource Block	TVWS	Television White Spaces
ProSe	Proximity Services	UCA	User-Centric Connection Area
PS	Public Safety	UDN	Ultra-Dense Network
PSM	Power Saving Mode	UE	User Equipment
PSME	Program Making and Special Events	UF-UFDM	Universal Filtered Orthogonal Frequency-
QCI	QoS Class Identifier		Division Multiplex
QoE	Quality of Experience	UL	Uplink
QoS	Quality of Service	UMTS	Universal Mobile Telecommunications System
QPS	QoS Parameter Set	UP	User Plane
RA	Random Access	UPF	User Plane Function
RACH	Random Access Channel	URC	Ultra-Reliable Communication
RAN	Radio Access Network	URLLC	Ultra-Reliable Low Latency Communication
RAO	Random Access Opportunity	V2X	Vehicular-to-everything
RAT	Radio Access Technology	VC	Virtual Cell
RB	Radio Bearer	VNF	Virtual Network Function
REST	Representational State Transfer	VoIP	Voice over IP
RLC	Radio Link Control	WP	Work Package
ROHC	Robust Header Compression	YANG	Yet Another Next Generation (modelling
RRC	Radio Resource Control		language)

PART I: ARCHITECTURE

1 Introduction

5G NORMA's main objective for its novel mobile network architecture is to enable the integration of different technologies and different use cases, concurrently for multiple tenants on a shared infrastructure. Different use cases respectively services create different requirements, which calls for different service-individual realizations. Together with the desire to separate different tenants' network resources according to contractual agreements, this necessitates to use the right functionality at the right place and time within the network and to assign resources to the right tenants.

In order to provide this flexibility, the network function virtualisation (NFV) paradigm is adopted in the mobile access and core network domain. Mobile network functionality is decomposed into smaller function blocks and flexibly instantiated and assigned to the tenants' network slices. This is complemented by the centralized control (software-defined networking, SDN) paradigm, broadened by 5G NORMA to become the comprehensive *software-defined mobile network control* (SDMC) concept, covering all aspects of mobile networks instead of just transport. The network programmability offered through SDMC enables mobile service providers to flexibly control and manage their networks, and tenants, through an exposed application programming interface (APIs), to customize their slices' behaviour to their needs.

5G NORMA proposed 5 key innovations, more precisely 3 innovative enablers and 2 innovative functionalities, that comprise such an ambitious mobile network architecture. While these 5 key innovations acted as general design guideline to all the work in 5G NORMA, each of the 5 key innovations is put specifically into focus in a different Chapter:

The "3" Innovative Enablers

- Adaptive (de)composition and allocation of mobile network functions between the edge and the network cloud depending on the service requirements and deployment needs is tackled with Chapter 3 *Use of functional decomposition for service adaptation* that presents ways of function selection and placement.
- **Software-Defined Mobile network Control (SDMC)**, which applies the SDN principles to mobile network specific functions, is described in Chapter 4 *Control and data layer* specifically in conjunction with centralized control.
- Joint optimization of mobile access and core network functions localized together in the network cloud or the edge cloud, is an important aspect of many of the novel functionalities presented in chapter 6 *Multi-technology architecture in HetNets*.

The "2" Innovative Functionalities

- **Multi-service- and context-aware adaptation of network functions** is covered again by Chapter 3 *Use of functional decomposition for service adaptation* with the introduction to multi-service radio access.
- **Mobile network multi-tenancy** is introduced in Chapter 2 *Flexible network design* with radio access network (RAN) slicing; related questions with respect to radio resource management and admission control are thoroughly examined in Chapter 5 *Multi-tenancy*.

1.1 Scientific highlights

Multi-tenancy support is the heart of 5G NORMA in general and of WP4 in special, where the focus lies on its realization within the RAN. A number of notable outcomes of WP4 form the basis for multi-tenancy as well as proof of its usefulness, to highlight some:

- A functionally decomposed control and data layer architecture to replace today's network of entities by a network of functions;
- RAN slicing to extend network slicing into the RAN for a tenant-individual customization of the RAN; and

• The proof that dynamic sharing of resources among multiple tenants never performs worse than static sharing.

The functionally decomposed control and data layer creates a network of functions. It replaces today's network of entities, in which each such entity is responsible for a pre-assigned set of mobile network functions. With5G NORMA, the only entities that exist are different kinds of cloud processing infrastructure, all interconnected by virtualized (SDN-enabled) transport networks. Only one non-virtualized physical layer function (PNF) must remain at the antenna sites, the PHY TP function block. It implements the very lowest level of the air interface, functions that inherently cannot be virtualized as they are non-digital (analogue and mixed signal processing). As presented in Section 3.1, the flexibility introduced by such a functionally decomposed control and data layer enables slice-individual and therefore tenant-individual service adaptation. Furthermore, as shown in Section 3.2, it can be used to realize a multi-service-capable RAN in a very flexible and future-proof way. Multi-service support within a single RAN instance, i.e. a single base station respectively cell, is a crucial functionality if the RAN is to be shared by multiple tenants with diverse and likely contradicting service requirements to the air interface.

Both tenant-individual service adaptation and tenant-aware multi-service support on common RAN infrastructure culminate into the unique functionality of RAN slicing, which is presented in Section 2.1. Up to now, tenant-individual customization is limited to the fully virtualized part of the network, where it is provided through network slicing. RAN slicing effectively extends network slicing into the RAN, allowing tenant-individual customization of the air interface. This includes, depending on the slicing option chosen, many of its non-virtualized physical network functions. Through RAN slicing, network slicing in 5G NORMA effectively becomes end-to-end, spanning the whole path from the data network (respectively end-user service) to the UE.

To proof the benefits of multi-tenancy, in Section 5.2.2, a criterion for dynamic resource allocation amongst tenants is proposed. One of the most important key properties provided is represented by the utility gain obtained over the static sharing (SS) approach. In the latter strategy, each tenant contracts for a fixed slice/fraction of the network resources at each base station for its exclusive use. In this scenario, each operator independently optimizes its users' associations and allocation of resources in order to maximize its utility. Under the Multi-Operator Resource Sharing (MORA) criterion, that jointly optimizes operators' users' association and resource allocation, the overall network utility is clearly larger than that under SS. Furthermore, we proved that for a given user association x, MORA's resource allocation achieves always a higher or equal utility than that of SS [CBV+16].

1.2 Outline and relation to other work packages

This document is structured as follows: This Chapter 1 introduces the overall objective of the WP4 work, highlighted three of its main outcomes and outlines the structure of this document as well as its relation to the other 5G NORMA work packages. The following Chapters 2, 3 and 4 cover general architectural aspects, before Chapter 5 and Chapter 6 dive into specific WP4 innovations and present representative evaluation results. Further background information and evaluation results to the specific WP4 innovations can be found in Part II. Throughout the document, we provided security considerations, which cover most of the architectural aspects and selected innovations with possible security implications.

In detail, Chapter 2 presents the control and data layer functional architecture in conjunction with three options for *RAN slicing*, namely slice-specific RAN, splice-specific radio bearer and slice-aware RAN, to exemplify the support of multi-tenancy by the 5G NORMA architecture. Then, Chapter 3 shows, by example of three prominent examples broadband, low latency and mission critical service, how the functionally decomposed c/d-layer is utilized to adapt a slice to specific service requirements through suitable *function selection and placement*. Second, for the common network functions that are shared by multiple slices, it shows how the functionally decomposed

c/d-layer helps to realize a *multi-service radio access* in a flexible and future-proof way (Annex A motivates why at all to service-specifically tailor transport over the air interface). In Chapter 4, the c/d-layer function blocks (a complete list of function blocks can be found in Annex B) are characterized and categorised into *centralized control*, *distributed control* and *data layer*, respectively, and specific aspects of each class are discussed: The SDMC concept is explained together with centralized control, why and where to deviate from the SDMC concept is motivated in conjunction with decentralized control and the limits of virtualization, which lead to the introduction of physical network functions (PNF), are discussed with the data layer.

Coming to the specific WP4 innovations, Chapter 5 presents solutions for *multi-tenancy radio resource management* (RRM) and *admission control*. Chapter 6 then presents several benefits of *multi-connectivity*, namely specific solutions to *increase reliability*, to *integrate mm-wave* and for realizing *virtual cells*, which increase efficiency and service quality in time-division duplex (TDD) networks. Further innovations bring core functions to the RAN, which has several benefits: The *user-centric connection area* (UCA) reduces mobility related signalling, *mobile edge computing* (MEC) reduces application latency and *reduced congestion in massive machine type communication* (mMTC) is achieved by clustering mMTC signalling with assistance of a central *geolocation database*. Security considerations for both architecture design and specific solutions complete the 5G NORMA view on flexible RAN design. Part I concludes with a summary of the WP4 outcomes in Chapter 7, Annexe with additional information on the general architecture design and references. Finally, Part II provides various additional background information to specific WP4 innovations including further evaluation results.

WP4 Flexible RAN Design only covers a subset of the overall 5G NORMA architecture and its related innovations, namely those of the RAN control and data layer. More specifically it focuses on the question how to provide flexibility through functional decomposition and how to use it to realize multi-tenancy and multi-service within a single RAN. In contrast, the main focus of WP5 Flexible Connectivity and QoS/QoE Management lies in shaping the SDMC concept, i.e. how control functionality can be implemented as SDM-C applications in a centralized way, specifically those for QoS/QoE and management [5GN-D52]. In WP3 Multi-service Network Architecture, the designs of WP4 and WP5 are integrated into a comprehensive control and data layer architecture and complemented by WP3 innovations on the management and orchestration (MANO) layer and the service layer. The interworking of WP4 with WP3 and WP5 is therefore primarily visible in Section 4.1 on centralized control, where the SDMC concept and specific SDM-C applications like QoS Control are presented and the connection to the MANO layer is made. Besides OoS Control, the User-centric Connection Area (UCA, cf. Section 6.4) relies on the SON application for determining a suitable UCA per UE, one of the functions discussed in WP5 connectivity management. Security considerations, which are provided throughout this document, are to be understood in relation to the primary work on security in 5G NORMA hosted by WP3 (Task 3.3 Security). From WP6 Demonstrator, which realizes selected innovations of WP4, first outcomes complement WP4 work in respect to implementation aspects, e.g. on southbound interface (SBI) protocol plugins for the SDMC controller (cf. Section 4.1) and with evaluation results, e.g. for virtual cells (cf. Section 6.3 and Chapter 16). Finally, WP4 results presented in this document contribute to the verification and evaluation of the overall 5G NORMA architecture design, which is carried out in WP3 [5GN-D33] and in WP2 Use cases, requirements and KPIs [5GN-D23], respectively.

2 Flexible network design

Traditionally, mobile networks implicitly group functions into **network entities** via specification of their interconnections respectively by absence of standardized interconnections between the grouped functions. Each entity is responsible for a pre-defined set of functions. Accordingly, the degrees of freedom for assigning network functionality to physical network entities are very limited. Replacing the traditional network of entities by a flexible "**network of functions**" allows for adapting the network to diverse services in a tailor-made way, by using different virtual network functions (VNF) rather than using just different parameterisations of a common VNF. Each block may be replaceable and individually instantiated for each logical network running on the same infrastructure. Depending on the use case, requirements, and the physical properties of the existing deployment, VNFs and (to some extent) physical network functions (PNF, cf. Section 4.3) are executed at different entities within the network. Coexistence of different use cases and services would imply the need for using different VNF allocations within the network as detailed in the Chapter 3.

The mobile network must further integrate also legacy technologies to guarantee that it can operate with existing networks. This is reflected in the following by using 3GPP EPS as the basis for the set of function blocks. 5G NORMA further adds blocks and amends existing blocks' functionalities to implement its aforementioned 5 key innovations. The involved NFs are described in Chapter 4.

This chapter focuses on RAN slicing as the central novel functionality and overarching structural element of 5G NORMA's flexible network design. The following section introduces three RAN slicing options, presents the according control and data layer functional architecture for each of them and shows the applicability to a set of current and future services and use cases. A discussion of standardization relevance and potential in reference to 3GPP concludes this chapter. Security considerations are provided with later chapters, when individual aspects of RAN slicing like function selection and placement, inter-domain interfaces or multi-tenant RRM are discussed.



2.1 RAN Slicing

Figure 2-1: Three RAN slicing options

In [5GN-D41], we introduced three options for RAN slicing¹ that are considered by 5G NORMA. Those three options differ by the degree of freedom offered for slice-individual customization, as well as the required complexity for implementation. In the following, all three options, jointly shown in Figure 2-1, are summarized; afterwards, the functional architecture for each RAN slicing option is detailed in the sub-sequent sections. All function blocks are described in detail later in Chapter 4 (also cf. <u>Annex B</u> and the List of Acronyms and Abbreviations above).

- Option 1: **Slice-specific RAN.** The first option in Figure 2-1 refers to the case where only transmission point specific functionality is shared among network slices while all other functionality is instantiated specifically for each network slice. In this case, the maximum degrees of freedom are achieved because each network slice may be customized down to the physical layer. On the other hand, this option requires a tight synchronization of the multi-tenancy policies applied to the common part and the perslice (dedicated) implementation, which may limit the achievable multiplexing gains considerably. Examples for this option include the possibility for implementing different radio access technologies within the same shared spectrum, e.g., 4G and 5G, or to separate two deployments while still exploiting multiplexing gains.
- Option 2: **Slice-specific radio bearer.** The second option in Figure 2-1 refers to sharing transmission point (cell) specific and user specific functionality, i.e., PHY and MAC in the data layer, and RRC in the control layer. This options slightly reduces the complexity because resource multiplexing would be implemented across all network slices and each network slice makes use of the same efficient flexible RAN implementation; on the other hand, each network slice may still customize the operation through configuration and parameterization based on the service requirements and each network slice may still implement its own QoS control (QoS prioritization). Hence, this option provides a reasonable trade-off of flexibility and complexity.
- Option 3: **Slice-aware shared RAN.** The third option in Figure 2-1 refers to a deployment where the complete RAN is shared by multiple tenants. This option is close to existing solutions such as eDECOR [23.711] although in our case, multi-slice connectivity is considered. Hence, one UE may be connected to more than one network slice. In addition, the SMDC Coordinator (SDM-X, cf. Section 4.1) is quite powerful in this option because a significant part of the RAN functionality may be implemented as applications on top of the SDM-X.

The above options are not meant to cover all possible RAN slicing options but they cover three setups of particular interest, each providing different benefits and requiring different degrees of complexity. In addition, the different options may co-exist, e.g., Option 1 may multiplex network slices customized down to the PHY layer with sets of network slices implemented using Option 2 and a set of Option 3 slices can be multiplexed either also with Option 1 or with Option 2, then sharing the lower part up to MAC with other Option 2 slices.

¹ Network slicing, in general, refers to sharing a common infrastructure between multiple logical network instances (see also [5GN-D32]). In this context, a tenant can be a mobile network operator, or companies from vertical industries requesting and using a network slice instance.

2.1.1 Slice-specific RAN (Option 1)



Figure 2-2: Functional control and data layer architecture, RAN slicing Option 1 (slicespecific RAN)

The functional architecture for the first option is shown in Figure 2-2. The depicted lowercase letters (m, r, c, e, n, s, p, g, q, t) indicate the exchange of information during operation between the respective function blocks. A letter at the bottom of one block relates to the same letter at top of other block(s), implying a controlling logic and controlled agent relation. The uppercase letters (C, U, M) represent proprietary control interfaces between distributed control and data layer (for the rationale and details cf. Section 4.2).

Transmission (and reception) point specific functionality is shared across network slices while all other functionality is implemented specifically by each network slice. This choice may coincide with the choice of functional split between the access point (AP) respectively more precisely the distributed unit (DU) located at the antenna site and the central unit (CU), i.e., the SDM-X functionality may be executed at the DU controlling the PNFs, while the slice-specific implementation would be implemented at the centralized processor using partly VNFs. Alternatively, the SDM-X may be more centralized controlling multiple DUs depending on the connectivity properties between SDM-X and PHY TP.

The implementation of Option 1, in particular the SDM-X, is very challenging because each slice may implement an own slice-specific scheduler. The resource management of these schedulers and the multi-tenancy policies enforced by the SDM-X need to be kept consistent, which requires closed-loop feedback between the SDM-X and per-slice schedulers. However, this may be alleviated by reserving fixed resources per slice, such as for legacy systems, which would also limit the multiplexing gain significantly.

Hence, the interface between common and dedicated part would mainly cover radio resource management information in order to allow for multiplexing the different network slices in the air interface, i.e. no post-processing of time/frequency domain symbols necessary. This covers also inter-cell coordination algorithms such as ICIC, which may be implemented in each slice but need to be coordinated with the SDM-X, e.g., ensuring similar resource allocation of the same network slice at different APs.

Furthermore, Figure 2-2 shows that RRC Cell is located in the common and dedicated part. This reflects the possibility of individual RRC Cell implementations in each network slice, while the set of configuration parameters pertaining to PHY TP (e.g. enabled antenna elements, total transmission) is under common control. Slice-specific RRC Cell may be necessary in the case of legacy systems that are not slicing-capable. In addition, the RRC User function block is implemented in each network slice, i.e., also user mobility would be implemented in a slice-specific way. Hence, if a user connects to multiple slices of Option 1, it must be able to support multiple RRC instances.

A key enabler for this option, and for the following RAN slicing options, is a flexible RAN numerology as investigated by H2020 FANTASTIC-5G, cf. Section 3 in [F5G-D32]. Such a flexible numerology allows for allocating radio resources and configuring their usage in a service-specific manner. Hence, each network slice representing a different service may use individual numerologies in order to adjust the air interface to its requirements. As such, using RAN slicing and this flexible numerology, very different networks, i.e. end-to-end logical networks, may be implemented in the same spectrum in order to optimize the radio resource usage.



2.1.2 Slice-specific radio bearer (Option 2)

Figure 2-3: Functional control and data layer architecture, RAN slicing Option 2 (slicespecific radio bearer) (function types)

The second RAN slicing option is illustrated in Figure 2-3 where both the transmission point and user specific part of PHY and MAC is shared across network slices, and the service (or bearer) specific part is implemented in each network slice. Hence, in this option, the individual network slices rely on the same radio access technology but customize their operation at lower layers through parameterization and at higher layers through customized implementation.

In this option, the MAC layer and with it the TTI-based scheduling is part of the shared functions. This tighter control of the radio resources may increase the achievable multiplexing gains and alleviates the consistency requirements compared to Option 1. On the other hand, the interface between the shared and dedicated part is more complex because more interactivity between both parts is required, e.g., the actual resource assignments in the shared part must be reflected by the

RLC layer in the dedicated part (segmentation, ARQ), channel measurements need to be exchanged, etc. All this information needs to be handled by this interface between shared and dedicated part. Again, this interface may coincide with the interface between the DU executing the common part and the edge/central cloud executing the dedicated part. Consequently, also the SDM-X in Option 2 is more complex than in Option 1 but it also offers more degrees of freedom for resource sharing among slices as well as opportunities for future evolution. While the MAC Scheduling function is element of the common part and therefore does not allow for customization by the tenant, the QoS Control function is element of the dedicated part and offers the possibility for slice-specific control based on policies and constraints customized by the tenant. The information about the individual services and their QoS Scheduling subblock, which then enforces those QoS constraints as part of its multi-service framework.

Compared to Option 1, the RRC Cell and RRC User function blocks are common to network slices of one UE and therefore the mobility handling is significantly simplified. The PDCP function block is part of the dedicated network slice implementation and as such also the security context in the RAN, i.e. the tenant may implement the security functions of each network slice. Furthermore, the RLC and PDCP function blocks are listed both on the common and dedicated part. The function blocks in the common part are used for the control signalling, which is user specific, and the RLC/PDCP function blocks in the dedicated part are used for the data layer, which is service/bearer specific.



Figure 2-4: Integration of RAN slicing Option 2 and multi-connectivity based on MAC and PDCP layer, respectively

Furthermore, in this option most ingredients of a flexible RAN are included in the common part such as the flexible numerology mentioned before, carrier aggregation, RRC state handling, etc. Hence, it is not necessary that each slice needs to reimplement this flexible RAN functionality but may use one efficient implementation while customization through parametrization is still possible. In addition, multi-connectivity can be well integrated with this RAN slicing option as illustrated in Figure 2-4.

2.1.3 Slice-aware shared RAN (Option 3)



Figure 2-5: Functional control and data layer architecture, RAN slicing Option 3 (shared RAN)

The third option is illustrated in Figure 2-5 and corresponds to the case where the RAN is a common resource for network slices. This case is similar to the ongoing discussion in 3GPP SA [23.799], which considers slicing only in core network and treats the RAN as a common resource [23.501]. In 3GPP LTE/SAE, eDECOR has been introduced which allows for implementing dedicated core networks such that each UE may be connected to a customized core network. Furthermore, 3GPP LTE/SAE allows for connecting a UE to multiple PDNs while using the same serving gateway (S-GW). This is the main limitation, which is alleviated by 5G where each UE may be connected to more than one network slice (i.e., core network). Hence, this option would provide seamless migration and requires minimal changes to current standards.



2.1.4 Exemplary use of RAN Slicing

Figure 2-6: Exemplary deployment using RAN slicing in the context of an industrial campus deployment

Figure 2-6 shows exemplarily different network slices that may co-exist within the same deployment such as an industrial campus network. In this particular example, we stagger different RAN slicing options, i.e. RAN slicing Option 1 multiplexing data from two network slices and the common part of RAN slicing Option 2.

In particular, the following slices are shown in Figure 2-6 (note that those are examples which may have different requirements or implementations in other scenarios):

- 4G (legacy systems): This network slice represents a legacy network such as 4G which is integrated into the same infrastructure. In existing deployments, it is unlikely the installed terminals could be replaced all at once, but those terminals will rather co-exist with new deployments utilizing 5G technologies enabling novel service. Hence, the architecture needs to integrate both existing standards as well as upcoming technologies, which is reflected by RAN slicing Option 1.
- Critical IoT slices may be implemented within factories for novel and critical IoT services [5GPPPWP] such as wireless robot control. They may also use different radio access technologies in order to satisfy the requirements and constraints imposed by the use cases and requiring customization down to the PHY layer. Compared to other slices, they may not require NAS signalling (or only very limited).
- Massive IoT slices may not need customized PHY and MAC layer which also increases the deployment flexibility of sensor nodes as well as reduces their costs because non-proprietary technologies are used. The required access may only be local at the factory in order to guarantee data privacy. The slice may be (partly) operated by the industrial tenant, e.g., data layer provided by tenant and control layer (partly) provided by third party (e.g. MNO).

- eMBB slice may be provided by MNO but operated only on the factory ground in order to provide data access into the public network, e.g., for video telephony and similar services.
- The two massive IoT slices at the right may be operated by the MNO in order to provide data access to third party companies operating equipment on the factory ground, e.g. logistics companies tracking goods. The slices may exploit additional infrastructure on the factory grounds in order to improve coverage while the factory owner may offer these services jointly with the MNO.

[5GN-D32] and [5GN-D33] provide more detailed analyses of business models and stakeholder roles which would affect the above architecture.

2.2 Standardization relevance and potential

This section does not aim to provide an architecture comparison between 5G NORMA and 3GPP 5G systems in general, rather to use 3GPP 5G as a best possible reference to show standardization relevance and potential of 5G NORMA innovations. Note that 5G NORMA is primarily a research oriented project and 5G NORMA does not aim to specify complete network systems for 5G at the first place.

The baselines of 3GPP next-generation radio access network (NG-RAN) are being captured in [38.300]. As depicted in Figure 2-7, NG-RAN consists of gNBs and/or eNBs, providing the uplane and c-plane (named d-layer and c-layer in 5G NORMA) protocol terminations for the radio-interfaces towards the UE. The gNBs and eNBs may be interconnected with each other via Xn interface. The gNBs and eNBs are connected to 5G core network (5GC) via NG interfaces, more specifically to AMF (Access and Mobility Management Function) via the NG-C interface and to UPF (User Plane Function) via the NG-U interface [23.501]. The functional split between NG-RAN and 5GC is given in Figure 2-8.



Figure 2-7: Overall architecture of 3GPP 5G network (Figure 4.1-1 of TS 38.300)



Figure 2-8: Overview of functional split between 3GPP NG-RAN and 5GC (Figure 4.2-1 of TS 38.300)

A realization of the 3GPP 5G system on 5G NORMA would have, for the 5GC part, the control functions of AMF and UPF implemented as a set of SDM-C applications, e.g. NAS Control would be part of such a set. The UPF is part of the NAS d-layer function block. The NG-RAN, i.e. the functions of the gNB, would be provided by the following 5G NORMA function blocks:

- Radio Admission Control: Multi-tenant Scheduling SDM-X Application (Multi-tenant Policy for RAN Slicing Option 1)
- Connection Mobility Control: Mobility-Management SDM-C Application [5GN-D52]
- RB Control and Measurement Configuration & Provision: RRC User and additionally RRC slice for RAN slicing Option 2
- Inter-cell RRM and Dynamic Resource Allocation (Scheduler): MAC Scheduling distributed control function, with tenant-specific control via SDM-X applications Multi-tenant Scheduling and SDM-C/X application QoS Control

Further aspects of NG-RAN architecture and interfaces can be found in [38.801]. In particular, the gNB architecture consists of a central unit (CU) and one or more distributed unit (DU) connected to CU via Fs-C and Fs-U interfaces for C-plane and U-plane, respectively, as shown in Figure 2-9.

3GPP has considered various options for reconfigurable RAN functional split between CU and DU, as shown in Figure 2-10. 5G NORMA omits split Options 3 and 5 within RLC and MAC, respectively, and High PHY, Low PHY and RF PHY relate to PHY User, PHY Cell and PHY TP of 5G NORMA, respectively.

For multi-connectivity architecture, 3GPP has been focusing on RAN level integrated dual connectivity between NR and LTE. RAN supports for network slicing are also addressed in [38.801], focusing on enabling selection of network slice and CN instance by a gNB. NG-RAN also incorporates supports for UEs in inactive mode, QoS flow management and mapping to data radio bearers as well as new challenging services such as ultra-reliable and low latency communications (URLLC).



Figure 2-9: The gNB architecture with CU and DUs (Figure 11.1.3.8-1 of TR 38.801)



Figure 2-10: Functional split options between CU and DU (Figure 11.1.1-1 of TR 38.801)

In a high-level comparison between 5G NORMA and 3GPP 5G access network architectures, the following can be observed:

- Both are sharing views in many aspects, including flexible RAN functional split or function decomposition, multi-connectivity with NR-LTE tight interworking and RAN level integration, supports for network slicing, UE inactive state, QoS flow and radio bearer mapping and management.
- 5G NORMA, however, has more extensive considerations on SDN/SDMC-enabled architectures and flexible multi-service and multi-tenancy supports via means of network slicing, resulting in a more centralized SDMC based architecture. This enables and facilitates different RAN slicing innovations with more flexible and dynamic adaptive function decomposition, as presented in Section 2.2.

The following Table 2-1 provides a more insightful view from 5G NORMA innovation perspective, which of the 5 5G NORMA overall objectives they address and how in this respect they compare to the approach 3GPP has taken. As can be seen, 3GPP in many aspects took similar approaches to 5G NORMA. Most notable differences stem from the different approaches of 5G NORMA and 3GPP. While 5G NORMA start from highest possible flexibility and reduces where technically necessary, 3GPP starts from with the legacy approach of a network of entities and a RAN core split and adds flexibility as needed for the short term targets. The results become visible especially at the RAN, where 5G NORMA applies its SDMC concept and offers RAN

slicing, while 3GPP, at least the first phase of 5G standardization, just adds support for network slicing to NG-RAN, i.e. a non-sliced RAN comparable to RAN slicing Option 3 but which lacks the programmability offered through SDMC and accordingly the means for slice-individual customization of the RAN.

Innovations	5G NORMA	3GPP 5G
Adaptive (de)composition and allocation of mobile network functions	 Flexible and cloud-RAN based approach Multi-RAT MC with tight 5G- LTE interworking More dynamic SON based 	 Flexible and cloud- RAN based approach Multi-RAT MC with tight 5G-LTE interworking Semi-static reconfigurable
Software-defined mobile network control (SDMC),	• NFV- and SDMC-based c/d- layer to realize a slice- individually customizable network of functions	 NFV is more visible in 5GC Certain legacy design principles are kept to allow RAN and 5GC evolve independently as much as possible
Joint optimization of mobile access and core network functions	 RAN support for in-bear QoS differentiation RAN support for context aware service delivery (multiservice radio access) UE in INACTIVE and RAN paging as defined by the UCA concept Integration of mmAP clusters 	 Through functional split between NG-RAN and 5GC, e.g., RB control and some QoS flow management is moved to gNB Feature specific optimizations, e.g., supports for MTC, UE in INACTIVE and RAN paging, RAN support for context aware service delivery
Multi-service- and context-aware adaptation of network functions	 URLLC support Multi-service radio access and high data rates via mmAPs 	URLLC support
Mobile network multi- tenancy	 Network slicing including RAN slicing Extensive considerations on characterization/categorization of tenants for future mobile networks 	 Network slicing for 5GC and NG-RAN support for network slicing

Table 2-1: 5G NORMA innovations vs. 3GPP 5G, grouped according to the 5 key
innovations

More specifically, 5G NORMA partners have (so far) submitted 21 contributions to 3GPP RAN2 and RAN3, the two 3GPP specification groups that the innovations of WP4 mostly apply to. The complete list of standards contributions will be documented at the end of the 5G NORMA project with [5GN-D72]. These standards contributions express 5G NORMA views and findings in some selected topics, including:

- Support of novel services in 5G beyond mobile broadband, in particular ultra-reliable low latency communications (URLLC) ([5GN-D41] Section 1.3 innovations 1, 3, 5 and 13);
- Requirements and principles for network slicing, including isolation between network slices and how these requirements can be supported in slicing the RAN ([5GN-D41] Section 3.3 and Section 1.3 innovation 2);
- RAN slicing, in particular radio resource isolation ([5GN-D41] Section 3.3);
- Multi-connectivity supports including tight interworking between LTE and 5G NR ([5GN-D41] Section 1.3 innovations 1, 3, 4, 5 and 7);
- RAN level QoS flow management and PDCP relocation ([5GN-D41] Section 1.3 innovations 11 and 12);
- Support of UE in an INACTIVE state with RAN level paging ([5GN-D41] Section 3.2.1.4 and Section 1.3 Innovation 6)
- Multi-connectivity support, as described in Section 6.1;
- RAN support for network slicing in RAN slicing Option 3, described in Section 2.1.3, with two innovations to enable the UE to access multiple network slices simultaneously:
 (i) PDCP multiplexing of service flows from different network slices and (ii) coordinated dynamic scheduling and semi-persistent scheduling based resource allocations. These are described in Part II Chapter 17;
- SON based flexible configuration of 5G RAN protocol stacks aims for possible support of more flexible and dynamic on-the-fly adaptation of the RAN function decomposition, as compared to rather semi-static reconfigurable options for RAN functional split between CU and DU (Figure 2-10) ([5GN-D41] Section 1.3 Innovation 9);
- UE agent based end-to-end connection decomposition concept, which is directly related to and goes beyond 3GPP Rel-14 study on context aware service delivery in RAN for LTE [36.933]. This innovation also provides an effective means for integrating and leveraging lower data rate satellite communications into 5G network systems ([5GN-D41] Section 1.3 Innovation 10).

3 Use of functional decomposition for service adaptation

A fundamental aspect of flexible architecture is its ability to adapt to multiple services with different requirements, yet with minimal or no infrastructure amendments. In this regard, function decomposition plays an important role in tailoring the network operation to its specific needs, thereby, attaining the desired adaptability in terms of allocating the network functions and resources based on the requirement and deployment characteristics.

The remainder of this section is divided into three major parts. First, the concept of flexible function selection and placement is put forward, thereby highlighting the fact that different slices may involve different network functions located at different places. Second, the notion of multi-service radio access is analysed, elaborating on the conceptual differences between multi-connectivity, multi-tenancy as well as multi-service. Third, the above analyses are complemented with security considerations.

3.1 Function selection and placement

Conceptual relation to network slicing

As discussed in Chapter 2, the concept of flexible network design entails the ability of network functions to be deployed in an adaptive fashion. This practically means that, in contrast to conventional network architectures where network functions are pre-configured, network functions are dynamically configured both in terms of functional operation and physical location. In fact, with the flexible design approach embraced in 5G NORMA, a new dimension in network architecture is offered, where functions are systematically extracted from a pool of resources and utilised exactly when and where they are needed.

Such network design is in line with the architectural approach of *network slicing*. In particular, the fundamental conceptual element of network slicing is the separation of services into independent logical networks. In this respect, the 5G NORMA architecture comprises decomposed functional elements, which can be allocated on demand, such that different slices are associated with different network function blocks. Moreover, the diverse requirements of network slices (for example, in terms of latency of end-to-end services) imply that the requirements on the physical placement of individual network functions also vary.



Figure 3-1: Function selection and placement (RAN slicing Option 2)

A closer look on flexible function allocation

We distinguish two basic levels of flexible function allocation, namely *function selection* and *function placement*. *Function selection* refers to the set of the network functions deployed in a slice, thereby reflecting the ability of the 5G NORMA architecture to construct the functional operation of a network slice via function blocks. *Function placement* refers to the physical network locations, where those functions are deployed, again on the basis of a particular network slice with given requirements.

The above two levels of flexible function allocation are illustrated in Figure 3-1. In particular, we map the concepts of function selection and function placement onto the principal types of network slices anticipated for 5G, namely mobile broadband, low latency and mission critical slice. As can be observed from Figure 3-1, different slices comprise different set of functions placed in different parts of the network, aiming to better meet the diverse requirements imposed. Next, we treat each of the three slice types separately and elaborate on the effect of function selection and placement on the considered architecture, focusing on the data layer.

QoS control per network slice

In line with the above consideration, the network will be designed as a multi-service adaptive mobile network, where the network resources fulfil the service requirements in a flexible and dynamic way per each network slice.

The QoS Control function block will be in charge of the network monitoring and configuration. For the QoS monitoring tasks, two approaches can be considered: 'intrusive' and 'non-intrusive'. The non-intrusive methods are purely based on monitoring the already available QoS parameters (i.e., latency, jitter, or packet loss). On the other hand, intrusive methods are based on installing specific purpose SDM-C applications to get additional QoS parameters. Other solutions are focused on including new network functions (e.g., network probes and analysers, deep packet inspectors, etc.) that are responsible for capturing the traffic from a certain service and analysing its performance.

A QoS agent (specific piece of software) in network function blocks (VNFs, PNFs) will monitor the QoS status at run time, e.g., connection speed or packet loss rate. When some monitoring parameter does not fulfil the proper values, an event will be generated and captured by the SDM-C/X and distributed control functions. Then, the QoS Control function block will receive the reports of these events from the controllers and will aggregate and analyse the data, enforcing new configuration actions (resource re-scheduling, new resource assignments) on the network function blocks to meet the QoS requirements of each mobile service according to the assessment results. The adjustment happens when the evaluation result violates a threshold that is defined by service requirements.

The broadband slice

Function selection aspect of flexible design: The broadband slice is designed to serve applications associated with high data rate transmissions. This means that the broadband slice involves network functions, the primal use of which is to increase the overall throughput. Hence, considering the network functions analysed in deliverable D4.1 [5GN-D41], the broadband slice would involve the PDCP split bearer as well as the MAC carrier aggregation (CA) function blocks from the data layer domain. Of course, any combination of the different transmission legs is possible. For instance, the MAC CA function can be applied to one or more components of the bearer, which is split in PDCP by means of the PDCP split bearer function. In Figure 3-1, an exemplary implementation of function selection for the broadband slice can be observed in the left part.

Function placement aspect: Beside function selection, a typical example of locating network functions based on the corresponding broadband service is also depicted in the left part of Figure 3-1. Specifically, we notice that the PDCP and PDCP split bearer blocks are located at the *edge cloud*. This location has been chosen to facilitate multi-connectivity as well as to minimise the

mobility signalling to the core network [Rav16]. We also notice that, in contrary to PDCP functionalities, RLC as well as lower layers of the protocol stack need to be co-located at the antenna site. This is due to their real-time operation, implying a synchronous interaction between one another and thus a low latency inter-layer communication.

The low latency slice

Function selection: Contrary to the broadband slice, the main requirement for the low latency slice type is not, in principal, high data rate. In fact, with the exception of virtual reality services where both latency and throughput need to be taken care of, low latency applications are usually associated with machine-type packets of bursty nature, the size of which is considerably smaller than that of the usual packet size of MBB applications. This means that multi-connectivity is not a crucial element of the low latency slice and is therefore not included. This is depicted in the middle part of Figure 3-1. It is also important to notice that for slices supporting ultra-low latency requirements, the outer Automatic Repeat request (ARQ) function is also not included in the RLC part of radio stack, thus the RLC operates in the unacknowledged mode. The RLC UM is thus included in this slice.

Function placement: The low latency slice entails special, tight requirements in terms of the endto-end latency. This demands for the realization of the network functions in a location which is as close to the antenna site as possible. As a result, there is no use of the edge cloud for this slice; instead, all network functions are placed in the (logical) edge cloud co-located at the antenna site (which may be a small NVFI integrated into the DU), including those supporting PDCP functionalities. In other words, the edge cloud is unused or "transparent" for low latency slices.

The mission critical slice

Function selection: This slice is designed to mainly support services with ultra-high reliability levels². Such reliability levels are considerably higher than the usual values in LTE, and are sometimes referred to as the "five nine reliability", implying that communication should be established without errors for at least 99.999 % percent of time. It is important to note that the above required reliability level refers not only to the time, for which a UE is within a given area with decent coverage, but it rather refers to the overall percentage of time, thus including also areas with limited coverage. In other words, the requirement for high reliability here encompasses the need for eliminating coverage holes, such that sufficient coverage is provided for 99.999 % of UE operation. Such ultra-high reliability levels are hard to achieve with existing standards, implying that new access techniques need to be employed. To this end, the data duplication method has been recently proposed as a special case of multi-connectivity. Contrary to multiconnectivity variations proposed for broadband services [36.808], data duplication involves replicating the same message over multiple links, so as to leverage the independent streams and increase the probability of correct reception [Rav16], [MVD16]. Of course, this would involve modifying the functionality of PDCP, so that data streams are not split into multiple legs (as is the case in the broadband slice) but duplicated instead, thus achieving redundancy. This implies that special coordination needs to take place to ensure that replicas of correctly transmitted messages are discarded, and that unnecessary further transmissions of already successfully received data are avoided whenever possible. It is worth noting that the missing critical slice involves the acknowledged mode in RLC (i.e., the RLC AM block in the right part of Figure 3-1), as opposed to the RLC UM block used for the low latency slice.

Function placement: As mentioned above, high reliability entails centralised functionalities, which are mainly associated with the coordination of the duplicate transmission. In particular, given that duplicated and redundant streams must be jointly coordinated, a distributed

² Although high reliability and low latency are sometimes studied on a common basis (known as ultra-reliable low latency communication services, URLLC), ultra-high reliability entails a different architecture design which is the focal point of this paragraph.

implementation of the RLC and PDCP functionality would entail a large signalling overhead between the involved nodes, which is in principle not desired. As a result, the realization of the functions in the PDCP and RLC stack is expected to take place in the edge cloud (c.f. Figure 3-1, right part).

3.2 Multi-service radio access



Figure 3-2: Multi-service vs. multi-tenancy and inter-tenant multi-connectivity vs (single tenant) multi-connectivity

Multi-connectivity is the technique to connect a single user (UE) to multiple distinct instances, i.e. cells, of a single service respectively slice, cf. Figure 3-2. *Multi-tenancy*, on one hand, means that multiple tenants share common infrastructure, including sharing a single air interface instance. On the other hand, with respect to a single UE, multi-tenancy means that the UE connects to the services of multiple tenants concurrently, i.e. in the sense of multihoming, where the UE connects to multiple data networks in parallel. This sharing by multiple and connection to multiple tenants may be for the same single service, i.e. is independent of whether the service offered by each tenant differs or not. A UE may need to concurrently support both, multi-connectivity, the connection to multiple instances, and multi-tenancy, the connection to multiple tenants, due to different availability of the tenants' services at different cells or due to tenant-individual mobility control. Such *inter-tenant multi-connectivity* therefore occurs, when a UE connects to multiple slices, but the connectivity to the slices is provided to the UE via different cells.

In 5G NORMA, each slice provides exactly one telecommunication service, e.g., eMBB or URLLC. Having just one service per slice best captures the benefits of slicing as it allows for optimizing a slice for its specific service as discussed in Section 3.1. Accordingly, a UE that utilizes multiple services in parallel connects to multiple slices in parallel. On fully virtualized network infrastructure, where specifically communication/transport resources are fully virtualized, the same means, namely *network slicing*, is used to implement both multi-tenancy and multi-service. This is typically the case in the non-access stratum. In contrast, communication resources in the access stratum are not virtualized into generic units of transport resources. The implementation of transport over the air interface itself is very much dependent on the service and deployment needs and thereby represents an important differentiator for tenants. Accordingly, implementation of transport differs between services and between tenants, notwithstanding the use of a common implementation where this is deemed sufficient.

The previous Section 3.1 showed how the functionally decomposed c/d-layer is leveraged to realize different services as well as how to support different deployments, considering isolated services. The following elaborates on how the functionally decomposed c/d-layer is utilized to implement multiple services concurrently within a single air interface instance, i.e. how to realize a *multi-service radio access*.



Figure 3-3: Example 5G New Radio multi-service RAT for RAN slicing Option 2

Figure 3-3 shows a possible multi-service capable future 5G RAT implementation using the functionally decomposed c/d-layer as proposed by 5G NORMA. In this example, the access stratum supports four services, eMBB, URLLC, mMTC and eMBMS, the former two additionally for D2D. Each service has a different realisation on the PHY layer (PHY User), but all share a common carrier, i.e. share the mixed signal and analogue processing part (PHY TP). In a simplified implementation, a subset or even all services may employ the same subcarrier spacing and symbol length (assuming an OFDM-based system). This would allow the lowest part of PHY Cell, namely the (i)FFT and CP insertion/removal, respectively the subband filtering for filtered waveforms like UF-OFDM, to be shared in addition to PHY TP. The upper part of PHY Cell and PHY User differs for each service, creating an optimized implementation for the specific service. Such service-specific PHY User implementations and their respective RRM may be (with references to selected results from H2020 FANTASTIC-5G):

- RRC Cell provides signals for synchronisation and channel measurement as well as initial access (IAC) to the system (system broadcast, contention-based UL access and basic DL channel), being simple and robust to be accessible by all UEs from low cost machine type to high performance broadband ([F5G-D31] Section 6.6.4.1 synchronization signal);
- eMBB is optimised for maximum spectral efficiency, utilising TTI duration and pilot patterns (DM-RS) adapted to the coherence time and bandwidth of the radio channel as well as multi-cell multi-user capable advanced spatial multiplexing and receive processing (CoMP) ([F5G-D32] Section 3.3.2 physical downlink control channel);
- URLLC is optimised for lowest latency at the cost of spectral efficiency, utilising very short TTIs and accordingly higher DM-RS overhead and an RRM able to pre-empt other services ([F5G-D42], Section 2.4 dynamic resource allocation);
- D2D conveys control information only (semi-persistent radio resource assignments) in a robust way, specifically for its respective service eMBB or URLLC [F5G-D42], Section 3.1 D2D);

- mMTC is optimised for processing massive amounts of sporadic small packets, employing contention based access, open loop link adaptation (MCS selected and signalled by UE) and autonomous synchronisation (ATA) for maximum energy efficiency (UE side) in UL and size-optimised and robust DL control to provide extended coverage ([F5G-D32] Section 6.1 waveform candidates, [F5G-D42] Section 3.2 efficient massive access protocols);
- eMBMS is adapted to realise a DL single-frequency network (MBSFN), employing intercell coordination and extended CP length and optimised subcarrier spacing for both control and data layer ([F5G-D42] Section 3.3 MBMS).

The distinct PHY layer implementation per service, i.e. the different PNF types used for PHY Cell and PHY UE, imply distinct PNF instances per each service. Accordingly, Layer 2 (MAC, RLC and PDCP) must employ different VNF instances per each service, too, but the same type of VNF may be used for each instance. Today's common practice and reasonable starting point is to use the same Layer 2 protocol implementation but with different sets of allowed options and parameter settings per each service. When the system evolves over time and especially when new not yet foreseen use cases emerge, a new VNF type can be introduced in any service without impacting the others, thereby guaranteeing a future proof overall system design.

The distributed control VNFs RRC Cell, RRC User and MAC Scheduling again differ in their type. Note that the RRC protocol that RRC Cell and RRC User use to communicate with UEs may be the same for several services similar as the Layer 2 protocols PDCP, RLC and MAC may be the same, but the control logic and relevant state is specific to each service and accordingly the NF type of RRC Cell and RRC User.

For RAN slicing Option 2 depicted here, the data layer implementation of the access stratum becomes slice-specific (and therefore tenant-specific) at RLC and above. The same considerations as for the (access stratum) control layer RLC and PDCP implementations hold, i.e. NF instances differ but NF types may be common among services respectively slices. Since implementation is now slice-specific, multi-service (and multi-tenancy), like in the non-access stratum, are provided jointly by means of network slicing.

So far, the access stratum provides multi-service differentiation. Furthermore, the common part of the access stratum needs to be *slice-aware* for multi-tenancy support. Common NF instances shared by multiple slices maintain the one-to-one mapping of traffic to individual slices and SDM-X provides each slice, via the 5GNORMA-SDMC-SDMX interface (cf. Section 4.1), the means to influence how traffic of their own slice should be processed. Thereby, both separation of services and separation of tenants become available in the common part as it is available in the dedicated part of each slice via network slicing. As a result, network slicing in 5G NORMA effectively becomes end-to-end, spanning the whole path from data network (respectively end-user service) to the UE.

3.3 Security considerations

Flexible function placement may pose a security issue, as a network function may be placed in different environments with different security threats. To cope with that, when designing and implementing a network function for which flexible placement is applicable, the security architecture of this network function must consider all possible placements, including the "worst case" placement, i.e. the most hostile or exposed environment. If a function is generally aware of its placement and designed to adapt to the actual placement, this may also apply to security measures implemented within that function.

Also, flexible function allocation may imply that two communicating network functions could in some cases reside on a single hardware platform, with few possibilities for external attackers to interfere with the communication. However, the same two network functions could in other cases reside on two geographically separated data centres, interconnected by inherently exposed wide

area connections. There can be even dynamic changes between these two setups during the lifetime of the communicating network functions.

Section 4.3 describes the general principles for securing interfaces in the 5G NORMA architecture. For inter-domain interfaces, cryptographic security associations are always the means of choice, independent of the location of the network functions. Intra-domain interfaces within the distributed cloud environment can be isolated and secured by general security mechanisms of this environment that can relieve network functions from the need to setup dedicated security associations between each other. As explained in section 4.4.2, this protection and isolation can be provided independent of the current location of each network function, implying that network functions in this case need not adapt the protection measures applied to their inter-communication in dependency of the current placement.

Clearly, that doesn't mean that each network function can simply execute on any incoming request – certain network functions may still need to be aware of their communication peers and treat incoming requests accordingly, e.g., distinguish between requests coming from different VLANs to which the function is connected. In particular, due to the flexible function selection, the number and nature of the communication peers of a network function may vary, depending which network functions have been selected to compose the specific network or slice. This may impact the way a network function must treat traffic received by its communication peers with respect to security, e.g. whether there is a need for filtering out specific messages or not. Clearly, this needs to be taken into account when designing and implementing network functions that must support different sets of communication peers in different network setups.

The support of multiple services in a single network is facilitated by the usage of network slices. Isolated network slices in turn allow the implementation of individual security policies per slice. In the RAN, this mainly affects access stratum security policies (see [5GN-D32] section 5.4.4 for the 5G NORMA access stratum architecture). As an example, different slices may use a different choice of crypto algorithms, or different preferences which algorithm to choose. As another example, some slices may enforce encryption and maybe even integrity protection of the data layer, while others may allow an unprotected data layer. The individual security setup per slice can be enabled by making the PDCP handling a slice specific function, rather than a common function.

4 Control and data layer

This chapter provides a concise characterization of the control and data layer design of the 5G NORMA architecture. In Chapter 2 we discussed options for which functions are common to all slices respectively tenants and which are dedicated, while Chapter 3 elaborated on which functions to select and where to place to adapt to individual service and deployments characteristics. In the following subsections, we now focus on the general defining characteristics and according classification of functions into either data layer, distributed or centralized control, which is independent of any multi-tenancy and multi-service aspects. The complete list of control and data layer function blocks along with their individual functionality can be found in Annex B.



4.1 Centralized control

Figure 4-1: 5G NORMA SDMC interfaces

Figure 4-1 shows all functions and interfaces of the centralized control layer of the 5G NORMA architecture. The centralized control layer consists of the SDMC Controller (SDM-C), the SDMC Coordinator (SDM-X), and SDM-C/X applications (App). It interfaces, exclusively through SDM-C and SDM-X, with the distributed control and data layer as well as with the management and orchestration (MANO) layer. The MANO layer is out of scope of WP4 and is not further considered here. For details on 5G NORMA MANO functions see [5GN-D33].

Distributed network functions dedicated to a specific network slice, the *Dedicated NFs* in the above figure, are under exclusive control of the slice's own SDM-C. The SDMC concept foresees that the control logic is implemented as part of one (or multiple) SDM-C application(s), while the controlled functions are called agents. The SDM-C controls distributed NFs, i.e. the agents, via its *5GNORMA-SDMC-NF* interface and in turn provides SDM-C applications the means to monitor and program the agents through its *5GNORMA-SDMC-App* interface. In SDN controller terminology, *5GNORMA-SDMC-App* represents the controller's Northbound Interface (NBI),
which exposes a well-defined Application Programming Interface (API) to the applications running "on top" of it. The interface *5GNORMA-SDMC-NF* represents the Southbound Interface (SBI) which a controller uses to access the different classes of distributed NFs under its control. Thereby, the SDM-C hides and abstracts specifics of the underlying distributed agents and their control interfaces by translating them into a well-defined API that it then exposes to the central control logic.

A novel aspect of the SDMC concept is that some of the distributed NFs are shared by multiple network slices. Such common distributed NFs are put under control of a single SDMC Coordinator, the SDM-X. It coordinates the control information coming from the different slices' SDM-Cs via the 5GNORMA-SDMC-SDMX interface. Through this interface, each SDM-C controls those NFs that are shared with other slices up to the extent exposed by the SDM-X through this interface, i.e. typically only to an extent that enables the SDM-C to control the user data traffic of its own slice. For this purpose, the SDM-X needs to execute specific rules to be able to make meaningful coordination decisions. For example, SDM-X must authorise the requests of SDM-Cs and reject requests that are not in line with the SLA between mobile service provider (MSP) and tenant. SDM-X Apps include the according policies and can provide such rules via the 5GNORMA-SDMX-App interface. This generates a single meaningful outcome from the various SDM-Cs' requests coming in through the 5GNORMA-SDMC-SDMX interface before being forwarded towards the agents in the data layer or distributed control layer. Hence, the SDM-X runs applications that exclusively control common network functions and resources shared by multiple slices. These applications are run by the stakeholder that determines the set of common NFs, i.e., usually the MSP.

Southbound Interface – 5GNORMA-SDMC-NF and 5GNORMA-SDMX-NF

A major achievement of 5G NORMA is the separation of control and execution of NFs and the centralization of the control part in SDM-C and SDM-X – with the exception of distributed control functions as detailed in Section 4.2. This separation implicates that both parts are connected through the *5GNORMA-SDMX-NF* interface for common NFs and *5GNORMA-SDMC-NF* interface for dedicated NFs. These interfaces are realized as dedicated Southbound Interface (SBI) protocol plugins. The kind of protocol plugin depends on the respective class of distributed NF:

- controlling packet forwarding among NFs, i.e. "classical" SDN,
- monitoring and controlling further NF behaviour specific to mobile networks, as newly introduced with SDMC.

For the first class, a set of state-of-the-art protocols have been proposed and studied in research projects on 5G such as 5G-Crosshaul, which evaluated OpenFlow, IETF ForCES and P4 protocols [5GC-D21]. OpenFlow operations are based on flow entries which are stored in flow tables (one or several) within the OpenFlow switch. An OpenFlow node is modelled by a number of network ports and a pipeline including a set of flow tables. The communication between the controller and network nodes is carried over SSL/TLS-secured transport channels. In OpenFlow, the set of match operations and actions, applied to the packets of the different flow tables, is fixed while P4 offers a higher level of flexibility at the cost of an increased processing. In the 5G NORMA architecture, some environments may have strict timing requirement, which should be considered by the forwarding abstraction. In addition, a common frame format should be preferred for the information exchanged at the SBI.

The second class of SBI plugins serves the NFs of which the control logic is implemented as part of an SDM-C application and a specific set of information is needed by this application to operate. The SDM-C extracts such information through the 5GNORMA-SDMC-NF interface.

The SDMC-enabled functions differ significantly from each other, but it should be possible to derive a common interface. In Section 3.2 of [5GN-D41], the individual properties of the SDMC-enabled control functions have been studied to obtain the requirements of this novel interface. In particular, the SBI should support the abstraction of the underlying topology and NFs for all SDM-C/X applications, the monitoring of control layer parameters, and the configuration of NF

related parameters (scheduling policies, QoS values, etc.). Some SDMC-enabled applications like SON, RAN Paging, NAS Control, or QoS Control require low latency communication (in the order of the envisaged 5G handover process or even lower), which is imposed as a requirement on the 5GNORMA-SDMC-NF SBI plugin as well. However, the properties of these NFs and their requirements on the SBI could not yet be investigated in detail in 5G NORMA and are up for further evaluation in 5G-PPP Phase 2.

According to [MYV+15], about 95 % of network devices are configured by proprietary Command-Line Interfaces (CLI). The Open Networking Foundation (ONF) developed the OpenFlow Management and Configuration protocol (OF-CONFIG 1.2) [OF-CONFIG] to provide the possibility of reconfiguring the network devices in order to add new and improved communication capabilities. The basic OpenFlow, in fact, does not provide this functionality. The OF-CONFIG protocol uses NETCONF [RFC6241] as transport protocol for programmatic management of network configuration. NETCONF is built on protocols such as SSH to provide secure reconfiguration of the network devices and XML-based data models. NETCONF exploits YANG, a language for developing standardized configuration data models [RFC620].

Also, the SBI candidate protocol plugins (e.g. OF-CONFIG, proprietary APIs) for monitoring and controlling the SDMC-enabled NFs may differ in flexibility (granularity in the packet classification) at the cost of an increased pipeline processing. However, the choice of a common frame format can optimize the packet processing speed.

For the identified SDMC-enabled NFs, a list of parameters that SDM-C/X applications require for configuration and monitoring need to be defined. The configuration parameters are those that the SDM-C/X can set trough the SBI (for example those related to the scheduling policy) while the monitoring parameters are the ones which need to be tracked by the SDM-C application where the control logic resides (for example the CQI). The selection of parameters comprises a tradeoff: a too small set of parameters could inhibit the efficiency of some applications while a too wide set could negatively affect solution cost and scalability.

In Table 4-1, an example of a subset of parameters is shown for the eMBMS Control NF.

Parameter	Configuration	Monitor	Notes
Data MCS	•	•	Formerly SCTP/IP on
			standardized M2 interface
MCH Scheduling		•	Formerly SCTP/IP on
Period			standardized M2 interface
CQI		•	To be communicated through
			SBI to SDMC

Table 4-1: Subset of	parameters obtained	for eMBMS Control
----------------------	---------------------	-------------------

Once the parameters are identified, they can be mapped in the protocol stacks of the SBI protocol plugins. In this case, some extensions to the existing state-of-the-art SBI protocols may be required for enabling the appropriate behaviour of the SDM-C application.

5G NORMA aims to integrate in its architecture also PNFs: one possibility is to develop a proprietary SBI plugin as done by 5G NORMA partners in WP6 [5GN-D62]. Here, in order to develop Demo 1, the project partners Azcom and Nomor designed a novel communication SBI protocol for the proprietary API of the PNF, which is designed for communication between the software SDM-C and the hardware legacy LTE eNB(s). This communication protocol is flexible and can also be extended for connecting any other network element beside hardware eNBs. All communication messages of this proprietary SBI protocol are structured according to a predefined fixed format that defines the message structure, according to the aforementioned SBI requirements. In Demo 1, the choice of a proprietary SBI protocol guarantees a low latency communication flow between the SDM-C application and the hardware eNB because a

reconfiguration of the scheduling policy is performed. For a detailed description of Demo 1 the reader is referred to [5GN-D62].

On the other hand, an abstraction layer can be used to turn non-SDN devices into, for example, OpenFlow controllable devices, cf. Figure 4-2 [TMM+16].



Figure 4-2: Abstraction layer turning non-SDN devices into SDN-controllable devices

The abstraction layer could provide a conversion layer between OpenFlow configuration messages and the native management interfaces. Again, some extensions to OpenFlow may be required to support the SDM-C application. Moreover, it will be important to understand to which extent the connection between the implementation of the SDM-C controller and the controlled NFs adversely impacts the network performance.

SDM-C/X applications

One of the innovations introduced by 5G NORMA is mobile network multi-tenancy. In the 5G NORMA vision, the multi-tenancy control function, *Multi-tenancy Policy* or *Multi-tenancy Scheduling*, depending on the RAN slicing options used, runs as application on top of SDM-X. The multi-tenancy resource management algorithms described in Section 5 utilize the interface with SDM-X in order to receive from the different SDM-Cs information regarding the slices respectively users that they control, for the latter, this includes the association of users to slices respectively tenants and services as well as the specific QoS requirements of their service flows. Exploiting the information received, the multi-tenancy function is able to derive and apply through the SDM-X the configuration desired in order to maximize the network utility.

Another crucial function is *QoS Control*, which oversees QoS throughout the network. The QoS Control is an SDM-C/X application that runs on both SDM-X and SDM-C to ensure QoS throughout the network slice end-to-end. It is composed of two basic algorithms that allow the QoS monitoring and the QoS enforcement. The first algorithm is in charge of configuring the QoS parameter set that has to be monitored and which receives the events captured when some value is not correct. The second algorithm is oriented towards evaluation and control of parameter settings during the service lifetime and to enforce actions if needed. Further details can be found in Part I Section 2.1.2 of [5GN-D52].

Further SDM-C/X applications that implement WP4 innovations are (for details see the referenced sections):

- *mMTC RAN Congestion Control*, which controls group management to reduce signalling respectively increase scalability of massive machine type communications, cf. Section 6.6;
- *RAN Paging*, which implements RAN-based mobility management with user-centric connection areas (UCA), cf. Section 6.4;

- *SON* (self-organizing networks), which recommends suitable UCAs (Section 6.4) and mm-wave access point clusters (Section 6.2);
- *GDB* (geolocation database), generally assisting in spectrum sharing, managing heterogeneity and providing equipment location awareness, cf. Section 6.7 (an exemplary usage of geolocation information in mMTC RAN Congestion Control is described in Section 19).

4.2 Distributed control

In general, control logic that is very time critical or not efficiently implementable at a central controller usually is implemented as distributed control NF. This section explains this design decision.

Each *MAC Scheduling* instance incurs (at least) one signalling exchange for each MAC and PHY User instance under its control per TTI, i.e. one exchange per each physical cell per millisecond for LTE and an order of magnitude more often for URLLC. Implementing MAC Scheduling as an SDM-C (RAN slicing Option 1) respectively SDM-X (RAN slicing Option 2 and Option 3) application may be feasible in case of a CRAN deployment, where only PHY TP and PHY Cell (in the case of advanced CRAN) are located at the antenna site, while the remainder including MAC and PHY User is co-located with SDM-X/C in an edge cloud within the access network. SBI traffic would be kept within the edge cloud and therefore the available communication bandwidth would be high and latency small. But even in this case, this poses stringent timing requirements on SDM-X/C NBI API to SBI plugin translation capacity respectively the involved (de)serialization/marshalling of control messages. In contrast, in a DRAN deployment, where up to PDCP all functions are implemented at the antenna site, communication bandwidth and especially latency between distributed antenna sites and their single central edge cloud become critical. This applies even more in case of high (user) traffic load and/or in case of URLLC.

This suggests the introduction of distributed control to complement the conceptually centralized SDMC control with the ability to efficiently support (TTI-)synchronous control functions. Implementing those as distributed control avoids the above described scalability issues and, respectively, the need to design SDM-X and SDM-C for synchronous sub-millisecond processing. MAC Scheduling, as a distributed control VNF, may be co-located specifically with those distributed data layer functions that it controls via its "M" interface (cf. Figure 2-2, Figure 2-3 and Figure 2-5). For a DRAN deployment, each MAC Scheduling VNF instance is co-located in the same (logical) edge cloud at the antenna site that hosts the distributed data layer VNFs/PNFs (up to PDCP) of the physical cells under its control. For a CRAN deployment, each MAC Scheduling instance can be executed within the same cloud node as the MAC and PHY User instances it controls, or at least on a node (very) few physical communication hops away, e.g. a node of the same cluster.

The "M" interface is implemented, in the classical case of single cell processing, as a one-to-one logical interface between a single MAC Scheduling instance and a single set of MAC, PHY User and PHY Cell instances under its exclusive control. Generally, though, the "M" interface is a oneto-many interface to support carrier aggregation and CoMP. Vice-versa, MAC, PHY User and PHY Cell also interface one-to-many with multiple MAC Scheduling instances. MAC Scheduling instances may interface among each other to support distributed coordination schemes, but if and how depends on the employed coordination scheme. For example, MAC Scheduling may just interface with SDM-X executing a CoMP-capable Multi-tenant Scheduling control application or a dedicated control application specifically for CoMP. Besides providing a clean separation of control and data layer functionality, CoMP (and under multitenancy also carrier aggregation) is the motivation for separating MAC Scheduling from MAC in the first place. Integrating both fits single (independent) cell processing, but otherwise artificially imprints a legacy architecture, thereby contradicting 5G NORMA's ambition of architectural flexibility and openness respectively future-proofness.

The drawback of implementing MAC Scheduling as a distributed control NF is that scheduling cannot be simply replaced by providing a different SDM-C/X application. The radio scheduler in MAC Scheduling is dependent on and specifically designed for particular transmit and receive processing capabilities of PHY User, i.e., we cannot replace one of without the other.

RRC Cell and *RRC User* incur a much lower signalling load and are less latency critical compared to *MAC Scheduling*. This is mainly because the signalling is asynchronous, periodic or event triggered over a longer period of time, e.g., tens of milliseconds instead of per TTI. Still, both are realized as distributed control, too. This benefits fast reconfigurations triggered by MAC Scheduling to adapt to the time variant radio channel and interference conditions, without the need for a possibly time consuming detour via SDM-X/C. RRC Cell and RRC User hold RRC state and generate according RRC messages to the UEs. RRC messages convey system wide respectively per user configuration for the data layer NFs on the network side is set via the "*C*" *interface* and "*U*" *interface*, which conceptually implies stateful data layer NFs. For a stateless data layer implementation, the control information of the "C" and "U" interface is repeated per each TTI and sent in conjunction with the per TTI control information from MAC Scheduling sends via the "M" interface.

Whether RRC is co-located with MAC Scheduling depends on the chosen deployment option. In case of DRAN, RRC is co-located with MAC Scheduling. If a CU-DU split between PDCP (in CU) and RLC (in DU) is used, RRC is naturally co-located with PDCP, since most signalling messages of RRC User are conveyed via the DCCH logical channel, which employs PDCP for integrity protection and ciphering³. For RRC Cell (BCCH, PCCH and MCCH logical channels), there is no "natural" placement as it interfaces both with SDM-X/C (SON, RAN Paging, eMBMS Control) and MAC Scheduling within control layer, and only employs RLC TM, thereby essentially directly conveying its RRC messages via MAC. If RRC is not co-located with MAC Scheduling, the latter caches RRC state to avoid per TTI signalling exchanges with RRC.

4.3 Data layer

Data layer network functions are inherently distributed due to their purpose of providing data forwarding end-to-end. There are two options to control data layer NFs: first, they may be controlled directly by control layer NFs as it is done in current mobile networks; second, they may be controlled by SDM-C or SDM-X, which allow for more degrees of freedom of programmability. As introduced in Section 2.1, shared (common) NFs are controlled by SDM-X while dedicated (customized) NFs are controlled by SDM-C.

The corresponding interfaces towards SDMC-controlled data layer NFs 5GNORMA-SDMX-NF and 5GNORMA-SDMC-NF, respectively. Parts of the data layer NFs have no direct interface to SDM-X/C but are instead controlled indirectly through distributed control NFs with proprietary (logical) interfaces "M", "U" and "C". The only exception is *Transport (SDN)*, which extends the classical SDN forwarding plane and accordingly is controlled through interface 5GNORMA-SDMC-SDN. The 5G NORMA functionally decomposed data layer can be further categorized into non-access stratum (NAS) and access stratum (AS), which is then further split into physical layer and link layer.

Physical layer

Physical layer processing is provided by three function blocks PHY TP, PHY Cell and PHY User. For an energy and cost efficient implementation, these blocks are generally realized as physical

³ Except the very first connection setup or reestablishment message, in each direction UL and DL, which is sent via CCCH using RLC TM.

network functions (PNF) employing dedicated hardware. Although the use of PNFs may limit the replaceability compared to VNFs, they may still be highly configurable in order to adopt to the diverse service requirements. The rationale behind grouping into three functions blocks as chosen by 5G NORMA is detailed in the following:

PHY TP subsumes all functions that map one-to-one to a single specific, spatially localized transmission (and reception) point. This includes antenna elements, amplifiers and analogue and mixed signal processing. These functions are sometimes implemented in distinct physical entities (e.g. into antenna(s) and RRH(s) in case of macro sites) but are regularly integrated for sites with lower transmit power (small cells, femto cells and access points in unlicensed spectrum). One PHY TP instance may only host a single PHY Cell instance, but regularly multiplexes several of them (i.e. several carrier frequencies). PHY TP is the only function block that inherently is (implemented as) a PNF. For 5G NORMA, which is investigating "softwarization" and virtualization of mobile networks, PHY TP is therefore out of scope and is provided for completeness only, representing the lowest/last part of the overall processing chain of the mobile network infrastructure towards the UE.

Next, *PHY Cell* subsumes all functions that map one-to-one to a single radio carrier, i.e. what 3GPP calls a physical cell. The main task is the (de)multiplexing of signals and physical channels that constitute a single carrier. Physical signals, e.g. synchronization (PSS, SSS), demodulation (CRS) and channel measurement (CSI-RS) are generated by PHY Cell itself, while each physical channel is generated by a distinct (logical) PHY User instance. Accordingly, one PHY Cell instance almost always multiplexes several PHY User instances, at least for multi-service-capable air interfaces. Today, dedicated hardware is employed but may also be virtualized as in software-defined radio.

Finally, *PHY User* subsumes all functions that map to a specific (group of) users. Here, the largest and most processing intense part of the PHY layer is carried out such as multi-antenna channel equalization/precoding, possibly across multiple PHY Cell (CoMP JT/JP), and FEC. A single PHY User instance may process several physical channels, grouping multiple local PHY User instances together, e.g. a channel for user data transport with its associated control channel. Again, like PHY Cell, PHY User is best implemented as PNF for efficiency reasons but may be implemented as VNF, too.

While the split between PHY TP and PHY Cell is predefined due to the nature of PHY TP, i.e. the split between functionalities that cannot be virtualized and those that can, the split between PHY Cell and PHY User is less strict. We have chosen this split because of the characterizing feature that functions of the former map one-to-one to a single air interface instance (i.e., radio carrier/cell), while the latter do not necessarily do so but rather relate to specific users respectively services. Although we aligned the functional decomposition of the control of data layer along 3GPP EPS, we consider the characterization valid for a multitude, if not all, radio access technologies. While the PHY TP to PHY Cell logical interface is in line with the approach to CRAN (carried out as CPRI or ORI physical interface), the PHY Cell to PHY User interface represents a valid alternative, which still centralizes most PHY processing and offers full flexibility at the benefit of lower bandwidth requirements and simpler switching within the CRAN cloud necessary to realize pooling gains.

The challenge of avoiding many additional interfaces may be addressed by a flexible container protocol on data and control layer. The main benefit of the flexible functional architecture is the possibility to exploit centralisation gains where possible, to optimise the network operation to the actual network topology and its structural properties, and to use algorithms optimised for particular services, i.e., optimise through dedicated implementations instead of parameters.

Link layer

The link layer⁴ or "Layer 2" is provided by the three main function blocks *MAC*, *RLC* and *PDCP*. These blocks are mandatory for user data transmissions of logical channels or radio bearers of Layer 2 in the presented example of 3GPP EPS and, respectively, LTE (E-UTRA) as its air interface. Additionally, there are three optional Layer 2 function blocks *MAC CA* for carrier aggregation, *PDCP Split Bearer* for multi-connectivity incl. multi-RAT, and *eMBMS* for multimedia broadcast and multicast services. These blocks and, respectively, the grouping of functions are quite specific to each RAT and, unlike the grouping in the physical layer, does not claim to be of universal applicability.

Non-access stratum

The access stratum, which includes the physical and link layer function blocks described above, specifically focuses on the realization of data transmissions over the air interface between UE and serving RAN, optimized to service, coverage (radio channel) and deployment characteristics (also cf. Section 3.2 and Annex A). In contrast, the non-access stratum abstracts from the underlying transport. For example, 3GPP EPS simply assumes a secured IP-based transport. Accordingly, communication can be easily and fully virtualized. In WP4, all functions are subsumed in a single function block *NAS*. For NAS functions (with a focus on their control through SDM-C) the reader is referred to WP5 [5GN-D52].

Finally, the three function blocks *MEC Application*, (end user) *Service* and (Packet) *Data Network* are considered. Strictly speaking, these three blocks are no functionalities of the mobile network itself but represent functions that utilize the communication link to the UE. Access to a data network is the most generic end user service provided to a UE, e.g. Internet connectivity and the access to services offered through it. Such (IP-based) services may also be integrated into the mobile network itself, e.g. the IMS (IP multimedia subsystem) for providing VoIP as its prime service. Last, connectivity to MEC (mobile edge computing) as a generic platform for service deployment within a mobile network near to the UE is shown. Details on the integration of MEC can be found in Section 6.5.

4.4 Security considerations

A general discussion of security aspects of the 5G NORMA architecture is given in [5GN-D32] and [5GN-D33]. This section focusses on specific security aspects in the context of the RAN control and data layer described in the preceding sections.

4.4.1 Securing inter-domain interfaces

The essential aspect in this context is the split of the network in common and slice specific parts. According to the 5G NORMA stakeholder models discussed in [5GN-D32] and [5GN-D33], the common parts are typically owned by an MSP (Mobile Service Provider), while slices may be operated by tenants, i.e. MSP customers who rent slices from the MSP and operate them on their own. So we call the interfaces between common parts and the slices "inter-domain interfaces" in the context of this section. Mechanisms discussed in [5GN-D32] and [5GN-D33] ensure that a slice is isolated against other slices. In contrast, the common parts need to interact with multiple slices and thus need to be accessible by slice-specific functions. The degree of control that can be exhibited by a tenant over its slice may vary, depending on the so called "service offer type" (see

⁴ More precisely, only the link layer of the access stratum, namely the radio bearer, is meant. From the (packet) data network respectively network layer point of view, e.g. from the Internet (protocol) point of view, the whole path to the UE represents a single link (layer), the EPS bearer, which is composed of the radio bearer, S1 bearer and S5/S8 bearer (in case of 3GPP EPS).

[5GN-D3.2] section 3.2), but cases where the tenant has a high amount of control and can use its own tailored VNFs, are supported. Consequently, a tenant/slice can use in arbitrary ways the interfaces exposed by the common part, without being restricted by the need to use predefined software images for its own functions.

Besides accessing the exposed interfaces, a tenant slice has no control whatsoever over the common parts. In the other direction, the common parts need not rely on any function provided by slices, and the operator of the common parts, will typically have control also over the slices – in its role as the party renting out these slices to tenants.

The operator of the common parts, the MSP, and the slice operators, the tenants, are contractual partners, with an SLA specifying the service that is offered by the MSP. The tenants need to trust the MSP anyway, as the MSP has technical means to access their slices, at least in all setups where the tenants cannot run their slice-specific functions on tenant-owned infrastructure. In contrast, the MSP cannot rely on tenants to behave correctly at all times. Rather, the MSP must anticipate erroneous or malicious behaviour and design and operate the common parts in a way that they are secure against such misbehaviour of the slice specific parts. This mainly affects the interfaces exposed by the common parts. (Resource isolation for the common parts is also an issue, but this is covered in the general security considerations in [5GN-D32] and [5GN-D33].)

In the following, control and data layer interface security is discussed separately. It is understood that in any 5G NORMA implementation, specific care must be taken that the interfaces in both layers are implemented in a sound and robust way, to minimize the amount of implementation flaws and potential vulnerabilities that may leverage successful attacks on the common parts.

Control layer interfaces:

The control layer interface to the common parts is provided by the SDM-X. While security for the SDM-X interfaces is described in [5GN-D33], it can be noted that state-of-the-art means can be sufficient to secure this interface. Typically, this would comprise the setup of dedicated security associations between the SDM-X and each slice-specific function (i.e., slice-specific SDM-Cs) accessing it. Such security associations need to be based on long term credentials such as private/public key pairs per party, plus certificates that assert the identities of the owners of the public keys. Access to the SDM-X is then protected by the security association (which could for example be implemented by means of TLS). Based on this, the SDM-X can verify the origin of any request and may implement arbitrary methods and policies how to authorize access of slice-specific functions to common functions.

Note that the primary purpose of these security associations is not confidentiality protection (as the traffic is mostly not visible to any third parties), but integrity protection and origin authentication, i.e. making sure that attackers cannot inject faked communication, or that one slice cannot impersonate another slice.

It should be further noted that securing the control layer interfaces may not always require cryptographic protection. For example, different slices may be connected via different VLANs to different ports of a common function, thus forming three isolated connections. However, it seems that such an approach is more fragile and is more prone to errors, such as configuration errors, that may lead to lack of isolation and security.

Data layer interfaces:

The functional architecture described in chapter 2 comprises different options. The location and nature of the data layer interface between common and slice-specific parts varies for the different options. It is beyond the scope of this document to specify all the interfaces in detail, so a detailed description how to secure them lacks the foundation. Still, by means of example we describe in the following how a data layer interface can be secured. For this example, we select the option 2 from chapter 2, where the interface between common and slice-specific parts is between the MAC and RLC layers.

The interface is assumed to be implemented via buffers located in memory that is accessible by both the common MAC function and the slice specific RLC function. For example, there could be one buffer per data radio bearer and per direction (uplink and downlink). The common MAC function ensures that tenants cannot access arbitrary buffers, but only those assigned to the tenant. This could be facilitated by providing addresses in common memory to slices only in form of offsets, relative to some base address that is known to the common MAC function. Clearly, the common MAC function must ensure that offsets used by slices do not exceed the allocated range, i.e. discard any request to access memory at out-of-range offsets.

For downlink traffic, the common MAC function processes all downlink buffers. It understands which buffers belong to which tenant and ensures that the SLAs of all tenants are met, e.g. in terms of the amount of radio resources agreed in each SLA. Similar, in uplink direction, the common MAC function understands which data radio bearers belong to which tenant and processes them in a way that the SLAs are met.

For the protection of the MAC function against abuse by tenants, it is essential that the procedures provided to slice-specific functions for accessing the buffers are implemented soundly, without flaws that could negatively impact the common MAC function. (Moreover, there must not be flaws that would endanger the correct assignment of radio resources to slices.) It is further essential that the common MAC function merely transmits the data, i.e. does not process the data in a way that would allow that maliciously crafted data harm the function itself.

Interfaces between control and data layer:

Basically, interfaces between the control layer and the data layer are control interfaces – the data layer function is controlled by the control layer function and must comprise a "control part" to be able to accept and execute control commands from the control layer. Consequently, the security measures of section 4.4 apply.

The overall architecture does not explicitly comprise inter-domain interfaces between control and data layer, but one could think about generalizations or variants of it that may comprise common data layer functions that accept commands from slice specific control layer functions. Note however, that such a setup may burden the data layer function with the task to distinguish these different domains, to establish and maintain security associations to all of them, and possibly even to perform authorization for requests. Therefore, such a function would need to comprise a significant control part and may thus no longer classify as a pure data layer function, blurring the separation of control and data layer. Consequently, such a function should be split into its control and its data layer part, with the inter-domain interface in the control layer rather than between control and data layer.

4.4.2 Securing intra-domain interfaces

Such interfaces are per definition either inside a single slice, or inside the common parts. As described in [5GN-D32] and [5GN-D33], within NFV environments isolation mechanisms are available that ensure that such internal interfaces are not accessible from the outside. For instance, internal functions of a network slice may be interconnected by an isolated virtual network spanning virtual switches inside hypervisors, but also parts of the physical network consisting of top-of-rack and other switches in data centres, as well as WAN switches and connections. Various existing techniques may be applied to achieve isolation, such as VLANs implemented by virtual switches or by physical Ethernet switches, or virtual routing and forwarding in IP routers, or isolation by means of dedicated MPLS label switched paths in WANs.

Assuming proper isolation of slices and mutual trust of functions inside a single domain (i.e. inside a slice or inside the common parts), there is no general need for establishing security associations and cryptographically protecting traffic at all interfaces. Indeed, using crypto at all interfaces would not reasonably scale in setups where messages or data packets pass chains of functions inside VNF environments and would require multiple encryption and decryption on their way through the network.

Cryptographic protection is however required when the communication uses physically exposed transmission media, such as fibres interconnecting distributed data centres, or any type of backhaul links, wired or wireless. However, such cryptographic protection need not be provided on a per interface base, separately in each domain. Rather, it may apply for the aggregated traffic on a physically exposed link. As an example, the operator of a distributed cloud infrastructure may apply wholesale encryption on the optical layer on all the fibres interconnecting data centres.

Assuming such a secure, distributed cloud infrastructure, intra-domain interfaces between VNFs are isolated and protected by default, without specific per domain measures. Interfaces involving PNFs are not covered by this. Here, it depends on the nature of each single interface and PNF, how the interface needs to be secured.

As an example, in the common RAN parts, there may be fronthaul interfaces between virtualized RAN functions and remote radio heads. User and control layer traffic to and from mobile devices on these interfaces may be protected anyway by security associations between the mobile devices and the PDCP function inside the network. If management communication between network and remote radio head is required, it may be protected by dedicated cryptographic management protocols.

5 Multi-tenancy

5.1 Multi-tenancy aspects of RRM

5.1.1 Multi-tenancy radio-resource management

Driven by the capacity requirements forecasted for future mobile networks as well as the decreasing margins operators are able to realize, infrastructure sharing is emerging as a key business model for mobile operators to reduce the deployment and operational costs involved in initial roll-out (capital expenditure (CAPEX) and operational expenditure (OPEX) of their networks).

The issue of dynamically sharing resources between operators has received substantial attention both from industry and standardisation as well as in the research community.

Network sharing solutions are already available, standardised, and partially used in some mobile carrier networks. These solutions can be divided into passive and active network sharing: passive sharing refers to the reuse of components such as physical sites, tower masts, cabling, cabinets, power supply, air-conditioning, and so on; active sharing refers to the reuse of backhaul, base stations, and antenna systems, and it's labelled as active radio access network (RAN) sharing. However, these sharing concepts are based on fixed contractual agreements with mobile virtual network operators (MVNOs) on a coarse-grained basis (monthly/yearly).

3GPP has recognised the importance of supporting network sharing since Release 6, and defined a set of architectural requirements and technical specifications that have been continuously extended since then. The latest activities have focused on the definition of new sharing scenarios and requirements [22.101], and the corresponding network management architecture and functionality extensions towards on-demand capacity brokering [32.130].

Based on the enhanced capabilities of dynamic sharing, new business models for infrastructure owners are expected to emerge, resulting in new revenue sources. Indeed, such an approach supports not only classical players (mobile operators) but also new ones such as Over-The-Top (OTT) service providers that may buy a share of a mobile network to ensure a satisfactory service for their users. Think, for instance, of Amazon Kindle support for downloading content from anywhere, or pay TV sports subscriptions including a premium for watching live games, to name just two examples.

However, substantial attention has been devoted to the architectural framework for multi-tenancy, but relatively little work has focused on the design of criteria and algorithms for this purpose. While some algorithms and criteria have been proposed in the literature, these either fail to meet the requirements for a practical solution or rely on criteria that have been proposed without proper justification on their optimality.

In order to extend the network sharing 3GPP standardisation to meet the 5G requirements, our contribution involves the design of:

- New sharing criterion
- New sharing mechanism

3GPP standardisation provides a static allocation that guarantees a minimum level of resources and limits the maximum amount of resources allocated to a tenant. These resources could be available for a specified period of time and a certain location. In the context of 5G networks, our goal is to design a new sharing criterion that allows for allocating the resources among tenants in a more flexible way.

The idea is to design a criterion that maximises the network utility while at the same time

• allocating computational resources fairly among operators;

- allocating computational resources of each operator fairly among its users; and
- taking into account the computational complexity required by each user's transmission through the modulation and coding scheme (MCS) selection.

The information involved include channel capacity (number of resource blocks of each eNB), the number of users in the network, the channel quality of each user and so the MCSs available.

When multiple tenants decide to share a network, they share all the network resources including the computational ones. Each mobile network is characterised by a well-defined set of functionality, timing requirements, and protocols. This imposes very precise requirements on the operation of each base station, including data processing requirements in order to maintain real-time properties. 3GPP LTE defines a set of MCSs whose choice depends on the signal-to-interference-plus-noise ratio (SINR) so, there is an inherent relationship between the channel quality and the computational requirements.

In future mobile networks, a significant deployment of small cells is foreseen. Since the local processing capability of each small cell is anticipated to be far less than that of a macro cell, the dimensioning of data processing resources needs to be revisited. Furthermore, very dense networks are subject to dramatic temporal and spatial traffic fluctuations such that many small cells may be strongly underutilized. It is therefore not economically viable to equip a small cell based on peak data processing requirements, yet under-dimensioning the computational resources limits its capabilities. Thus, in the design of a criteria for multi-tenancy radio resource management, also the computational complexity required in order to process each user's transmission has to be taken into account.

We define the network utility as the sum of operator's utilities (which depend on the users' allocation x and the fraction f of resources allocated to a user) weighted by the operator's share as:

$$W(x,f) = \sum_{o \in O} s_o U_o(x,f)$$

where s_o is the network share of operator o, in terms of resources assigned to each operator. The operator utility is given by

$$U_o(\mathbf{x},\mathbf{f}) = \frac{1}{|U_o|} \sum_{u \in U_o} \log \left(R_u(\mathbf{x},\mathbf{f}) \right)$$

where R_u is the average throughput of user u of the operator o, and U_o the set of users belonging to operator o.

With the above we can formulate the multi-operator computationally-aware optimisation problem that will drive the scheduling as follows: at each TTI a central scheduler decides (*i*) the allocation of users to resource blocks of the associated eNB, and (*ii*) the users' MCS, that will determine the users' transmission rate. The goal is to maximize the network utility that corresponds to maximizing performance in terms of proportional fairness. The scheduling decision is subject to the following constraints:

- the sum of all operator's users' computational complexity cannot exceed the share of computational resources assigned to each operator;
- in each TTI the aggregated computational load cannot exceed the computational capacity;
- a user can use only one MCS in all resource blocks allocated to him;
- each user can be associated with at most one eNB per TTI.

The proposed criterion allocates resources across operators dynamically and fairly, tracking changes in the numbers and locations of operators' mobile users and the associated transmission rates and so the computational complexity required.

5.1.2 Reinforcement learning for network slice resource management



Figure 5-1: Reinforcement learning for slice admission control

While the flexibility brought into multi-tenant systems with the network slicing concept pushes for a rapid network virtualisation evolution, infrastructure providers do not quantify the tangible benefit on their current business cases. Assessing and brokering network slicing operations appears to become crucial while developing new architectures supporting network slicing. In this direction, 3GPP has already standardised a centralised entity applying admission control policies to incoming network slice requests, acting as capacity broker [22.852] and residing within the infrastructure provider's network. In [SCS16] it has been further enhanced in order to map incoming slice request SLA requirements to wireless physical resources, e.g., resource blocks. In this way, tenants can directly obtain a "slice" of the radio access network (RAN) elements. Although conservative mappings may be used for mission critical services (that need ultra-high availability), enhanced admission control algorithms that leverage multiplexing gains of traffic among slices are key to the optimisation of network utilisation and monetisation. To this end, the ability to predict the actual footprint of a particular network slice is essential to increase the maximum number of slices that might be run on the same infrastructure. Building on this idea, we designed and interconnected three building blocks, as depicted in Figure 5-1: (i) a learning module in charge of predicting network slices' traffic based on past traffic and user mobility, (ii) an admission control policy and (iii) a slice traffic scheduler in charge of fulfilling the agreed SLAs and feeding back (reinforcement) anomalies to the traffic prediction module.

Slice traffic forecasting:

k	$\mathbf{T}^{(\mathbf{k})}$	Туре	QCI
0	10 ms	GBR	-
1	50 ms	GBR	3
2	100 ms	GBR	1
3	150 ms	GBR	2
4	300 ms	non-GBR	6
5	1000 ms	non-GBR	-

Table 5-1: Traffic Class Requirements (similar QCI from [23.203])

Traffic predictions are computed on an aggregate basis for every tenant. Each tenant *i* might ask for a different network slice request $\sigma_i^{(k)}$ tailored for specific service requirements. Indeed, the forecasting process can easily categorise the traffic requests based on related service requirements, thereby performing a prediction separately per slice. We assume different classes of traffic based on specific SLAs as shown in Table 5-1. We denote the traffic volumes of tenant *i* for traffic class *k*, e.g., satisfying given service requirements, as a realisation of a point process $\sum_T \delta_t r_i^{(k)}(t)$, where δ_t denotes the Dirac measure for sample t. We express traffic requests $r_i^{(k)}(t)$ in terms of required physical resources but they can be easily translated into different metrics, such as latency or throughput demands while applying the same algorithmic approach. We use the Holt-Winters (HW) forecasting procedure [KSO01] to analyse and predict future traffic requests associated to a particular network slice. We rely on the additive version of the HW forecasting problem as the seasonal effect does not depend on the mean traffic level of the observed time window but instead it is added considering values predicted through level and trend effects. After properly setting the HW parameters (α , β and γ), we define the one-step forecasting error which can be obtained during the training period of our forecasting algorithm, i.e., when predicted values are compared with the observed ones. We can then derive the prediction interval wherein future traffic requests lie for that particular network slice with a certain probability. Due to the penalties imposed by traffic SLAs, we focus only on the upper bound of the prediction interval, which provides the "worst-case" of a forecasted traffic level. Interestingly, a larger prediction time window results in a reduced accuracy making the system behaving closer to the real network slice demand with limited multiplexing gains. Conversely, an accurate forecasting with a lower error probability can result in higher gains while still guaranteeing the traffic SLAs. Therefore, we adjust the forecasting error probability according to the service requirements and to the number of prediction points the forecasting process needs to perform. For instance, best effort traffic requests having no stringent requirements can tolerate a prediction with a longer time pace resulting in unprecise values. Hence, we might select for this service type a low forecasting error probability. On the other hand, when guaranteed bit rate traffic is considered, the corresponding SLA must be fulfilled in a shorter time basis, which makes our forecasting process much more complex requiring significantly more predicted values. Therefore, our system models such a type of traffic with a higher forecasting error probability. Finally, forecasting error probability values are monitored and adjusted through a reinforcement learning process, based on the SLA violations experienced during the scheduling phase.

Slice admission control:



The 5G°NORMA Inter-slice Broker, together with SDM-X, might decide on the network slice requests to be granted for the subsequent time window based solely on the current resource availability. However, if forecasting information is considered, network slice requests might be accurately reshaped to fit additional slice requests into the system, as shown in Figure 5-2.

Let us assume a rectangular box with fixed width W and height H representing the resource availability within a fixed time window. Let us assume a set of items I, where each item i corresponds to a network slice request having width u_i corresponding to slice duration L_i and height h_i corresponding to the amount of resources R_i . In addition, each item is provided with a profit c_i corresponding, in our case, to the amount of resources needed. This assumption relies on the fact that every slice request pays the same amount of money proportional to the number of resources granted. The objective of the admission control problem is to find a subset of items which maximises the total profit, e.g., the total amount of used resources, as shown in Figure 5-2.

In this illustrative example, different amounts of needed physical resource are forecasted for a single network slice request. It may be observed that when the forecasting phase is accurate, more room can accommodate more slices, as the slice 6 admitted into the system. Please note that in our case the (flexible) geometric two-dimensional knapsack problem is constrained by the orientation law of the considered items. In particular, each item i has a fixed orientation, which cannot be changed to fit in the box. We can formulate our admission control problem as follows:

Problem ADM-CONTROL:

maximize
$$\sum_{i \in \mathcal{I}} c_i \cdot x_i$$

subject to
$$\sum_{i \in \mathcal{I}} w_i \cdot x_i \leq W; \quad \text{(relaxed)}$$

$$\mathcal{S}(x_i) \cap \mathcal{S}(x_k) = \emptyset, \quad \forall i \neq k;$$

$$\mathcal{S}(x_i) \subset \mathbb{S}, \quad \forall i \in \mathcal{I};$$

$$x_i \in [0, 1], \quad \forall i \in \mathcal{I};$$

where $S(x_i)$ depicts the geometrical area of the item *i* (either rectangular or irregular defined) whereas S is the area of the box. The first constraint refers to the weight of each item (w_i) . For the sake of simplicity, we consider the weight capacity of our box as infinite to neglect the item weight. The next two constraints state that items cannot overlap with each other and must be contained within the total space of the box. The solution of such a problem provides a set of x_i , which is a binary value indicating whether the item *i* is admitted into the system or rejected for the next time window.

Heuristics:

Algorithm 1 Network Slices Packer: Algorithm to admit network slice requests $\sigma_i^{(k)}$ within the system capacity Θ for the next time window T_{WINDOW} .

Input: $\Sigma = \{\sigma_i^{(k)}\}, \Theta, T_{\text{WINDOW}}, \mathbb{S}$ Initialization: $\mathcal{C} \leftarrow \emptyset, \mathcal{F}_1 \leftarrow \emptyset, \mathcal{F}_2 \leftarrow \emptyset, \mathcal{E} \leftarrow \emptyset$ Procedure 1: for all $C_l \leftarrow {\binom{\Sigma}{2}}$ do 2: if C_l fits into \mathbb{S} then 3: $\mathcal{C} \leftarrow \mathcal{C} \cup C_l$ end if 4: 5: end for 6: for all $C_l \in C$ do $\{v(C_l \cup B_l), s(C_l \cup B_l)\} \leftarrow$ Solve the knapsack problem $P(C_l)$ 7: 8: end for 9: $l* = \arg \max_{l \in \mathcal{C}} \{ v(C_l \cup B_l) \}$ 10: if $v(C_{l*}) \ge \frac{v(C_{l*} \cup B_{l*})}{2}$ then 11: return C_{l*} 12: else $\mathcal{F}_1 \leftarrow C_{l*}$ 13: $\mathcal{F}_{2} \leftarrow B_{l*}$ if $s(\mathcal{F}_{1}) \geq \frac{|\mathbb{S}|}{2}$ then return B_{l*}^{l*} 14: 15: 16: 17: else 18: Sort \mathcal{F}_2 in non-increasing order of their profits and traffic class k while $s(\mathcal{F}_1) < \frac{3}{2}$ do 19: $e = pop(\mathcal{F}_2)$ 20: $\mathcal{F}_1 \leftarrow \{\mathcal{F}_1 \cup e\}$ 21: end while if $v(\mathcal{F}_2) \ge \frac{v(C_{l*} \cup B_{l*})}{v(\mathcal{F}_2)^2}$ then return $v(\mathcal{F}_2)^2$ 22: 23: 24: 25: else $\mathcal{E} \leftarrow \max\{v(\mathcal{F}_1 \setminus e); v(\mathcal{F}_2)\}$ 26: return *E* 27: 28: end if 29: end if 30: end if

Figure 5-3: Pseudocode of the admission control algorithm

We assume rectangular shapes for network slice requests with different traffic requirements. Considering the traffic classes introduced in [23.303], when traffic class k = 0 the regular shape of the network slice is hardly defined and no flexibility is allowed for allocating the traffic requests. Conversely, when less-demanding slice requests k > 0 are considered, the slice might be reshaped, delaying the slice traffic, to efficiently fit into the network. We rely on the assumption that each tenant is not allowed to ask for more than half of the available resources of the infrastructure provider. This implies that at least 2 network slices can be accommodated. The algorithm is listed in Figure 5-3. Among all possible pairs of network slice requests, only those fitting the available system capacity are taken into account. For each 2-slice set we formulate a 0-

1 knapsack problem to maximize the total profit assuming a single weight (the area of the slice) per item (line 7). The item set to evaluate for the knapsack problem includes the 2-slice set and all the other slices, while considering C_l as already allocated slices. Based on the Fully Polynomial Time Approximation Schemes (FPTAS) proposed in [KP99], we retrieve the best solution, i.e., a set of network slice requests among all knapsack problems. If the total profit v(.) assigned to the 2-slice set requests C_l is greater than the half of the best profit retrieved after running all knapsack problems, we keep C_l as the best feasible set (line 10). Otherwise, we split the optimal set into two subsets F_1 and F_2 . If the total space (s(.)) covered by the items in F_1 is greater than the half of the optimal solution (line 24). Therefore, the subset F_2 could be packed into the system capacity in polynomial time. Otherwise, we move the item with the greatest profit and the highest traffic class k from F_2 to F_1 until the space of F_1 is greater than the half of the system capacity. Then, if the total profit of F_2 is greater than half of the optimal one, the algorithm ends and we keep F_2 as the optimal set. Otherwise, we choose as output the set providing the best total profit after comparing F_2 , without the latest added element, with F_1 .

Slice traffic scheduling phase:

We generalize the scheduling model for accounting different traffic SLAs. We assume a traffic request from tenant *i* for traffic class *k* as $r_{i,z}^{(k)}$. We consider 6 traffic classes. Each traffic class is characterized by a time window *z* identifying the offset between two consecutive resource requests, shorter for high-demanding traffic requirements and larger for best effort class. The scheduler ensures that the whole amount of required resources is served for any given time window. The key-objective of this novel network slice traffic scheduler is to minimize the amount of resources scheduled while guaranteeing the traffic SLAs within a network slice. When forecasted information is available, the scheduler expects slice traffic levels below the predicted traffic $\hat{R}_{i,z}^{(k)}$ bounds such that $r_{i,z}^{(k)} \leq \hat{R}_{i,z}^{(k)}$. If forecasted traffic bounds are underestimated and the traffic demands exceed the expected values, traffic requests are automatically capped at the original amount of resources agreed during the slice request admission, i.e., $R_{i,z}^{(k)}$. Hence, slice slocations may overlap and traffic class requirements might not be fulfilled incurring in slice SLA violations. We model the scheduler problem as a general minimization problem addressing any traffic class SLA and providing $s_{i,j}^{(k)}$ as the amount of resource served per time *j* upon the list of admitted slices $x_i^{(k)}$ is available from the admission control phase. We introduce the scheduled traffic representing the real amount of resources assigned per time *j* upon the list of admitted slices is available from the admission control phase. The problem is formulated as follows:

$$\begin{array}{ll} \text{minimize} & s_{i,j}^{(k)} \\ \text{subject to} & \begin{pmatrix} zk + \bar{t} + T^{(k)} \\ \sum \\ j = zk + \bar{t} \end{pmatrix} \geq r_{i,z}^{(k)} x_i^{(k)}, \ \forall z \in \left[0, \left\lceil \frac{L_i}{T^{(k)}} \right\rceil - 1\right]; \\ & \sum_{i \in \mathcal{N}} s_{i,j}^{(k)} \leq \Theta + P_{i,j}^{(k)}, \quad \forall j \in \mathcal{L}; \\ & s_{i,j}^{(k)} \in \mathbb{R}_+, \qquad \forall i \in \mathcal{N}, j \in \mathcal{L}, k \in \mathcal{K}; \end{array}$$

where θ is the total capacity of the system expressed as the total amount of resource blocks whereas $P_{i,j}^{(k)}$ is the penalty incurred for not having satisfied a particular tenant slice traffic SLA, namely SLA violation. The network slice scheduler keeps track of SLA violations to promptly trigger dynamic forecasting parameters adjustments.

5.2 Multi-tenancy in multi-RAT environments

On-demand network sharing provides a new degree of flexibility for multi-tenancy systems compared to the first generation of network sharing concepts, which were based on long-term contractual agreements.

Resources are acquired on a short-term scale (minutes) leaving the actual allocations to signalling feedbacks. The synchronisation in resource sharing is guaranteed by a central resource management entity, which is represented by the capacity broker, within the MNO infrastructure. A tenant request reaches the capacity broker, which has a global view of the network resource utilisation. Based on such information, the capacity broker decides whether to accept or reject the tenant request aiming at optimizing the resource utilisation while maximizing the overall profits.

In order to alleviate the spectrum scarcity problem, the future 5G networks will leverage on multiconnectivity supporting simultaneous connectivity across different technologies such as 5G, 4G, and Wi-Fi, multiple network layers, such as macro and small cells, and multiple RATs. This introduces a higher complexity in the management of the resources because the different layers and radio access technologies present different characteristics.

Now we have to consider more possible allocation solutions because we share not only the spectrum resources but also the different technologies. This obviously introduces more complexity but even more flexibility, i.e., our algorithm has to decide not only about rejecting and accepting a request but even which technologies (among the available ones) are the most suitable for the service that the tenant wants to deliver to its users in respect of the QoS requirements (we could assign to a tenant even different technologies simultaneously). For these reasons the algorithms developed and illustrated in the following paragraphs are agnostic to the technologies utilised and so they can perfectly work independently of the different technologies shared.

In order to design new algorithms for multi-tenant approaches we have focused on 2 different scenarios:

- 1. Algorithm for handling resource requests at the infrastructure provider;
- 2. Algorithm for dynamic resource sharing among operators.

5.2.1 Admission control

One of the key novel concepts of the 5G architecture is network slicing: the infrastructure can host different slices each of which can provide different services. This opens the mobile network ecosystem to new players:

- infrastructure Provider (InP), which operates the infrastructure;
- Mobile Service Provider (MSP), which offers the (mobile) telecommunication service (realized by a network slice); a Mobile Network Operator (MNO) can be considered to combine the roles of InP and MSP; and
- tenants, which acquire a network slice from the MSP to deliver a specific applicationlevel service to own subscribers.

In this new ecosystem, MSPs issue to the InP requests for spectrum and computational resources in order to set up their slices, which are finally used by subscribers of the tenant. Since spectrum is a scarce resource, for which overprovisioning is not possible and its availability heavily depend on SLAs and users' mobility, the InP cannot apply an "always accept" strategy for all the incoming requests from MSPs. In the same way, MSPs cannot serve all incoming requests from tenants. Thus, the new 5G ecosystem calls for novel algorithms and solutions for the allocation of network resources among different tenants.

Our idea is to design a network capacity brokering algorithm executed by the MSP in order to decide whether to accept/reject a request from a tenant with the goal of maximizing the MSP revenue, satisfying the service guarantees required.

In our model the MSP receives requests from tenants characterised by:

- amount of resources to be reserved,
- starting and end times for the reservation,
- type of traffic (elastic or inelastic) that imply,
- required quality of service/SLA, and
- its price ρ (amount of money per time).

The MSP needs to decide which requests to accept knowing that by accepting a request with a small bid, it may lose future opportunities to involve higher bids (not enough resources available), but rejecting a request the MSP loses the corresponding bid. Our algorithm leverage on semi-Markov decision process (SMDP) theory, that models the resource allocation to network slices as a Markov chain in which the next state depends only on the actual state, the decision taken and the transition probability function.

The algorithm requires full knowledge of the system parameters like

- the inter-arrival time λ of requests,
- the request duration μ ,
- the transition probability function,

and that the system is memory less. By applying decision theory, it is possible to find the decision policy that maximises the MSP's revenue. While SMDP provides the optimal policy, it implies very high computational cost as the state space is large, so for practical purposes we need an adaptive algorithm, but we will use SMDP as a benchmark to evaluate the performance of the proposed adaptive algorithm.

The adaptive algorithm is based on Q-learning, a reinforcement learning tool that learns about the system behaviour by taking non-optimal decisions during the learning phase. After an initial learning phase, by evaluating the best possible action starting from a certain state, the algorithm is able to find the decision policy that maximises the MSP's revenue. The best advantages of this tool are:

- it does not need any knowledge of the system parameter (λ , μ and transition probability function);
- it works even if the system is not memory less; and
- it's an online algorithm that can react to system perturbations (it just needs a short learning phase).

In order to evaluate the adaptive algorithm, we compared, for various price ratios ρ_i/ρ_e of inelastic to elastic traffic, the revenue obtained by applying the policy derived by Q-learning with

- the revenue obtained applying the SMDP approach,
- the revenue obtained by accepting all requests (we reject them only if there are no resources available), and
- the revenue obtained by rejecting all elastic requests.

We simulated two classes of incoming requests: inelastic that demands a certain fixed throughput which needs to be always satisfied with a fixed outage probability, and elastic that requires average throughput guarantees. Each class has different traffic parameters (λ , μ).



Figure 5-4: Revenue vs ρ_i/ρ_e

As we can conclude from Figure 5-4, SMDP (implemented through Value Iteration algorithm) always converges to the optimal policy, the one that provides the maximum revenue for the MSP for each experiment. Furthermore, we can see that with our adaptive algorithm we can obtain close to optimal performance even without the knowledge of all system parameters.

5.2.2 Dynamic resource sharing

Another important challenge in multi-tenancy is the definition of a sharing criterion and the design of an algorithm that follows it in order to enable statistical multiplexing of spatio-temporal traffic loads.

With the definitions given in section 5.1.1,we can formulate the Multi-Operator Resource Allocation (MORA) optimisation problem as follows:

$$\max_{\mathbf{x},\mathbf{f}} W(\mathbf{x},\mathbf{f}) \coloneqq \sum_{o \in O} \sum_{u \in u_o} w_u \log(r_u(\mathbf{x},\mathbf{f}))$$
s. t:
$$\begin{cases} r_u(\mathbf{x},\mathbf{f}) = \sum_{b \in B} f_{ub} x_{ub} c_{ub}, & \forall u \\ \sum_{b \in B} x_{ub} = 1, x_{ub} \in \{0,1\}, & \forall u \\ \sum_{u \in u_0} f_{ub} x_{ub} \le 1, \quad f_{ub} \ge 0, & \forall b \end{cases}$$
(5.2-1)

where the second equality and the last inequality correspond respectively to *user association* and *base station resource allocation* constraints, and w_u is the user weights.

The proposed criterion allocates resources across operators dynamically, tracking changes in the numbers and locations of operators' mobile users and the associated transmission rates. Furthermore, the MORA criterion satisfies some desirable properties both in the way base stations' resources are allocated to associated users, and the way users are associated with base stations:

- 1. given fixed user associations, MORA allocates base station resources to the associated users proportionally to their weights;
- 2. the resulting resource allocation is Pareto-optimal, which means that if under some other user association choice a user sees a higher throughput than that under MORA then there must be another user which sees a lower throughput allocation;

3. MORA is not harming any operator for the global benefit.

Compared with the static slicing (SS) approach, where each operator contracts for a fixed slice of the network resources at each base station for its exclusive use, MORA provides a higher overall network utility and a higher operator utility for a given user association. For different user associations, there may be cases in which an operator sees a higher utility under SS than MORA, but the additional utility cannot be more than log(e), i.e. has an upper bound [CBV+16]. Another important result is represented by the capacity saving resulting from operators sharing infrastructure: sharing the infrastructure with MORA dynamic sharing provides a capacity saving that will be highest when infrastructure is shared by a large number of operators, each with a small number of users per base station. With current trends towards small cells, the number of users per base station is expected to be small, suggesting that infrastructure sharing may be particularly beneficial.

The optimisation problem underlying MORA is a *non-linear integer programming problem*, which can be shown to be NP-hard, so an algorithm that provides the exact solution is not feasible.

Thus, we developed an approximation algorithm which performance is close to the optimal one, is semi-online (trigger a reassociation of a limited number of users upon a user joining, leaving or performing a handover) and distributed (due to the amount of information involved including the channel quality of each user).

In designing the algorithm, we need to decide

- where the users should be (re)associated,
- in which order they should be reassociated and
- how many reassociations are needed.

The proposed algorithm named Greedy Local Largest Gain (GLLG) is a modification of Distributed Greedy (DG), a simple distributed greedy algorithm that requires too many handovers and incurs too high overhead. In particular, with GLLG,

- the reassociation is done based on a largest gain policy, i.e., reassociations are needed because using an online algorithm (upon a user joining the network, it only decides how to associate the new user, without triggering any reassociations of existing users) the performance can be arbitrarily bad;
- the number of handovers is limited by a parameter *m* that represents the maximum number of handovers allowed, in order to meet the best trade-off between the performance of the algorithm and re-association overhead; and
- the eligible users to be reassociated is restricted locally to the ones within two base stations (the ones involved in the previous reassociation).

In terms of network utility our approach performs very close to the benchmark given by a centralised algorithm and DG and it outperforms static slicing (SS) very substantially. The results are shown in Figure 5-5.



Figure 5-5: Utility gains for different approaches as a function of network size

We can conclude that dynamic resource sharing among tenants can be very beneficial. We have then devised a practical algorithm with limited complexity and overhead and which performance is very close to the optimal one.

5.3 Security considerations

5.3.1 Isolation between multiple tenants

When a network supports multiple tenants by creating tenant-specific network slice instances, it is an obvious requirement to isolate these slice instances in a way that one tenant is not aware of the other tenants and has no means to access or even modify information in the other tenants' slices. In NFV environments, this type of isolation is a basic feature that also includes the capability to limit the resource usage of each slice instance in a well-defined way, to prevent a tenant from using up so many resources that other tenants cannot get resources anymore and thus experience a denial of service (DoS).

Tenant isolation in NFV environments is somewhat endangered by vulnerabilities in the NFV software, for example in hypervisors. Assuming however that the relevant NFV software is designed, implemented, configured and operated with highest care to minimize the number of errors and thus the vulnerability, tenant isolation can be achieved in NFV environments (i.e. with respect to the RAN: in the edge cloud and in NFV environments at access points).

Multi-tenancy is not restricted to infrastructure that provides an NFV platform, but also affects "bare metal" RAN equipment. Depending on the nature of the equipment, it may or may not be aware of the different tenants. In the former case, equipment specific mechanisms need to facilitate multi-tenancy and provide proper isolation. For example, a radio scheduler implemented on bare metal equipment may be configurable to ensure certain amounts of radio resources for each of several different slices. (Note that such radio scheduler implementations exist already today to support RAN sharing between different operators.) Naturally, the radio scheduler will not mix up data belonging to different radio bearers, so it maintains isolation between those radio bearers and consequently between the different slice instances.

5.3.2 Trust relationships

A likely multi-tenancy scenario is that a mobile network operator owns the network infrastructure and provides individual network slice instances to tenants. But there are also other models, including scenarios with the use of tenant-owned radio infrastructure. The trust relationships between the different stakeholders involved in such scenarios, and the necessary security measures to protect the different stakeholders in those scenarios are analysed in [5GN-D32] and [5GN-D33]. Trust among multiple tenants on a NFV infrastructure is however not required – with proper isolation, for each tenant the other tenants are invisible.

6 Multi-technology architecture in HetNets

6.1 Multi-connectivity algorithm

6.1.1 The Inter-RAT Link Controller

It is anticipated that future networks will be heterogeneous consisting of several small cells and legacy systems that are seamlessly integrated and delivering data cohesively. Multi-connectivity has recently been a hot topic of research due to its ability to enhance throughput, capacity, coverage and reliability, thereby increasing the overall Quality of Service (QoS). Multi-connectivity refers to a situation where a device's radio modem(s) connect to at least two different access points (AP) associated to different radio bands, as well as different Radio Access Technologies (RATs).

LTE-A already supports the concept of dual connectivity that allows simultaneous connection of a UE to two eNBs, Master eNB (MeNB) and Secondary eNB (SeNB), which are connected to each other via the X2 interface. According to the LTE-A standards, data flow in dual connectivity is split at the MeNB with some data transmitted from MeNB to the UE, while the remaining data is transferred over the X2 interface to the SeNB and then transmitted to the UE via the SeNB.

In the context of this work, we extend the concept of multi-connectivity to small cells without the use of X2 interface in order to achieve low latency and high reliability. In this regard, we introduce a new function block in the AP called Inter-RAT Link Controller (IRLC) for controlling packet duplication and distribution at the PDCP layer. In 5G NORMA, the IRLC is implemented by the *RAT/Link Selection* subblock of the *RRC User* distributed control function block. The purpose of IRLC is to identify and assign a set of RATs to a UE that maximizes the system throughput while simultaneously satisfying the UE service requirements.



6.1.2 Slice specific constraints

Figure 6-1: Slice specific requirements and UE mapping

The 5G services have varying QoS requirements. For example, the mMTC service has very strict requirements on coverage but is lenient with respect to latency and data rate constraints. On the other hand, V2V communication requires ultra-low latency and ultra-high reliability. Hence, the traditional "one-size-fits-all" approach to wireless networks for all use cases and services to every device is no longer viable.

In order to support the new communication demands, mobile network architecture must evolve from the current network of entities to a network *of capabilities* architecture. Network slicing is an end-to-end logical functionality of the mobile network architecture which, among others,

enables operators to provide networks as a service. The network slices will be implemented on the network operators' physical resources classifying the slices based on the service requirements, as well as satisfying the basic business needs. This implies that the UE in the network will itself have varied QoS requirements. To satisfy these per user QoS requirements while simultaneously utilizing the network resources efficiently, we classify the users and map them accordingly to the available slices. Therefore, the QoS of every user will be dependent on the network slice it belongs to. Also in the heterogeneous networks, exploiting all the available RATs to support QoS requirements is a promising technique. According to the requirements, users can be connected simultaneously to multiple RATs.

Figure 6-1 depicts an example where two network slices with different requirements are implemented in the same infrastructure. Network slice 1 requires very high reliability but low availability; network slice 2 has high reliability and availability but is less time critical. While providing multi-connectivity to the user, it is necessary to consider the reliability, latency and availability constraints of that UE.

Currently the RAT selection function is responsible for initiating handover if another RAT appears to have a stronger radio channel than the current one. The RAT selection function consists of two stages: initiation and decision making. In the initiation phase, the RAT selection function gets periodic inputs, such as measurement data, available bandwidth and user preferences. In fact, there exist many algorithms such as context-aware multiple attribute decision making (CAMADM), fuzzy logic algorithms, radio base station efficiency algorithms and Access Network Discovery and Selection Function (ANDSF) rules that contribute to decision making of the RAT selection function. Based on these algorithms, the RAT that suits best according to these algorithms is selected and handover is initiated.

Contrarily, according to the proposed algorithm, the connection of a UE to multiple RATs will be governed by a controller. To execute the service request, the Inter-RAT link selection function will select the set of links to which a UE can be connected to. Once the multi-connectivity link is established between AP and user, the controller selects the operating mode. In fact, the controller decides either to send duplicate data across multiple RATs to increase the reliability by redundant links or to split the traffic across RATs so as to increase the throughput and essentially decrease the end-to-end latency of the service.



Figure 6-2: Data duplication mode (left) and data split mode (right)

The two different modes are depicted in Figure 6-2. The UE is connected with two links, one to AP1 and the other to AP2. In addition, the UE is assumed to have a multi-connectivity support and can decode packets from both links simultaneously. On the left-hand side of Figure 6-2, data packets are duplicated and then transmitted to the UE over both links in parallel. As the data is transmitted over two separate links, the probability of failure dramatically decreases thereby hence the reliability increases. However, the time required to receive the last data packet, in first approximation when neglecting the positive effect of less retransmissions, is the same as if only a single link is used.

On the other hand, if a single link is able to provide the required reliability, in order to reduce the end-end latency of delivering a packet of given size, the data packets are distributed by PDCP and sent over multiple APs simultaneously. The packets are received and reordered at the receiver for decoding. In this approach the data is inversely multiplexed over multiple links as shown in the right-hand side of Figure 6-4. This implies that the links are not redundant anymore, i.e., the failure over even one link will affect the whole transmission and decrease the reliability.

It is important to note that, in the ongoing 5G standardisation discussions, reliability and latency are the two key indicators for the application of new services. The concept of multi-connectivity is used so that such reliability and latency constraints are satisfied. The focus of the designed algorithm is to provide reliability to users efficiently. The end-end reliability of any service depends on the selected links and is highly affected by the channel conditions, including fading and shadowing, as well as the amount of resources allocated to the user. As of now, research has only been done to study multi-connectivity within a single RAT. In our work, however, a centralized controller is introduced that supports multi-connectivity across different RATs.

6.1.3 **Proposed algorithm**

The proposed algorithm consists of two steps. Initially, based on the received signal strength, the users are assigned to the k best APs. Based on the reliability constraints of the UE, an operating mode is selected in the second step as shown in the flow diagram described in Figure 6-3 below. Further details on the overall process for inter-RAT link selection can be found in Part II Section 13.



Figure 6-3: Flow diagram of the proposed algorithm

To identify the optimal number of links per UE as per the reliability constraints, we designed an optimization problem for the cell assignment as

$$\operatorname{arg\,max}_{n_{u,b}} \sum_{b=1}^{B} n_{u,b} \log(1+\gamma_{u,b})$$

s.t $1 \le \sum_{b=1}^{B} \lceil n_{u,b} \rceil \le M \quad \forall u = 1 \dots U$
 $1 - R_u \le \epsilon_u \quad \forall u = 1 \dots U$

In the optimization problem above, the parameters used are defined as follows:

 $\gamma_{u,b}$: SNR that can be achieved over the link between user *u* and AP *b*,

 $n_{u,b}$: assignment operator (0/1),

M: maximum number of links per user,

 $R_{\rm u}$: reliability of user u,

 ϵ_u : maximum tolerable error on the link,

U: number of users,

B: number of APs.

The objective of the optimization problem is to maximize the overall spectral efficiency for the user. The reliability achieved per user depends on the operating mode.

If the data is duplicated, the multi-connected links are redundant. Hence the reliability of the user is given as (assuming statistical independence)

$$R_u = 1 - \prod_{b=1}^{B} \lceil n_{u,b} \rceil P_{u,b}^e$$

with

 $P^e_{u,b}$: probability of error on the link.

If the data is split across different links, the reliability of the user is given as

$$R_{u} = \prod_{b=1}^{B} \left\lceil n_{u,b} \right\rceil \left(1 - P_{u,b}^{e}\right)$$

In general, the best user-to-cell assignment is carried out by the Hungarian Algorithm. The Hungarian algorithm represents the received SNR for all the users, from all the APs in a matrix, and then performs row operations to get the best assignment for all the users. In the proposed algorithm we calculate M-best assignments for all the users using the greedy algorithm.

		APs			M-best APs
	γ_{11}	γ_{12}		γ_{1B}	$\{1,B\}$
users	γ_{21}	γ_{22}		γ_{2B}	$\{2,4\}$
Ļ	:	÷	·	÷	$\{1\}$
	γ_{U1}	γ_{U2}		γ_{UB}	

The reliability constraints are checked and accordingly the number of links over which the data is duplicated are increased/decreased after every time slot.

6.2 Clustering of mm-wave access points controlled by a 5G low band coverage layer

Extreme mobile broadband in 5G networks will be provided by a millimetre wave radio access technology. To provide a reliable and seamless data transmission to UEs a multi-connectivity scheme consisting of a UE specific cluster of a 5G low band node and several mmAPs is required. Such concept is illustrated in [Figure 6-2A, 5GN-D41], where it is shown how UEs can detect mmAPs and then set up such a UE specific cluster.



Figure 6-4: Deployment of a UE specific multi-connectivity cluster of mmAPs on a 5G-LB coverage layer

The proposed multi-connectivity works on a higher radio layer, i.e. using the Packet Data Convergence Protocol (PDCP) level, which is characterized by an asynchronous control and data layer and a distributed MAC layer scheduling.

Different architectural approaches which take into account backhaul limitations, processing power, the placement of RRC and the split of data flows have been investigated and analysed with respect to the required messages for a given scenario. For all such architectural approaches, it is assumed that the RRC messages towards the UE are transmitted via the 5G low band, in order to guaranty the required reliability inside the defined UE cluster consisting of the 5G LB and mmwave APs (mmAPs). The description and further details of the architectural approaches can be found within Part II of this deliverable.

The 5G NORMA functional control and data layer architecture includes the sub function RRC mmW within the RRC User function block (see Figure 2-3). This sub function provides functionalities which are not available in LTE dual connectivity, such as:

1) Activity Management within MC-Cluster (AMMC)

All 5G mmAPs engaged in the MC framework can be considered as a MC-cluster. Each cluster is UE specific and will be controlled by RRCmmW. Each mmAP in the cluster can be assigned to different activity levels by RRCmmW, see Figure 6-5. We define three main activity levels for each mmAP which are based on the following attributes: i) The UE-context (e.g. air interface related identifier);ii) data in the buffer, and iii) active transmission towards the UE. The activity levels can be set and modified based on the service requirements and the radio conditions of the UE. An early prepared mmAP with context or in addition with buffered data minimizes the interruptions times which may occur due to unpredicted mmW link blockages. Using flags, the mmAP are triggered to transmit user data or just to store the data in a buffer. Consequently, the PDCP protocol has to be enhanced and be able to delete data in a buffer in case it was already successfully transmitted by another mmAP.

2) Link Management within MC-Cluster (LMMC)

Using a new non-conventional physical layer trigger from the mmW-link, a link degradation is detected by transmitting mmAP. This mmAP will inform the RRCmmW functionality which acts pro-actively and triggers the UE to carry out radio measurements towards other mmAP, for a handover to another mmAP. This minimizes the probability of interruptions as this handover will be faster compared to a handover request triggered by the UE.

- a. Using the above functionalities, there are two main benefits which are attained, as explained below. A seamless data transfer based on the proposed clustering scheme can be ensured
- b. The mobility robustness for mmW connectivity is improved, compared to a standalone mmWAP, which is not belonging to a mm-wave cluster.



Figure 6-5: Activity Management within MC-Cluster, triggering of data transmission with flags

6.3 Virtual cells and multi-cell coordination

6.3.1 Novel virtual cell algorithm

The algorithm for formation of the virtual cells has been introduced in [5GN-D41]. This section presents a brief review of this concept and the key enabler in the framework of 5G NORMA. The concept of forming virtual cells has been proposed in [SSP+14] as a technique that allows the utilisation of resources (i.e., allocation of subframes) from multiple base stations.



Figure 6-6: An example of the virtual cell concept (based on [CSS+16])

The introduction of the Software-defined Mobile Network Controller (SDM-C) is one of the key enablers. SDM-C for slice-specific functions and SDM-X for common functions jointly provide the concept of network programmability for 5G networks [5GN-D32]. Each is a unique logically

centralised network function that abstracts and homogenises different network technologies. SDM-C (for RAN slicing Option 1) respectively SDM-X (for RAN slicing Option 2 and 3) enables controlling and coordinating multiple base stations, which is the cornerstone for forming virtual cells. As it is illustrated in Figure 6-6, the cells with different TDD patterns form a new logical cell offering a different TDD pattern compared to the patterns of the forming cells.

In addition, multi-connectivity of a single UE to multiple access points is another key enabler. Multi-connectivity supports simultaneous connectivity and aggregation across different technologies [5GN-D41]. The common MAC approach is the candidate solution for the virtual cell realisation. Each (physical) cell has its individual PHY layer while all of them have the coordinated MAC scheduling the radio resources. Using the common MAC approach will solve the synchronisation and packet numbering for the UEs connecting to the virtual cell while enabling centralised radio resource management techniques. It is noting that the common PDCP approach can be used for this situation instead of a common MAC. The advantages and drawbacks of any of these two approaches are comprehensively discussed in the previous deliverable [5GN-D41].

The formation of the virtual cells can enhance the network performance and the QoS offered to the users. The cells can choose the patterns based on the delay requirement of service, for example, by selecting the patterns with more frequent resources in uplink. The total network throughput can be improved by allocating resources to UE from neighbour cell(s). In other words, the extra flexibility offered by the formation of the virtual cells can be used to provide a low latency platform while balancing the loads between cells to compensate any possible throughput-loss.

The main drawback of the formation of virtual cell, however, is the increased interference in the network and the less efficient usage of radio resources as the result of transmission to/from the UEs in the virtual cell. Hence, the formation of virtual cell can be considered as the trade-off between the increment of throughput for UEs placed in the virtual cell and the decrease of throughput for the other UEs (as the result of increased interference).

The proposed algorithm for formation virtual cells limits the search for the eligible UEs in the edge of each cell. This way the formation can be achieved in shorter time and be practical. Using the threshold for considering the edge of the cells are the key parameter in this approach. The bigger the threshold, the larger is the virtual cell. Also, it is expected to have higher interferences on the reset of the UEs when the virtual cell size is increased and eventually the total network throughput decreases. Hence, finding the optimal (or near optimal) value for the aforementioned threshold is very important in formation of virtual cells, which is non-trivial in dynamic network and moving UEs.

In the first step, as it has been described, is detecting the UEs in the edge of cells. The goal is to serve the UEs in the cell centre in the first step. The UEs near the cell edge are served in the second round of resource scheduling. In order to differentiate the edge UEs from the centre ones, the proposed algorithm considers the difference of Reference Signal Received Power (RSRP) from the Master eNB (MeNB) and Secondary eNB (SeNB). Based on this concept, the edge user are the users, which meet the following condition:

where:

$$|P_{MeNB} - P_{SeNB}| \le \tau_{\rm p}$$

- P_{MeNB} : RSRP from the MeNB,
- *P*_{SeNB}: received power from the SeNB,
- τ_p : cell edge threshold.

After detecting the edge users, in the second step, the UEs are sorted based on their channel qualities and a two –round scheduler serves the UEs with higher SINR with priority. If the cell is not congested and still has available resources, the UEs in the VC region are allocated resources in the second round of scheduling. By allocating a single Physical Resource Block (PRB) to each UE *i* from eNB *k* the throughput is:

$$R_{b[kbps]} = BW_{PRB[kHz]} \min\left(B_{eff} \log_2\left(1 + \frac{\gamma_{i,k}}{\gamma_{eff}}\right), S_{eff}\right)$$

where:

- BW_{PRB} : bandwidth of a PRB,
- *B_{eff}*: bandwidth efficiency (cf. [PDM+11]),
- *S_{eff}*: spectral efficiency (cf. [PDM+11]),
- γ_{eff} : SINR efficiency (cf. [PDM+11]),
- γ_{ik} : SINR of UE *i* connected to eNB *k*.

The interference is the summation of a) received power from other eNBs, b) received power from other UEs, which are working in uplink. It is worth noting that the eNBs operating in the UL direction or UEs in the DL should not be considered in the interference calculations. The received power in uplink mode can be relatively be calculated. The key parts of the proposed algorithm are a) detecting the congested cells, b) detecting the cell edge UEs, and c) alteration of TDD subframe patterns. Table 6-1 summarizes the constants configured during simulation.

Scenarios and numeric results

A practical reference scenario for studying the performance of the proposed algorithm for forming virtual cells has been considered. The key elements in the scenario are the base stations, the traffic profiles and UEs' placement.

For the sake simplicity in analysis of the results, only two base stations with omnidirectional antennas and an inter-site distance of 1 km are assumed. It is worth noting that the interference of the other cells was neglected. Each base station has a bandwidth of 10 MHz with 50 available PRBs. Terrain is assumed to be flat without buildings. Table 6-1 summarises the related details of scenario.

Parameter	Value
Spectral efficiency	4.0
Bandwidth efficiency	0.65
SINR efficiency	0.95
Noise power	-121.45 dB
Antenna gain of eNB	15 dBi
Antenna gain of UE	0 dBi
PRB bandwidth	180 kHz
Number of eNB	2
Inter-site distance	1000 m
Operating frequency	2000 MHz
Bandwidth	10 MHz
eNB DL power per PRB	26 dBm
UE UL power	22 dBm
Number of PRBs over BW	50
Traffic profile	constant 300 kbps

Table 6-1: Virtual cell scenario configurations

In addition, the scenario considers UEs with Constant Bit Rate (CBR) of 300 kbps. UEs are placed in the coverage area of the base stations. The number of UEs varied between 50 (low load) and 250 (high load). Three different distributions of the UEs have been assumed, as follows:

- 1. Scenario A: UEs are uniformly distributed over the coverage area of eNB and 50 % of them are active in the uplink and the reset in downlink. The ratio of UL/DL kept the same in other two scenarios.
- 2. Scenario B: The UEs are not distributed uniformly instead this scenario consider the situation where the distribution of UEs in the outer ring is increased. Thus, 60 % of the

UEs are placed in the outer ring of the cells. The remaining 40 % are placed in the inner rings of the cells (i.e., are cell centre UEs). Cell edge UEs are UEs that fulfil the following property:

$$0.35eNB_{ISD} \le d \le 0.5eNB_{ISD}$$

Cell centre UEs are UEs that fulfil the following properties:

$$0 \le d \le 0.35 eNB_{ISD}$$

where:

- *eNB_{ISD}*: inter-site distance (ISD) between eNBs,
- *d*: distance between UE and (nearest) <u>eNB</u>.
- 3. Scenario C: In the final scenario setup, one of the eNBs is populated with 30 % of the total UEs and the second eNB with the remaining 70 %. Distribution of cell centre and cell edge UEs have been kept as it is as scenario B.

The proposed algorithm for forming the virtual cells has been implemented in NOMOR MxART [NOMOR17] for the aforementioned scenarios. For all three scenarios, the gain of virtual cell formation has been plotted with varying number of UEs and cell edge thresholds. The cell edge threshold has been varied from 3 dB to 27 dB with step size of 3 dB.

Figure 6-7 shows the changes of the total network throughput as the result of implementing virtual cells, exemplarily for Scenario A and different number terminals (50 up to 250). Further numeric results for the other scenario are presented in Part II Section 16). According to the figure, as the cell edge threshold increases, the total network throughput increases by about 30 % compared to the initial value. However, as the threshold increases beyond 23 dB, when there are 150 UEs deployed, the total network throughput decreases again. Although the presented results are the average over 10 simulation runs, the change of the optimal value for the cell edge threshold is due to the placement of the UEs in the virtual cell as well as in other cells. However, the behaviour of system and result pattern for all of the scenarios stays the same. When the threshold value increases, the total network throughput increases. After reaching the maximum total network throughput, a further increase of the threshold leads again to a decrease of the total network throughput.



Figure 6-7: The network throughput increases as a function of the cell edge threshold (Scenario A)

Comparing the optimal value for the cell edge threshold of the three scenarios unveils its high dependency on the placement of UEs. A practical algorithm for forming virtual cells should therefore be able to adopt itself to changes of the user distribution and traffic load. In the next section, the proposed algorithm is extended by adding reinforcement learning techniques.

6.3.2 Forming virtual cell using Q-learning

As it is shown through the numeric results in the former section, the optimal cell edge threshold and the size of the virtual cell is highly dependent on the distribution of users and their traffic demands. In practice, the threshold for forming the virtual cells has to be updated based on the observation and state of the network. In this section, a new approach to form the virtual cells based on Q-learning is proposed.

Q-learning is a model-free reinforcement learning technique [Wat89] to find an optimal actionselection policy for any given (finite) Markov decision process (MDP). It works by learning an action-value function that ultimately gives the expected utility of taking a given action in a given state and following the optimal policy thereafter. A policy is a rule that the agent follows in selecting actions, given the state it is in.

It is assumed that there is learner agent, which take action a_t at the time t even system in state s_t [PW96]. As the result of the action reward r_t is achieved and system moves to the next state, s_{t+1} . The learner agent tries maximise discounted cumulative rewards over multiple iteration. The total discounted return (or simply return) received by the learner starting at time t is given by:

$$r(t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n r_{t+n} + \dots$$

where γ is the discount factor in range of [0,1].

The discount factor γ is used to weight near term reinforcement more heavily than distant future reinforcement. The closer it is to unite the greater the weight of future reinforcements.

The objective is to find a policy π , or rule for selecting actions, so that the expected value of the return is maximized. It is sufficient to restrict attention to policies that select actions based only on the current state (called stationary policies). For any such policy π and for any state *s* we define the value of the policies as:

$$V^{\pi}(s) = V^{\pi}(s) + \alpha (r + \gamma V^{\pi}(s) - V^{\pi}(s))$$

where:

- $\alpha \in [0,1]$: learning rate,
- s: current state,
- \dot{s} : next state after taking the action based on the π policy.

Many dynamic programming-based reinforcement learning methods involve trying to estimate the state values $V^{\pi}(s)$ for a fixed policy π .

Developed from the theory of dynamic programming [Ros83] for delayed reinforcement learning, Q-learning is an incremental algorithm. In this algorithm, the policies and the value function are represented by:

$$Q^{*}(s, a) = R(s, a) + \gamma \sum_{s \in S} T(s, a, s') \max_{a'} Q^{*}(s', a')$$

where:

- *T*: the transition probability function,
- *R*: reinforcement function, also known as reward,
- $\gamma \in [0,1]$: discount factor, the closer γ is to 1 the greater the weight is given to future reinforcements.

At each iteration, we will choose to make either a random action or the optimal action. We will take the first possibility (the random action) with a probability of ϵ . At the beginning of the

learning ϵ must be huge (near 1) in order to visit different states and experiment the effect of different actions. As the agent visits the states and takes actions, it learns. In the next steps, this value is reduced to take actions based on the learned experiences rather than random actions.

In the proposed approach for forming virtual cells using the Q-learning approach the state of the system is defined by the cell edge threshold, which directly indicates the size of virtual cells. The threshold is also changing in steps of 3 dB, similar as in the previous section. Hence, the action set for the operating Q-learning will be:

- Increase: the cell edge threshold is going to be increased by 3 dB up to 30 dB,
- Decrease: the cel edge threshold is going to be decreased by 3 dB up to 0 dB,
- NON: no change applies to the cell edge threshold.

The reward is the differentiation of the average total network throughput, as follows:

$$R(n) = \bar{R}_{b[\text{Mbps}]}(n) - \bar{R}_{b[\text{Mbps}]}(n-1)$$

where:

- R(n): reward in the *n*-th observation window,
- \overline{R}_b : average total network throughput given by:

$$\bar{R}_{b[\text{Mbps}]}(n) = \frac{1}{N_o} \sum_{i=0}^{N_o} \sum_{j=0}^{N_{cell}} \left(R_{b_j[\text{Mbps}]}^{DL}(t_n - i) + R_{b_j[\text{Mbps}]}^{UL}(t_n - i) \right)$$

where:

- N_o : number of frames in the observation interval,
- $R_{b_j}^{DL}$: downlink throughput of cell *j*, $R_{b_j}^{UL}$: uplink throughput of cell *j*,
- t_n : starting frame number of the *n*-th observation interval.

It worth noting that the approximate Q-learning with similar assumptions has been implemented. However, the studies have shown that the system cannot be approximated by the linear functions.

The Q-learning function has been added to the simulator. The learning factor α , is heuristically chosen to be 0.95 and the discount factor, γ , is set to 0.2. The exploration takes place in three phases: i) The epsilon-greedy, ϵ , is set to unit at the beginning of the simulation, ii) the value is then exponentially reduced to 0.1 in the second phase before iii) the value is fixed at 0.1 in order to cope with the network changes. Finally, the Mx-ART simulator [Nom17] is used to simulate a scenario with two eNBs and 500 UEs per eNB (high load). The UEs in cell 1 are active in downlink direction while the UEs in cell 2 are demanding traffic on the uplink. The rest of the parameters for the simulations has kept like in the previous section. Figure 6-8 shows the total network throughput when the Q-learning was active. It can be seen that the Q-learning manages to have significant improvement of the total network throughput and sets the network parameter into almost optimal values.



Figure 6-8: The total network throughput and the three phases of the Q-learning.



6.3.3 Implementation of virtual cells in Demo 1

Figure 6-9: The key elements of Demo 1 of WP6.

The concept of a virtual cell along with the Q-learning algorithm presented in the previous section is going to be implemented in Demo 1 from WP6 [5GN-D62]. The aim of this demo is to present the proof-of-concept for a flexible, adaptive, intelligent, and service-aware (de)composition of NFs and services. This demo is a combination of the software simulations running in parallel to the hardware eNB and SDM-C demo as it is shown in Figure 6-9. The demo is composed by two parts, software (SW) and hardware (HW), developed by 5G NORMA partners Nomor and Azcom, respectively.

The software part of the demo shows that the SDM-C can be coupled with a mobile network simulator. This enables the SDM-C to control software simulated eNBs and reconfigure them into central or edge cloud based on the network parameters and SDM-C's decision logic. The main motivation for a software demo is to highlight the effects of 5G NORMA innovations on a larger scale, i.e. with many eNBs and UEs, and to investigate the gains on network level. The details of the demo are presented in Chapter 16 as well as [5GN-D62].

The Demo 1 is planned to be extended by adding the concept of a virtual cell. The demo is going to have the option added to SDM-C to form a virtual cell either by manually selecting the cell edge threshold or by using Q-learning to set the threshold automatically.

6.4 User-centric connection area

One of the 5G NORMA innovations presented in [5GN-D41] is the so-called User-centric Connection Area (UCA). The UCA concept defines a new RRC state, a new mobility concept inside the UCA and a RAN based paging which should reduce the signalling towards the functional entity inside the 5G core network, which is responsible for mobility management.

The UCA, i.e. the number of radio cells belonging to the RAN based paging area is defined by

- the RRCmmW module, based on UE measurements,
- the SON module, based on cell neighbour relationship and anticipatory knowledge.

In case the location of the UE inside the UCA is no longer known, the RAN paging module is triggered by the RRCmmW module to initiate a paging within all radio cells belonging to the

UCA. For this purpose, the interfaces "s" and "p" have been defined, see Figure 2-3. The information exchanged via these interfaces is described in [5GN-D41].

We note that in Part II Chapter 17 the proposed reduction of messages is confirmed based on results of a detailed simulation study.

In this section, the LTE procedure with the transitions between RRC_CONNECTED and RRC_IDLE is compared with the UCA concept in which the UE remains within the UCA_Enabled state. During the course of the 5G NORMA project the main proposals of the UCA concept have been adapted by 3GPP for the New Radio Access Technology [38.804].

In [38.804] the UCA is defined as RAN-based notification area. In particular, a notification area can cover a single or multiple cells, and can be smaller than a CN area, e.g. a tracking area. In addition, a UE does not send any "location update" indication when it stays within the boundaries of the notification area. Moreover, leaving the area, a UE updates its location to the network. Furthermore, a new RRC_INACTIVE state, comparable with the 5G NORMA state UCA_Enabled, will be introduced, characterized by five major elements, namely:

- cell re-selection mobility, i.e. the CN is not informed about the mobility inside the RANbased notification area;
- the UE AS context is stored in at least one gNB and in the UE;
- paging is initiated by NG-RAN;
- RAN-based notification area is managed by NG-RAN; and
- NG-RAN knows the RAN-based notification area which the UE belongs to.



Figure 6-10: Signalling messages towards the core based on 4G LTE
With the introduction of the new NORMA RRC state "UCA_Enabled", the signalling load towards the core network is drastically reduced compared to state of the art in LTE, especially for short data packets (SDP) in uplink and downlink. The transmission of SDP in downlink generates several cycles from RRC_CONNECTED to RRC_IDLE and back to RRC_CONNECTED transitions, these transitions are depicted in Figure 6-10. Signalling messages for the S1 release procedure, the paging procedure and bearer establishment have to be exchanged with the EPC.

As long as the UE does not leave the coverage of the UCA no signalling messages towards the EPC are required, this is depicted in Figure 6-11. In case the UE leave the UCA, a UCA update is required, which generates signalling messages to the EPC.



Figure 6-11: Reduction of signalling messages towards the core with the NORMA UCA concept

The gain with respect to LTE in terms of the amount of signalling messages towards the core network was evaluated based on a detailed system simulation and is depicted in Figure 6-12.

The gains are depending on

- RRC timer: for larger timer values the UE might remain in the RRC_CONNECTED state as new uplink or downlink packets might arrive while the RRC timer is running. This prevents the transition to the RRC_IDLE state, i.e. reduces the amount of signalling towards the MME/EPC,
- UE speed: in LTE a larger amount of UE handovers will occur for a higher UE speed. After the handover the UE remains in the RRC_CONNECTED state. This has to be compared against the increased number of UCA updates for higher UE speed. It is

assumed that a UCA update requires the same amount of signalling than a UE handover in LTE.

The gain of the UCA framework as compared to LTE is calculated as follows:

$$\Omega = 100 * (1 - \frac{\alpha}{\beta})$$

where, Ω is the percentage gain over LTE, α is the MME/EPC signalling due to UCA updates and β is the total MME/EPC signalling in LTE due to the events: UE originated Idle-Active transition, S-GW originated Idle-Active transition, eNB terminated Active-Idle (see Figure 6-10) and intra MME X2 handover.

The UCA size is defined by the anchor node for the UE by applying a window size in dB: a neighbour radio cell which Reference Signal Received Power (RSRP) is not lower than "x" dB compared to the RSRP of the UE in the anchor node will be included within the UCA of this UE.

Interestingly, we note that, depending on the UE speed and the RRC timer values, gains of up to almost 100 % can be reached.



Figure 6-12: Gain over LTE for different RRC timer values and different UE speeds

In addition, the UCA concept enables considerable gains on the air interface by reduction of paging messages. If a UE is paged on LTE basis within the complete paging area, e.g. all 90 radio cells within the simulation environment are included, the gain can reach up to 100 % as depicted in Figure 6-13.

3GPP has defined a LTE paging optimization, i.e. only a few cells near to the last known UE location are paged in a first round. If this paging is not successful, a second round might span over a larger area, e.g. the complete paging area. Even with this optimization high gains can be achieved which are shown in the left part of Figure 6-13. For higher UE speeds the gain decreases, as also more UCA updates will occur, which also generate load on the air interface.



Figure 6-13: Gain over LTE for different paging are sizes

Further simulation results, the simulation environment and simulation parameters can be found within Part II of this deliverable.

6.5 Mobile edge computing

Mobile edge computing concept provides the technological means for tenants and external providers to install own application for customized services. The mobile edge computing (MEC) platform can be envisaged as an IT platform where tenants' applications are installed and managed by the infrastructure (or platform) provider. The MEC architecture is depicted in Figure 6-14. This is a draft picture derived from [MEC003] to show the concept of our MEC applications. In particular, we highlight the Slice QoS Monitoring functional block as an MEC application installed at the Mobile edge host level. Such application can be installed on platforms residing on (or close to) eNBs to provide high accuracy during the monitoring phase. This application exploits the MEC services, such as Radio Network Information Service (RNIS) and could be dynamically configured on different/multiple MEC platforms based on traffic conditions. However, it can cooperate with the QoS Control App, as shown in Figure 2-2. Each network slice may issue its own monitoring application based on different KPIs. Please note that the MEC functional block depicted in Figure 2-2, Figure 2-3 and Figure 2-5 refers to the Mobile edge host block and mobile edge platform manager (depicted in Figure 6-14), where different MEC applications can be deployed. Communication between MEC platforms and 5G NORMA MANO layer is performed by means of Mm2, though it is still under definition.

The Slice QoS Monitoring block is logically attached with the 5G NORMA Network Slice Broker (described in Section 5.1.2), as it continuously provides an accurate feedback on the QoS guarantees. The 5G NORMA Network Slice Broker may dynamically change slice configurations (e.g., it may change the forecasting parameters to force the system to be more conservative as to avoid slice SLA violations) by means of ad-hoc Slice QoS Monitoring applications, such as the Mobile Throughput Guidance Service. A message sequence chart for this particular MEC service is provided in Figure 6-15. Specifically, radio conditions are monitored through a Radio Network Information Service (RNIS) specified per user (or per network slice). Different actions might be triggered: the SDM-X controller may directly adjust network resources assigned to different network slices in order to efficiently handle slice SLA violations.



Figure 6-14: Mobile Edge Computing platform supporting Slice QoS Monitoring



Figure 6-15: Message sequence chart of MEC slice QoS monitoring/enforcement

In the following Table 6-2 we provide an overall description of the process. The MEC block will monitor the radio channel condition by means of a Radio Network Information Service (RNIS). Monitoring values are retrieved per network slice (or per UE). Additionally, the MEC might monitor the slice SLAs requirements, such as latency or throughput. As soon as slice SLAs (or user SLAs) are not fulfilled, the MEC function block may trigger an SDM-X application, e.g., QoS Control App,

which will dynamically adjust the network resources assigned to that particular network slice through the 5GNORMA-SDMC-SDMX interface.

Actors:	Triggering actor:	
	• RRC User	
	Involved actors:	
	SDM-X/Multi Tenancy Scheduling	
	• SDM-C	
	MAC Scheduling	
	• MEC	
	QoS Control App	
Preconditions:	UE is in RRC_CONNECTED	
	• The Radio Network Information Service (RNIS) is available at the	
	MEC block.	
	• MEC App is running to monitor end-to-end network performance	
	(throughput or latency)	
Postconditions:	UE is in RRC CONNECTED	
	• Network slice resources readjusted to satisfy overall end-to-end	
	slice requirements	
Frequency of	Triggered per signalling radio bearer (SRB) per UE; on average expected	
Use:	to be much less.	

Table 6-2: Process of MEC triggering	MAC Scheduling to adapt scheduling
--------------------------------------	------------------------------------

Normal Course	1. An SRB is established through the control layer between RRC User
of Events:	blocks.
	2. UE is in RRC_CONNECTED. Therefore, the RRC User block
	performs measurement reporting.
	3. The Radio Network Information Service (RNIS) is enabled and available at the MEC.
	4. <i>MEC</i> will enable the Mobile Throughput Guidance signalling protocol to assist TCP while ensuring high utilization and high service delivery performance. Alternatively, the <i>MEC</i> could enable a Slice Latency Monitoring service for monitoring and trigger slice resource adjustments on-the-fly, e.g., by means of the QoS Control App.
	5. If the Service requirements are not satisfied, the QoS Control App by means of the MEC block will communicate through the <i>5GNORMA-SDMC-SDMX</i> interface with the <i>SDM-X/Multi-tenancy Scheduling</i> for dynamically adjusting the network resources dedicated to a particular network slice.
	6. The <i>Multi-tenancy Scheduling</i> block will instruct the <i>MAC Scheduling</i> block to optimally arrange the radio resources based on the management policies.
Alternative Courses:	• If no network slice is considered or only a single slice is considered, the QoS Control App will directly communicate with the <i>MAC Scheduling</i> block.
Exceptions:	
Assumptions:	• A system could have multiple MEC blocks providing the same services.
Notes and Issues:	

6.6 Massive machine-type communication RAN congestion control



Figure 6-16: Different congestions and their sources in LTE-A networks

In current Long Term Evolution (LTE) and LTE-Advanced (LTE-A) systems, network congestions mainly happen either in the radio access network (RAN), at the Mobility Management Entity (MME) or at the core network gateways (GWs), as shown in Figure 6-16. RAN congestion is sensitive to the device density. MME congestion is sensitive to the device mobility. GW congestion is sensitive to the overall user data rate. Differing from human-type communications (HTC), machine-type communication (MTC), especially massive MTC

(mMTC) in the upcoming 5G networks, is much more sensitive to RAN congestions rather than the other two kinds. The reason is that many MTC devices (MTCDs) are immobile or with low mobility, and send only small data packets in each transmission. Additionally, they are usually synchronized to each other, e.g. all the sensors in a same sensor network can be usually set to upload measurements at the beginning of every minute/hour. Hence, mMTC applications can easily create bursts of random access (RA) collisions without triggering congestion at the MME or GWs.



Figure 6-17: Comparing the performances of different RAN congestion controlling methods

Many different solutions have been designed to control RAN congestions for mMTC, which can be briefly classified into some categories, as shown in Figure 6-17. A state-of-the-art survey on existing mMTC RAN congestion control approaches is available in [AHK17]. The D2D-based grouped RA approach here is developed based on the concepts proposed in [CKS+15]. Compared to the other methods, D2D-based grouped RA shows ubiquitous advantage in energy efficiency and access delay. The principle, as shown in Figure 6-18, is to cluster MTCDs into groups, each group consists one and only one group coordinator (GC), which is connected to other group members (GMs) over D2D links. GC accesses the base station for the entire group, aggregates UL data from GMs, and distributes DL data to them. Devices are also labelled with different device classes (DCs), each group consists of only devices of same DC. RACH resources can be either dedicated to different DCs, or shared by them.



Figure 6-18: The RAN topology in D2D-based grouped RA

It is clear, that increasing the group size leads to a decrease in RA request density, and hence reduces the RAN congestion rate. However, it should be considered that the intra-group D2D links are not always available and reliable. For instance, some links may be assumed by the base station as available and utilized for the intra-group communication, while the channel is actually strongly attenuated by shadowing objects or interferers nearby. Even an established D2D link can also fade due to the mobility of device, the time variation of channel condition or out-used device batteries. The reliability even drops when the group size grows. And the link failures will lead to extra signalling for reports and exception handling, this will compromise the gain if these messages are transmitted through extra RA processes. Assuming an exponentially decreasing average D2D link lifetime w.r.t. the group size, unreliable D2D links can even lead to a loss in case of large groups. To take the unreliable D2D links into account, we have carried out the following tasks:

- designing a grouped RA protocol with D2D link failure report;
- embedding D2D controlling messages as much as possible into d-layer messages;
- selecting the optimal group size and the grouping map; and
- selecting the optimal RACH resource allocation.

First, three different processes are designed in the protocol, including global group update, group joining, and group leaving. For details, including a new frame structure, cf. Part II Chapter 19.

Subsequently, we investigated the RACH resource allocation in D2D-based grouped RA. Generally, there are three strategies for sharing random access opportunities (RAO) among device classes:

- full sharing, where every RAO is available for all DCs,
- full dedication, where every RAO is dedicated to one specific DC, and
- partial dedication, where every RAO can be either dedicated to one DC, or shared by multiple DCs.

The full sharing strategy returns average collision rate for all DCs, full dedication is able to isolate DCs from each other, so that RA request burst from one DC will not affect other DCs. It shall be optimized to maximize the RAO utilization. Partial dedication is too complex to optimize.

Collision density:

Following [37.868], in grouped RA, approximately assuming large amount of RA requesters, the collision rate for each request in DC *i* under full dedication can be expressed as

$$p_{\text{FD},i} = 1 - e^{-\frac{\gamma_i}{L_i}},$$

Where γ_i is the RA density (per second) in DC *i*, and L_i is the number of available RAO dedicated to DC *i* per second. Further approximately assuming that all collision events are independent from each other, the overall collision density is

$$C_{\rm FD, cell} = \sum_{i=1}^{N} \gamma_i p_{{\rm FD},i},$$

and the overall probability of collision occurrence

$$p_{\text{FD,cell}} = 1 - \prod_{i=1}^{N} [\exp(-\frac{\gamma_i}{L_i})]^{\gamma_i} = 1 - \prod_{i=1}^{N} e^{-\frac{2\gamma_i}{L_i}}$$

Aiming at minimizing the overall collision occurrence probability, there is the optimal allocation

$$\frac{\gamma_i}{L_i} = \frac{\gamma_j}{L_j} \quad [i, j] \in \{1, 2, \dots, N\}^2,$$

where *N* is the number of DCs.

Simulations show that the approximate model of independent collisions matches the numerical results in the collision density, the optimal allocation can be estimated by this model with satisfying accuracy as Table 6-3, Figure 6-19 and Table 6-4 show.

Class	Accessing devices	Avg. access frequency	RA density
1	3000	1/60 Hz	50 Hz
2	30 000	1/300 Hz	100 Hz
3	30 000	1/60 Hz	500 Hz
4	30 000	1/30 Hz	1000 Hz

 Table 6-3: Simulation specification



Figure 6-19: Collision densities with respect to RACH resource dedication in a 2-DC-case (DC1+DC2), numerical results obtained from 500 iterations of Monte-Carlo test

DC Combination	1&2	1&3	1&4
Optimal. L ₁ (estimated)	3600	982	514
Optimal. L ₁ (simulated)	3460	1048	534
Collision density at est. opt. L ₁ (Hz)	1.766	26.320	95.690
Collision density at sim. opt. L ₁ (Hz)	1.998	26.780	97.042
Collision density under full sharing (Hz)	2.166	26.940	97.634

Then we also designed a reserve-and-allocate method to implement DC preference in the RACH resource allocation, in order to provide special DCs with guaranteed QoS, as shown in Figure 6-20.



Figure 6-20: Performances of different allocation methods, when DC1 is required to have an average collision rate of 0.02

6.7 Geolocation database

Geolocation and associated management capabilities, administered through a GDB or hierarchy of GDBs operating in government/regulatory and operator contexts, will be highly relevant in 5G communication contexts. This is for reasons such as:

Spectrum Sharing

5G technologies will require spectrum sharing to achieve sufficient spectrum availability in some scenarios, e.g., at lower frequencies for coverage/reliability and signalling/control purposes to realize URLLC. Such spectrum sharing will often be between very different services and owners of the spectrum, such that the spectrum sharing must usually be approved by the regulator—typically through automated regulatory-driven or certified GDBs.

A number of regulatory initiatives/trials are being undertaken or have been completed showing such GDB concepts in action, in the contexts of TV white space (TVWS) [Hol16], Licensed-Shared Access [MOP+14], and the Citizens Broadband Radio Service [CBRS15]. Many light-licensing regimes, such as Program Making and Special Events (PMSE) licensing in the UK, are also effectively administered through regulatory GDBs.

Heterogeneity

One key aspect of 5G is heterogeneity; moreover, common management among the heterogeneous systems/elements in 5G will likely not preexist. A GDB can assist in managing connectivity, QoS, and other aspects in heterogeneous scenarios, also integrated with spectrum sharing/management at higher levels, through combination with regulatory-run or approved GDBs. In essence, this means that the GDB could operate in a distributed way partly within (and run by) the mobile network and partly within (and run or certified by) the regulator; alternatively, a purely mobile network/operator-run GDB could ask for regulatory certification. There is already precedent for similar such things happening in TVWS and CBRS contexts, for example. This is inherently possible because the GDB or GDBs will (at least partially) operate at higher (e.g., regulatory—transcending spectrum services and owners) levels for abovementioned reasons.

Such databases (e.g., TVWS and CBRS databases, etc.) under regulatory approval often incorporate aspects that can be reused for other such management purposes. They include advanced propagation and context (e.g., transmitter and receiver locations and characteristics) knowledge, facilitating the management of QoS and associated allocation, connectivity, and other aspects in heterogeneous networking scenarios. The GDB will therefore have all information to assist connectivity opportunities that would otherwise not be spotted without all radios continually

being switched on a scanning in a heterogeneous networking context, which would be very resource.

Infrastructure/device and other equipment location awareness

In many 5G and future communications contexts it will be necessary to consider the precise locations of equipment in the management of resources. This applies not only to spectrum resources, but also to computational resources achieving network functions through virtualization. One example here is the context of latency reduction for URLLC, requiring careful geographical/location placing/instantiation of virtualized equipment by network management to minimize propagation delay by achieving a direct as possible propagation path between the communication endpoints.

It is particularly noted that as well as the purposes given above, a key topic of interest is the potential to assist mMTC as covered in Section 6.6 above, hence the interface with mMTC Congestion Control in Figure 2-2, Figure 2-3, and Figure 2-5. Moreover, in the context of 5G NORMA, the GDB also can particularly serve optimisation based on localised (i.e., geolocation-based) traffic demand vs. radio resource availability. In that context, the GDB also interacts with RRC, both in this context, and in terms of the use of radio resources that it unlocks through spectrum sharing as discussed above.

7 Conclusions

With this final deliverable of 5G NORMA WP4, we presented the outcome of the control and data layer architecture design of the third and last design iteration along with a concise summary of individual innovations in that area. We showed that the architecture design offers great architectural flexibility by supporting concurrently all of multi-tenancy, multi-service, multi-RAT (by RAN slicing Option 1) and (intra- and inter-tenant) multi-connectivity. In detail, the major WP4 outcomes with respect to the architecture design that we presented are:

- the functionally decomposed control and data layer architecture, yielding a network of functions as shown throughout chapters 2, 3 and 4;
- classification and characterization of these network functions into one of i) centralized control, for control functions employing the SDMC concept, i.e. implemented as SDMC applications (Section 4.1), ii) decentralized control, for time critical and signalling intense real time control (Section 4.2), and ii) the inherently distributed data layer functions;
- characterization of data layer functions in the physical layer, (radio) link layer ("Layer 2") and the non-access stratum (Section 4.3), of which only the physical layer NFs are partly realized as PNFs, either because the function inherently cannot be virtualized (PHY TP) or implementation in a non-virtualized manner provides energy and cost efficiency benefits (PHY Cell, PHY User);
- how to adapt through suitable function selection and placement to service and deployment requirements for three major services broadband, low latency and mission critical (Section 3.1);
- how to realize a flexible and future proof slice-aware and multi-service-capable shared RAN instance, i.e. physical cell (Section 3.2);
- three RAN slicing options 1) slice-specific RAN (Section 2.1.1), 2) slice-specific radio bearer (Section 2.1.2) and 3) slice-aware shared RAN (Section 2.1.3, for which Section 17.1 adds a suitable slice multiplexing), and representative use cases for each;
- how different RAN slicing options integrate within a single RAN instance to best support the service mix of tenants that share the RAN (Section 2.1.4); and
- how different options of multi-connectivity, namely common PDCP and common MAC, combine with any RAN slicing option for full service flexibility (using Option 2 as example, Section 2.1.2).

For several specific innovations presented in chapters 5, 6 and Part II, we provided numerical results. In detail, we showed the

- benefits of coordination, enabled by relocation of functions (scheduling coordination) to the edge cloud (WP6 Demo 1, Section 6.3.3 and Section 16.3);
- benefit of multi-connectivity when used to realize virtual cells, which improves overall cell throughput (Section 6.3.1 and Section 16.2);
- benefit of multi-connectivity in conjunction with PDCP-based packet distribution, where a centralized placement of PDCP offers increased efficiency and thereby throughput with respect to the decentralized X2 interface (Chapter 14);
- benefit of the machine learning technique Q-learning for i) learning the optimal cell edge threshold that determines assignment of UEs to virtual cells (Section 6.3.2 and Chapter 16) and for ii) learning the best (close to optimal) performing policy for multi-tenancy admission control (Section 5.2.1 and Section 11.1);
- reduction in signalling messages both over the air and towards the core through introduction of a new third RRC state and RAN-controlled mobility (UCA) (Section 6.4 and Chapter 18);
- reduced collision probability in mMTC signalling through GDB assisted clustering and accordingly the improved scalability to a larger number of MTC devices (Section 6.6 and Chapter 19).

Finally, we not only showed by simulation (Section 5.2.2 and Section 11.2) but proofed mathematically the benefit of dynamic sharing over static sharing (applicable to RAN slicing options 2 and 3) and thereby the usefulness of the multi-tenancy approach at large [CBV+16].

8 Annex A: Why no virtualized connectivity on the air interface?

As already stated, the non-access stratum implementation of the c/d-layer utilizes network slicing to realize both, multi-tenancy and multi-service support. Both terms may be used interchangeably here. As explained, the primary enabler is the fully virtualized infrastructure that abstracts the underlying hardware into generic units of networking/communication, computing/processing and storage/memory that are sufficiently equal and independent of each other. Why doesn't the access stratum follow the same approach?

The non-access stratum primarily employs optical transport for communication. Latency is dominated by the signal propagation time (5 us/km for optical fibers). Packet time and unavoidable queueing delay (switching delay in non-overload state) in packet switched networks become negligible with typical line rates of 10 Gbit/s (e.g. 1.2 us per max-MTU-sized Ethernet frame) and beyond. Similarly, packet processing time becomes negligible with multi/many-core CPU/ASICs at multi-GHz clock rate with multi-channel memory of multi-Tbit/s channel bandwidth. Accordingly, any deviation from the idealized assumption of equal and independent computing, storage and networking resource units becomes sufficiently negligible that there is no benefit and therefore neither a need nor any interest to influence or adapt the underlying computing, storage and networking implementation for different tenants or services.

With respect to the access stratum, i.e. the air interface, this idealized assumption does not hold. Spectrum is a shared resource, which creates interdependencies (interference) among concurrent channel uses that, although limited to a spatially correlated set of users, is only imperfectly known a priori and accordingly hardly to predict and compensate for. Spectrum usability differs due to physics and regulations. Spectrum is a scarce resource, partly due to regulation and former (technology-limited) less efficient use, partly due to technical feasibility (especially macro base stations are typically bandwidth and transmit power limited).

Foremost, though, the air interface has the highest cost per bit of all links in the end-to-end path and consequently represents the bottleneck link with just as much capacity as needed. Given certain minimum requirements in one or multiple of service requirements 1a) rate, 1b) latency and 1c) reliability, coverage requirements 2a) range/penetration, 2b) user speed and 2c) supported scattering environments and deployment requirements 3a) site, 3b) connectivity and 3c) power availability (at a certain cost), a compromise is made in the others to keep the cost per bit as low as possible. Accordingly, there is no one type of realization that supports all services, as that would be cost inefficient, and it becomes beneficial to individually adapt the underlying implementation of the packet transport to each specific service (and tenant).

9 Annex B: List of function blocks

Network Functions in the 5G NOMRA c/d-layer have been named such that their respective purpose can be inferred from their name. To further assist correct understanding, the following subsection shortly summarize the main tasks of each function block and its subblocks, if any. An in-depth characterization of all listed functions blocks can be found in [5GN-D41] Annex A.

9.1 Data layer

- **PHY Transmission Point.** The analogue and mixed signal processing for all signals transmitted (received) via one transmission (reception) point.
- **PHY Cell.** (De-)multiplexing of *PHY User* and baseband signal generation including common PHY signals for one RAT or slice.
- **PHY User.** The generation of the baseband signal (in frequency domain for OFDM-based systems) from user data (DL) and decoding of baseband signals into user data (UL), respectively.
- MAC. (De-)multiplexing of radio bearers (MAC SDUs, including concatenation in case of 5G NR) to/from RLC and control information to/from MAC Scheduling (MAC control elements) out of/into PHY transport blocks.
- **MAC Carrier Aggregation.** Coordinates the exchange of scheduling information as well as feedback information corresponding to the aggregated legs.
- **RLC.** Unacknowledged Mode (UM) provides segmentation and reassembly as well as (re)concatenation (for 4G only) of RLC SDUs to adapt the amount of user data according to the PHY TB (resp. MAC PDU) size. Acknowledged mode (AM) additionally provides a re-transmission process (ARQ). Transparent mode (TM) provides none of the above functions, only forwarding RRC messages of BCCH, PCCH, MCCH (RRC Cell) and of CCCH (RRC User) logical channels unmodified between RRC Cell/RRC User and MAC.
- **PDCP Split Bearer.** Executes the functionalities of routing, reordering, and reordering timer associated with PDCP level multi-connectivity.
- **PDCP.** Carries out data transfer functions including functions sequence number maintenance, RoHC, (de-)ciphering, integrity protection, and verification.
- **eMBMS.** Performs the transmission of MBMS application data using the IP multicast address with the addition of SYNC protocol to guarantee that radio interface transmissions stay synchronised. Multimedia Broadcast Multicast Services (MBMS) offer support for broadcast and multicast services enabling the transmission of multimedia content (text, pictures, audio and video) and utilizing the available bandwidth intelligently [23.246].
- NAS. Performs routing functionality (S-GW, P-GW) and service delivery.

9.2 Distributed control

• **RRC Cell.** Handles control layer signalling protocols associated with broadcasting system information, including NAS common information and information relevant to UEs in RRC_IDLE, e.g., cell (re-)selection parameters, neighbouring cell information, and information (also) applicable for UEs in RRC_CONNECTED, e.g., common channel configuration information.

- **RRC User.** Handles the UE management and control including radio bearer setup. It includes the subblocks
 - **RRC mmW.** It adds functionality related to mm-wave transmission points controlled by 5G coverage cell and UCA for short data packet transmission.
 - **RAT/Link Selection.** Enables link selection and packet scheduling if the UE is simultaneously connected to two or more RATs.
- MAC Scheduling. Scheduling the transfer of user data and control signalling in DL and UL subframes over the air interface, including scheduling HARQ (re)transmissions. It comprises the radio scheduler, which provides time-frequency resource allocation, MCS selection and computation of precoding adapted to the estimated instantaneous (fast fading) radio channel. Further tasks are: transmission mode and CSI reporting selection and configuration, (UL) power control, to adapt to the (environment dependent) channel statistics. It does UL power control and configures discontinuous reception (DRX) to improve UE battery life.
 - **QoS Scheduling.** It adds functionality for prioritization among users to fulfil their QoS requirements and preselect users for the radio scheduler.

9.3 Centralized control

- **Multi-tenancy Scheduling.** Coordinating resource sharing among multiple tenants. It includes functions such as scheduling and ICIC schemes.
- **mMTC RAN Congestion Control.** Grouping mMTC devices into context-based clusters and schedule their RAN procedures in sub-frames, to reduce the RAN congestion rate.
- **QoS Control.** Network monitoring and configuration in real time through open interfaces that interact with the SDM-X and the control radio stack.
- **SON.** It covers *i*) Self-configuration, *ii*) Self-Optimisation, *iii*) Self-Healing, and *iv*) User Centric Connection Area (UCA) and mm-wave cluster configuration.
- **RAN Paging.** To reduce signalling messages on the air interface and towards the CN, 5G NORMA employs a RAN paging approach, i.e. inside a User Centric Connection Area (UCA), in addition to paging in a larger tracking area. Further, it sets up and updates connectionless transmission inside a UCA.
- **RRC Slice.** Handles the UE management and control related to the slice specific part of the RAN protocol stack.
- **eMBMS Control.** Performs the admission control and allocation of the radio resources, UE counting procedure, MBMS session management (initiating the MBMS session start and stop procedures), allocation of an identity and the specification of QoS parameters associated with each MBMS session.
- NAS Control. It includes the subblocks
 - NAS UE specific. Refers to the user specific NFs and procedures related to the data layer, which are triggered by the NAS UE-specific control layer functions. For example, a mobility management (MM) application that controls the behaviour of gateways following a software defined principle
 - NAS UE Specific and Data Layer. Refers to the user specific functions and procedures related to the d-layer, which are triggered by the NAS UE-specific c-layer functions.

- **NAS UE Specific and Control Layer.** Refers to the user specific functions and procedures related to the non-radio signalling between the UE and MME.
- **NAS Event-Control Layer.** Refers to the network-side c-layer functions and procedures including those of the interface between RAN and CN (S1-C in 3GPP LTE [36.300]), which are provided to facilitate mobility management.
- **GDB.** Geolocation Database (GDB) stores information linked to geolocation, and makes decisions based on that geolocation information.

10 References

- [22.852] 3GPP TR 22.852, "Study on Radio Access Network (RAN) Sharing enhancements (Release 13)", September 2014.
- [23.203] 3GPP TS 23.203, "Policy and charging control architecture (Release 14)" June 2017.
- [23.501] 3GPP TS 23.501, "System Architecture for the 5G System; Stage 2 (Release 15)", June 2017.
- [23.711] 3GPP TR 23.711, "Enhancements of Dedicated Core Networks selection mechanism (Release 14)", September 2016.
- [23.799] 3GPP TR 23.799, "Study on Architecture for Next Generation System (Release 14)", December 2016.
- [36.300] 3GPP, TS 36.300, "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2 (Release 14)", June 2017.
- [36.808] 3GPP TR 36.808, Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Carrier Aggregation; Base Station (BS) radio transmission and reception (Release 10), July, 2013.
- [36.933] 3GPP TR 36.933, "Study on Context Aware Service Delivery in RAN for LTE; (Release 14)", March 2017.
- [37.868] 3GPP TR 37.868, "Study on RAN Improvements for Machine-Type Communications (Release 11)", September 2011.
- [38.300] 3GPP TS 38.300, "NR; NR and NG-RAN Overall Description; Stage 2 (Release 15)", June 2017.
- [38.801] 3GPP TR 38.801, "Study on new radio access technology: Radio access architecture and interfaces (Release 14)", April 2017.
- [38.804] 3GPP TR 38.804, Study on New Radio Access Technology; Radio Interface Protocol Aspects (Release 14), March 2017.
- [5GC-D21] EU H2020 5G-CROSSHAUL, "Detailed analysis of the technologies to be integrated in the XFE based on previous internal reports from WP2/3", Deliverable 2.1, June 2016.
- [5GN-D23] EU H2020 5G NORMA, "Evaluation architecture design and socio-economic analysis – final report", Deliverable D2.3, December 2017. Available at: https://5gnorma.5g-ppp.eu/wpcontent/uploads/2017/12/5g_norma_d2-3.pdf
- [5GN-D31] EU H2020 5G NORMA, "Functional Network Architecture and Security Requirements", Deliverable D3.1, December 2015. Available at: https://5gnorma.5g-ppp.eu/wp-content/uploads/2016/11/5g_norma_d3-1.pdf
- [5GN-D32] EU H2020 5G NORMA, "5G NORMA network architecture intermediate report", Deliverable D3.2, January 2017. Available at: https://5gnorma.5g-ppp.eu/wp-content/uploads/2017/03/5g_norma_d3-2.pdf
- [5GN-D33] EU H2020 5G NORMA, "5G NORMA network architecture final report", Deliverable D3.3, September 2017. Available at: https://5gnorma.5gppp.eu/wpcontent/uploads/2017/10/5g_norma_d3-3.pdf
- [5GN-D41] EU H2020 5G NORMA, "RAN architecture components intermediate report", Deliverable D4.1, November 2016. Available at: https://5gnorma.5gppp.eu/wpcontent/uploads/2016/12/5g_norma_d4-1.pdf
- [5GN-D52] EU H2020 5G NORMA, "Definition and specification of connectivity and QoE/QoS management mechanisms – final report", Deliverable D5.2, June 2016. Available at: https://5gnorma.5gppp.eu/wpcontent/uploads/2017/07/5g_norma_d5-2.pdf
- [5GN-D62] EU H2020 5G NORMA, "Demonstrator design, implementation and final results", Deliverable D6.2, November 2017. Available at: https://5gnorma.5g-ppp.eu/wp-content/uploads/2017/12/5g_norma_d6-2.pdf

[5GN-D72]	EU H2020 5G NORMA, "Communication and dissemination – final report", Deliverable D7.2, December 2017. Available at: https://5gnorma.5g-
[5GPPPWP]	ppp.eu/wpcontent/uploads/2017/12/5g_norma_d7-2.pdf 5G_PPP_"5G_empowering_vertical_industries" White_Paper_April 2016
	Available at: https://5g-ppp.eu/wp-
	content/uploads/2016/02/BROCHURE_5PPP_BAT2_PL.pdf
[AHK17]	M. S. Ali, E. Hossain, and D. I. Kim, "LTE/LTE-A Random Access for Massive
	Machine-Type Communications in Smart Cities," IEEE Communications Magazine, vol. 55, no. 1, pp. 76–83, 2017.
[CBRS15]	"In the Matter of Amendment of the Commission's Rules with Regard to Commercial Operations in the 3550-3650 MHz Band, Report and Order and Further Notice of Second Proposed Rulemaking," Technical Document, Federal
	Communications Commission (FCC), April 2015.
[CBV+16]	P. Caballero Garces, A. Banchs, G. de Veciana and X. Costa-Perez, "Multi-
	Tenant Radio Access Network Slicing: Statistical Multiplexing of Spatial Loads", CoRR, July 2016. Available at: http://arxiv.org/abs/1607.08271
[CKS+15]	K. Chatzikokolakis, A. Kaloxylos, P. Spapis, N. Alonistioti, C. Zhou, J. Eichinger, Ö. Bulakci, "On the Way to Massive Access in 5G: Challenges and Solutions for Massive Machine Communications", in International Conference
	on Cognitive Radio Oriented wireless Networks, pp. 708–717, Springer
[000 16]	International Publishing, April 2015.
[CSS+10]	S. Costanzo, R. Snrivastava, K. Sarndanis, D. Xenakis, X. Costa-Perez, and D.
	Grace, Service-oriented resource virtualization for evolving IDD networks
	towards 5G, in 2016 IEEE Wireless Communications and Networking
	Conference, Doha, pp. 1–6, April 2016.
[F5G-D31]	EU H2020 FANTASTIC-5G, "Preliminary results for multi-service support in link solution adaptation", Deliverable 3.1, May 2016. Available at: http://fontastic5g.ou/wp.content/uploads/2016/06/EANTASTIC
	5C D21 public pdf
[E5G D32]	EU H2020 EANTASTIC 5C "Final report on the holistic link solution"
[F30-D32]	Deliverable 3.2, April 2017. Available at: http://fantastic5g.eu/wp-
[E5G D41]	EU H2020 EANTASTIC 5G. "Technical Results for Service Specific Multi
[150-041]	Node/Multi-Antenna Solutions", Deliverable 4.1, June 2016. Available at:
	http://fantastic5g.eu/wp-content/uploads/2016/06/FANTASTIC-
	5G_D41_public.pdf
[F5G-D42]	EU H2020 FANTASTIC-5G, "Final results for the flexible 5G air interface multi- node/multi-antenna solution", Deliverable 4.2, April 2017. Available at: http://fantastic5g.eu/wp-content/uploads/2017/05/FANTASTIC-
	5G D42 final ndf
[Hol16]	O Holland "Some Are Born With White Space Some Achieve White Space
	and Some Have White Space Thrust Upon Them" IEEE Transactions on
	Cognitive Communications and Networking vol 2 no 2 nn 178–193 2016
[MEC003]	TSLGS MEC 003: "Mobile Edge Computing (MEC): Eramework and Reference
[MLC005]	Architecture" v1 1 1 March 2016
[MOP+14]	M Matinmikko H Okkonen M Palola S Yriola P Abokangas and M
	Mustonen, "Spectrum sharing using licensed shared access: the concept and its workflow for LTE-advanced networks," IEEE Wireless Communications, vol.
	21, no. 2, pp. 72–79, 2014.
[MVD16]	D. S. Michalopoulos, I. Viering, and L. Du, User plane aspects of multi-
	connectivity in 5G, IEEE International Conference on Telecommunications
	(ICT), Thessaloniki, Greece, 2016.
[MYV+15]	A.Martinez, M. Yannuzzi, J. E. L. de Vergara, R. Serral-Gracià and W. Ramírez, "An Ontology-Based Information Extraction System for bridging the

configuration gap in hybrid SDN environments," 2015 IFIP/IEEE International Symposium on Integrated Network Management (IM), Ottawa, ON, , pp. 441-449, 2015.

- [Nom17] Mx-ART, NOMOR Research GmbH, June 2017. Available at: http://www.nomor.de/
- [OF-CONFIG] ONF TS-016, "OF-CONFIG 1.2: OpenFlow Management and Configuration Protocol 1.2", 2014. Available at: https://www.opennetworking.org/images/stories/downloads/sdn-resources/onfspecifications/openflow-config/of-config-1.2.pdf
- [PDM+11] A. Prasad et al., "Distributed capacity based channel allocation for dense local area deployments", IEEE VTC-Fall, San Francisco, CA, USA, September, 2011.
- [PW96] J. Peng, R. J. Williams, "Incremental Multi-Step Q-Learning", Journal of Machine Learning, Vol. 22, Issue 1-3, pp 283–290, January 1996.
- [SCS16] K. Samdanis, X. Costa-Pérez, and V. Sciancalepore, "From Network Sharing to Multi-tenancy: The 5G Network Slice Broker," IEEE Communications Magazine, July 2016.
- [Rav16] A. Ravanshid et al., "Multi-connectivity functional architectures in 5G", IEEE International Conference on Communications Workshops (ICC), Kuala Lumpur, Malaysia, 2016.
- [RFC6020] M. Bjorklund (editor), "YANG A Data Modeling Language for NETCONF" IETF RFC 6020, October 2010. Available at: http://tools.ietf.org/html/rfc6020
- [RFC6241] R. Enns, M. Bjorklund, J. Schoenwaelder, A. Bierman, "Network Configuration Protocol (NETCONF)", IETF RFC 6241, June 2011. Available at: http://tools.ietf.org/html/rfc6241
- [Ros83] S. M. Ross, "Introduction to Stochastic Dynamic", Academic Press, Cambridge, MA, USA, 1983.
- [SSP+14] K. Samdanis, R. Shrivastava, A. Prasad, P. Rost, D. Grace, "Virtual cells: Enhancing the resource allocation efficiency for TD-LTE", In Vehicular Technology Conference (VTC Fall), 2014 IEEE 80th (pp. 1-5), IEEE, September 2014.
- [TMM+16] A. S. Thyagaturu, A. Mercian, M. P. McGarry, M. Reisslein and W. Kellerer, "Software Defined Optical Networks (SDONs): A Comprehensive Survey," in *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, pp. 2738–2786, Fourthquarter 2016.
- [Wat89] C. Watkins, "Learning from Delayed Rewards", King's College, University of Cambridge, England, May 1989.

PART II: INNOVATIONS

This second part of the final WP4 deliverable provides additional information and evaluation results for the specific WP4 innovations. The following chapters are roughly ordered by topic: Chapters 11 and 12 add further considerations and results for multi-tenancy and resource sharing, Chapters 13 through 16 focus on the various uses of multi-connectivity, Chapters 18 through 19 provide details on protocol and signalling aspects for service flow differentiation, efficient mobility support and massive machine type communication, and Chapter 20 concludes Part II with background information on geolocation databases.

Note that some WP4 innovations are not covered here, either because they have already been conclusively presented before in Part I or in the previous WP4 deliverable, namely:

- RAN support for optimised on-demand adaptive network functions and services ([5GN-D41, Section 7.2]),
- Mobile Edge Computing and Network Resource Allocation for Multi-Tenancy (Part I Section 6.5 and [5GN-D41, Section 7.3]),
- Multi-tenant dynamic resource allocation (Part I Section 5.1 and [5GN-D41, Section 7.5]),
- Multi-RAT Integration ([5GN-D41, Section 7.6]), Load Balancing of Signalling Traffic ([5GN-D41, Section 7.13]).

11 Multi-tenancy in multi-RAT environments

In this section, following the algorithm described in Section 5.2, we are going to describe further obtained results that confirm our theoretical analyses. Additionally, the derived results shed light onto the benefit associated with given approaches.

11.1 Admission control

While the first result given in Section 5.2.1 shows that the proposed algorithm performs close to the optimal one, it has been only compared against two naive policies and thus does not give an insight on the revenue gains that could be achieved over smarter yet not optimal policies. To this end, we compare the performance of our algorithm against a set of "smart" random policies. The fundamental feature of such policies is that inelastic network slices requests are always accepted while the decision of rejecting an elastic request is set randomly. Then, by drawing a high number of random policies, we expect that some of them provide relatively good performance.



Figure 11-1. The distribution of the revenues obtained by random smart policies compared to the proposed algorithms

Figure 11-1 shows the comparison against 1000 different random policies. The results confirm that none of the random policies outperforms our approach. Thus, the optimality of such approach is confirmed, such that substantial gains (around 20 %) are obtained over the random policies.

This result also confirms that a smart heuristic is not effective in optimizing revenue, and very substantial gains can be achieved by using a close to optimal policy such as our adaptive algorithm.

It is worth noting that the previous results have assumed that arrivals and departures follow Poisson process with exponential times, and that the optimal algorithm has a perfect estimation of the statistics of this process. In this evaluation, we address a more realistic case in which neither of these assumption holds. In particular, we introduce two modifications on the arrival and departure model:

- 1. arrivals and departures are Pareto-distributed
- 2. we let the real arrival process λ deviate from the estimated one:

$$\hat{\lambda}(j) = \frac{\lambda}{j+1}$$

That is, the optimal policy obtained by SMDP under the original assumptions is computed offline, with the estimated parameter, and applied to the real system. Note that for negative j values, the

system receives a number of request per time unit higher than the estimated, while positive j values indicate a lower requests arrival rate. The results, depicted in Figure 11-2, show that our adaptive algorithm, which automatically learns the network slice behaviour on the fly and hence is not affected by possible estimation errors, substantially outperforms the optimal policy built upon flawed assumptions and estimations.



Figure 11-2. Revenue in perturbed scenario, $\rho_i / \rho_e = 5$

In conclusion, it can be claimed that our adaptive algorithm approximates the performance of the optimal policy, and provides substantial gains in revenue over potentially smart heuristics. With the proposed approach, we have addressed the need of new resource allocation mechanisms in the new network slicing business model.

11.2 Dynamic resource sharing



Figure 11-3: Normalised utility gain G_W as a function of the maximum allowed number of handovers m

Following the main results given in Section 5.2.2, additional results are provided in this section. As shown in Figure 11-3, the normalised utility gain obtained with m reassociations increase very

sharply with GLLG. Furthermore, its complexity is very small compared to MORA non-linear solver, centralised algorithm with performance guarantees and GD, as shown from the results in Figure 11-4. The results confirm that non-linear solver and centralised algorithms induce considerably higher complexity than DG and GLLG, especially considering that they have to be triggered each time the channel quality of a user changes. By contrast, the execution time for DG is very low and it's even lower for our GLLG approach.



Figure 11-4: Computational complexity of GLLG and SoA algorithms

It is important to note that substantially capacity savings are also obtained, as shown in Figure 11-5. Such capacity gains increase with the number of operators as well as with the density of base stations.



Figure 11-5: Capacity saving

To evaluate the gains from a user perspective, we compare the per-user throughput achieved by our approach against the static slicing (SS) with SINR-based user association (Baseline 1) and SS with enhanced user association (Baseline 2). We observe that our approach provides substantial gains both in terms of the median values as well as the various percentiles. Similarly as above, we observe that such gains increase with the number of the operators.



Figure 11-6: Improvement on the user throughput

In conclusion, it follows from the above analysis that the proposed algorithm guarantees a higher network utility compared with Static Slicing approach. In addition, it provides better performance in terms of capacity saving and requires low computational complexity. The above features render it usable in realistic scenarios.

12 Inter-slice resource sharing

12.1 Motivation and problem statement

The potential price to be paid for enabling network slicing in multi-tenant virtualized mobile networks is the underutilization of the scarce wireless and/or network resources. One way to increase overall network utilization would be to allow inter-tenant sharing, i.e., sharing of resources between different network slices. Such sharing might entail a more complex control layer but has the potential to further push attainable performance gains on the data layer. Therefore, the underlying trade off on the increased computational and communication complexity of the control layer versus the performance gains on the data layer is a goal worth pursuing. To this end, this work tries to shed further light into this issue by discussing different possible degrees of network sharing together with the associated linear integer mathematical programs that allows to investigate upper bounds on the achievable performance improvement. Furthermore, in order to realize real time and adaptive sharing of resources a scale-free heuristic is also presented that is amenable for real-time implementation. The work falls within the remit of dynamic resource slicing.

Based on the "depth" of the aforementioned multi-tenant sharing, we propose hereafter two sharing schemes named tight coupling (TX) and loose coupling (LX). Under 3GPP based scenarios and various assumptions, a set of numerical investigations have been carried out demonstrating the potential for significant gains in aggregated network throughput compared to conventional fully isolated network virtualization methods.

Below we outline different possible sharing configurations between tenants in different slices.

Non Sharing: This can be considered as the case where there is no shared resource between tenants and therefore all resource pooling/management is taking place by an intra-slice controller. This is the nominal use case where a pre-defined slicing (aka partitioning) is taking place according to various business and optimization criteria.

Loose Coupling: This case refers to the situation where resources per tenant are predefined as well as the amount of shared resources are also predefined but without specifying which resources from the available resource pool will constitute the shared ones. Therefore, inter-slice control functionalities need to be in place to support the algorithmic framework of optimal and/or near optimal sharing of the predefined resources between tenants.

Tight Coupling: The last case is when the resources for individual tenant as well as the shared ones are optimized jointly. Computationally this is the most complex one and in terms of architecture this case requires that some of the functionalities of the SDM-C and SDM-X overlap. In other words, if we would like to pool and share resources at the same time the roles of those two elements that would require a joint operation since the SDM-X would require more detailed information per tenant/slice from the SDM-C.

Note that in the loose coupling the SDM-C controller of each tenant provides a fixed set of shared resources and then the role of the SDM-X is to optimize their shared use. On the other hand, in tight coupling these resources are not predefined and the operation of sharing takes into account possible all resources by trying to find the optimal subset for sharing (which might change over time as network conditions evolve). This distinction is depicted via a toy example under the assumption of two tenants in Figure 12-1. In essence, tight coupling can be considered as a generalization of the loose coupling. In addition, if the available set of resources for sharing in the tight coupling is the whole set of resources allocated to a tenant then we have the case of full sharing between tenants.



Figure 12-1: Comparison between loose and tight coupling using a toy example scenario

We note that the above example, as well as the mathematical problem formulation in the sequel, assume the case of PRB inter—slice sharing between tenants but the different degrees of coupling can be equally used in a number of different possible resources between tenants such as spectrum, computational power, memory and communication links.

12.2 Related work

The essence of virtualization is to decouple infrastructure provider (InP) and Service Provider (SP) then multiple SPs can coexist and serve their subscribers/users utilizing the same substrate infrastructure networks [MP97]. In addition, it is expected that network function virtualization (NFV) can largely reduce infrastructure operation expenditure (CAPEX and OPEX, respectively) [LY15], binded with software defined network (SDN) that allows a decoupled control layer and data layer where control layer can be centralized and functional programmable [ONF12]. The centralized controller can be located in either core network resembling a data centre [LMR12], however, with the trend of densely deploying small base station (SBS) in network, the control layer is also considered to be based in the edge (cloud edge) or even at the macro base station (MBS) [IKT12]. Within this context, OpenFlow is an example of SDN providing a flexible network and path control of data forwarding via router softwarization [McK13]. On the control layer, the so-called wireless virtualization controller is implemented with main functionality to orchestrate slices and enforce isolation between them based on the constraints imposed by the physical substrate network to multiple tenants [NEC13], [RBS+16].

12.3 A mathematical programming formulation approach

To model the resource sharing and the different variations for inter-slice sharing as have been detailed in the previous sections and to formulate them via a mathematical programming framework, we proceed by defining a set of binary decision variables. The assumption, without loss of generality, is that we have one macro base station (MBS) and several small base stations (SBSs) based on a fully enabled C/U split architecture. The physical resource blocks (PRB) pool is created and assigned by MBS centrally and SBSs forward data rates to users based on the central assignment. The binary decision variables are as follows,

$$x_{rt} = \begin{cases} 1 & if PRB \ r \ is \ assigned \ to \ tenant \ t \\ 0 & otherwise \end{cases}$$

(1)

$$y_{irt} = \begin{cases} 1 & if \text{ user } i \text{ in tenant } t \text{ uses } PRB r \\ 0 & otherwise \end{cases}$$
(2)

Notations I, C and T are used to indicate the set of user, PRB and tenants respectively. In addition, we also define the following indicator functions to allow for a compact mathematical programming formulation which shows the relationship between each pair of users (UE),

$$\emptyset_{ij}^{t^{i}t^{j}} = \begin{cases} 1 & \text{if users } i, j \text{ of tenants } t^{i} \text{ and } t^{j} \text{ in different cells} \\ 0 & \text{otherwise} \end{cases}$$
(3)

$$\omega_{ij}^{t} = \begin{cases} 1 & if \text{ users } i, j \text{ of same tenant in different cells} \\ 0 & otherwise \end{cases}$$
(4)

In that respect, the SINR (signal to interference and noise ratio) for a given PRB-UE connection in the optimization problem will be estimated using the following expression,

$$\gamma_{irt} = \frac{g_{irt}^{i'} P_t y_{irt}}{\sum_{j \in I'} g_{irt}^{i'} g_{jrt}^{j'} P_t y_{jrt} + \sum_{j \in I'} \omega_{ij}^t g_{irt}^{j'} P_t y_{jrt} + I_{noise}}$$
(5)

In (5), we denote by $g_{irt}^{i'}$ the link gain between the serving base station i' and user i using PRB r. Also, with $g_{irt}^{j'}$ we denote that link gain between the serving base station of user j, which is base station j', and user i (as the interference source). Then the data rate for each association can be calculated as,

$$R_{irt} = \Delta f \log_2(1 + \gamma_{irt}) \tag{6}$$

Based on the above defined preliminaries the problem of maximizing reuse of PRBs via inter tenant resources sharing between different tenants is practically defined as follows,

$$\max \sum_{i \in I} \sum_{r \in C} \sum_{t \in T} \mathcal{Y}_{irt} \tag{7}$$

s. t.
$$\sum_{r \in C} x_{rt} \ge n_t, \forall t \in T$$
 (7a)

$$\sum_{t \in T} x_{rt} \le \beta, \forall r \in C$$
(7b)

$$\sum_{r \in \mathcal{C}} y_{irt} \le \delta, \forall i \in I, \forall t \in T$$
(7c)

$$\sum_{r \in \mathcal{C}} y_{irt} \ge 1, \forall i \in I, \forall t \in T$$
(7d)

$$\sum_{r \in C} \sum_{t \in T} x_{rt} \le \sum_{t \in T} n_t + \alpha \tag{7e}$$

$$y_{irt} \le x_{rt}, \forall i \in I, \forall r \in C, \forall t \in T$$
(7f)

$$\gamma_{irt} \ge \gamma_{th}, \forall i \in I, \forall r \in C, \forall t \in T$$
(7g)

$$y_{irt} + y_{jrt} \le 1 + \omega_{ij}^t \cdot \mathbb{V}, \forall i \in I, \forall r \in C, \forall t \in T$$
(7h)

$$y_{irt} + y_{jrt} \le 1 + \phi_{ii}^{t_i t_j} \cdot \nabla, \forall i \in I, \forall r \in C, \forall t \in T$$
(7i)

$$x_{rt} \in \{0, 1\} \tag{7j}$$

$$y_{irt} \in \{0, 1\} \tag{7k}$$

Constraint (7a) ensures that every tenant will be allocated at least n_t orthogonal PRB. Constraint (7b) denotes that across all tenants a PRB can be reused up to β times, this is to limit maximum aggregated co-channel inter-tenant interference across the available pool of PRBs. Constraint (7c) limits the intra-tenant interference by ensuring that only δ PRB can be used by a user in a tenant. Constraint (7d) indicates that for each user it should at least take one PRB. Constraint in (7e) ensures that only up to α PRBs (a fractional number) can be shared between tenants, this limits the number of PRBs that can be reused. The binding constraint between the decision variables y_{irr}

and x_{rt} is expressed in (7f). The constraint in (7g) expresses the SINR threshold that need to be satisfied in order to reuse. Constraints (7h) and (7i) define that users from the same tenant and different tenants should be allocated a different PRB if they connect to the same cell (same cell avoidance); *V* is an arbitrary large integer that ensures constraints are always satisfied. Constraints (7j) and (7k) indicate the binary decision variables.

12.4 Evaluation

For experiment set-up, different scenarios have been considered to provide a holistic analysis. The local traffic of arriving users is ranged from 50–90. The number of SBSs in MBS controlling area varied from 3 to 15. Due to space limitations results are shown for 9 SBSs (pico cells used in this investigation) scenario since the trends are similar. We also assume that each PRB has the same transmit power, therefore the achieved rate per PRB depends only on the interference level. We use 10 MHz bandwidth system which provides 50 PRBs in total. By ignoring guard band resources we assume all 50 PRBs can be assigned to users. Without loss of generality, the numerical investigations are conducted regarding a substrate network supporting two tenants (named A and B). The simulation parameters used in the numerical investigations presented in the sequel are summarized in Table 12-1 below.

Parameters	Values / Assumptions
Network layout	1 MBS with
	3, 6, 9, 12, 15 SBSs
MBS cell radius (km)	1.5
SBS radius (m)	200
Carrier frequency (GHz)	2
MBS antenna gain (dBi)	14
SBS antenna gain (dBi)	10
Attena configuration	1 Tx for BS, 1 Rx for UE
Thermal noise (dBm/Hz)	-174
System Bandwidth (MHz)	10
No. of PRBs in the pool	50
SBS Path loss (dB)	140.7+36.7log ₁₀ (D) (D in km)
MBS Path loss (dB)	128.1+36.7log ₁₀ (D) (D in km)
Shadowing standard deviation (dB)	6
Max SBS TX power (dBm)	24
MBS TX power (dBm)	43
Number of UEs	50 - 90

Table 12-1: Simulation parameters used for numerical investigations

The simulations are carried out in a use case environment where the number of UEs for each tenant is varied for each Monte-Carlo iteration. More specifically, among different congestion levels of the network (ranging from 50 to 90 UEs), tenant A and B are assumed to have slightly different number of users. Figure 12-2 depicts the sum rate (network throughput) performance of all 5 methods in different congestion level with 2 PRBs inter-tenant sharing ($\alpha = 2$). As can be seen, TX optimal solution provides the best sum rate for all network environment followed by TX heuristic solution. LX optimal and heuristic solutions both have weaker performance compared to TX because the fixed sharing PRBs owned by them might suffer from high interference in some iterations and there is no flexibility for LX to avoid such interference. As expected, the performance of Non Sharing starts degrading after 86 users scenario since the resources in isolated slices become very competitive with the increasing congestion. The resources in fully isolated slice might not even satisfy QoS threshold of some users. In Figure 12-2, the biggest gain compared to Non Sharing method is 20.1 %, 14.7 %, 10.5 % and 10.2 % for TX Optimal, TX Heuristic, LX Optimal and LX Heuristic, respectively.



Figure 12-2: Aggregate rate vs. number of user





Considering the same network environment ($\alpha = 2$), a cumulative distribution function (CDF) of achieved rate for individual UE is studied for 90 users scenario (shown in Figure 12-3). Indicatively, in a relative congested traffic environment, TX Optimal still provides the best performance in satisfying the data requirement of individual user. For the sake of comparison, note that at 1 Mbps point on the horizontal axis, we have that 40 % of users for TX Optimal has data rate lower than 1 Mbps compared to 50 %, 55 %, 58 % and 61 % for TX Heuristic, LX Optimal, LX Heuristic and Non Sharing, respectively. This results indicates that TX (both optimal and heuristic) not only improve network throughput aggressively but has the potential of improving data rate for more users individually in both tenants.

12.5 Concluding remarks

A key challenge in the area of network slicing for future virtualized 5G networks is the potential underutilization of wireless resources, especially during network congestion episodes where full capacity is needed. To this end, and in order to increase utilization of the scarce available substrate wireless network resources optimal and near-optimal inter slice sharing schemes between tenants are proposed in this paper. The proposed optimization framework aims to optimally reuse radio resources via a cross slice SDN-like controller (SDM-X). Two main schemes are detailed, the socalled loose coupling and tight coupling between slices that require different orchestration between the intra and inter-slice controllers. In addition to the mathematical programming formulations, efficient heuristic algorithms are also provided to offer practical solutions amenable for real-time operation. Via a wide set of numerical investigations the gains of the proposed sharing schemes compared to the nominal network slicing approach where resources are predefined and allocated solely to each tenant have been detailed. There are various avenues for extending this research by considering in a more explicit manner other potential resource sharing elements such as for example spectrum. Also, the issue of sharing versus partitioning of resources might strongly depend on the particular characteristics of the resource element and therefore different optimal policies might prevail (for example it might be beneficial for a specific resource to be partitioned whereas for some other sharing might increase the performance).

12.6 References

[MP97]	C. Morin and I. Puaut, "A survey of recoverable distributed shared virtual memory systems", IEEE Transactions on Parallel Distributed Systems, vol. 8, no.
	9, pp. 959–969, September 1997.
[LY15]	C. Liang and F. R. Yu, "Wireless Network Virtualization: A Survey, Some
	Research Issues and Challenges", IEEE Communications Surveys & Tutorials,
	vol. 17, no. 1, pp. 358–380, 2015.
[ONF12]	Open Networking Foundation, "Openflow white paper: Software-defined net-
	working: The new norm for networks", Palo Alto, CA, USA, 2012.
[LMR12]	L. Li, Z. Mao, and J. Rexford, "Toward software-defined cellular networks",
	2012 European Workshop on Software Defined Networking (EWSDN), pp. 712,
	October 2012.
[IKT12]	H. Ishii, Y. Kishiyama, H. Takahashi, "A novel architecture for LTE-B: C-
	plane/U-plane split and Phantom Cell concept", 2012 IEEE Globecom
	Workshop, pp. 624–630, December 2012.
[McK13]	N. McKeown et al., "Openflow: Enabling innovation in campus net- works",
	SIGCOMM Comput. Commun. Rev., vol. 38, no. 2, pp. 69–74, April 2008.
[NEC13]	NEC Coperation, "RAN sharing – NEC's approach towards active radio access
	network sharing", White Paper, Tokyo, Japan, 2013.

[RBS+16] M. Richart, J. Baliosian, J. Serrat, J. L. Gorricho, "Resource Slicing in Virtual Wireless Networks: A Survey", IEEE Transactions on Network and Service Management, vol. 13, no. 3, pp. 462–476, September 2016.

13 Multiple connectivity at the different layers

13.1 Motivation and problem statement

In recent years, the advent of new services like massive Machine Type Communication (mMTC), Internet of Things (IoT), Tactile Internet, Industrial automation and V2V communications has led to a demand surge for mobile and wireless connectivity. The proposed 5G mobile communication system will have to cope not only with the increased traffic demand and service types, it will have to do so within the specified service requirements in terms of latency, reliability and availability. Lately, there has been an increased focus of research efforts to develop solutions to satisfy these high profile service requirements. Novel technology innovations and concepts like direct deviceto-device (D2D) communication, carrier aggregation, microcells, software-defined networking (SDN), coordinated multipoint transmission and reception (CoMP) and dual-/multi-connectivity are some of the promising candidates. As we foresee the future, the network will be heterogeneous consisting of several small cells and legacy systems that are seamlessly integrated, and delivering data cohesively. Multi-connectivity has recently been a hot topic of research due to its ability to enhance throughput, coverage and reliability thereby increasing the overall Quality of Service (QoS). Multi-connectivity refers to a situation where a device's radio resources connect to at-least two different network points and can encompass different radio bands, different Radio Access Technologies (RATs).

LTE-A already supports the concept of dual connectivity that allows simultaneous connection of UE to the two eNBs, Master eNB (MeNB) and Secondary eNB (SeNB) via X2 interface. The data flow is split at the MeNB with some data transmitted from MeNB to the UE, while the other data are transferred over the X2 interface to the SeNB and transmitted to the UE via the SeNB.

We extend the concept of multi-connectivity to small cells without the use of X2 interface, in order to achieve low latency and high reliability. In this regard, we introduce a new function called Inter-RAT Link Controller (IRLC). The purpose of IRLC is to identify and assign a set of RAT's to a UE that maximizes the system throughput while simultaneously satisfying the UE service requirements.

In 5G scenario, every user will have different reliability, latency requirements according to the slice that it will belong. Some users might require ultra-high reliability, while other users might not have strict reliability constraints. Therefore, to utilize the network efficiently, more resources should be provided to the user that require ultra-high reliability and less to the user with lenient reliability constraints. Hence while providing multi-connectivity support, it is necessary to investigate if the user requires multi-connectivity. An algorithm is developed to connect the users to multiple APs simultaneously to ensure the required reliability constraints for a given user is satisfied. The algorithm also provides number of multi-connectivity links for every user, and the operating mode for connected links.

The major contributions of this work are:

- identification of new functionalities for multi-connectivity support,
- reliability analysis with multi-connectivity and
- appropriate operating mode selection for users according to their radio channel condition, to increase network utilization efficiency.

13.2 Architecture and process

Inter-RAT Link Controller (I-RLC):

Currently, the RAT selection function is responsible for initiating handover if another RAT fits better than the current one. The RAT selection function consists of two stages; initiation and

decision making. In the initiation phase, RAT selection function gets periodic inputs like measurement data, available bandwidth, user preferences etc. Many algorithms like Context Aware Multiple Attribute Decision Making (CAMADM), Fuzzy logic algorithms, radio base station efficiency algorithms and ANDSF rules contribute to decision making of RAT selection function. The RAT that suits best according to this algorithm is selected and handover is initiated.

The purpose of IRLC is to identify and assign a set of RAT's to a UE that maximizes the system throughput while simultaneously satisfying the UE service requirements.

The following functionality is also enabled by IRLC:

- AP discovery
- Load and capabilities (Available bandwidth, AP ID, backhaul capacity) of APs.
- Gathering UE information: UE ID, velocity, location, multi-connectivity support, preferences, signal strength, slice to which UE belongs etc.
- RAT selection (cell Assignment: considering UE slice/service requirement).
- UE to multiple AP connection establishment-management.
- Periodically update the UE information.
- Dynamically select operating mode.

Functional architecture:

In 5G NORMA, the IRLC is implemented by the *RAT/Link Selection* subblock of the *RRC User* distributed control function block. It is asynchronous with respect to TTI and controls the PDCP function blocks. The block is activated once per multi RAT connected UE. The block controls the PDCP and PDCP Split bearer data layer function blocks. The prerequisite for its operation is the presence of a common PDCP layer across RATs and UE support for multi-connectivity. Function placement is typically in the edge cloud, cf Section 3.1 and [5GN-D41].



Figure 13-1: C/d-layer architecture for multi-connectivity support in a multi RAT environment

Process description:

Figure 13-2 shows the complete process and interaction of the inter-RAT link selection function block with other function blocks. The selection algorithm, which is executed during this procedure was already described in Section 6.1.3. The respective steps are:

- 1. RRC Cell provides the information of all the neighbouring cells and user preferences if predefined.
- 2. RRC User function requests UE the measurement report about the radio signal conditions from the set of APs, user preferences etc.
- 3. The UE sends the measurement report to RRC User that can be later used for multiconnectivity connections establishment. RRC User now has the requested measurement

report for a set of APs that may belong to different RAT. The measurement information can be categorised into both, intra-RAT as well as inter-RAT.

- 4. RRC User function block selects the RAT by mapping it on to the QoS requirements of UE.
- 5. Depending on the inter-RAT selection algorithm a decision about most feasible RATs connection to UE is made.
 - a. The algorithm selects the inter-RAT link by considering the RAT load conditions, User's preferences towards particular RAT, RAN parameters etc.
 - b. For example, a modified proportional fair algorithm can be implemented to manage the resources across the different RATs.
- 6. UE connection is then establish according to the selected RAT AP list.
- 7. RRC Cell transfers the encryption keys to the PDCP function block.
- 8. PDCP Split Bearer function block is instantiated by RRC User (RAT/Link Selection subblock).
- 9. The RAT/Link Selection subblock of the RRC User function block then selects the operating mode either reliability mode that duplicates the data over multiple connections or diversity mode that splits the data over multiple connections. Depending on the load condition and QoS requirements, the data is either duplicated over the split bearer or multiplexed that is similar to the concept of dual connectivity in LTE.
- 10. As per the operating mode selected, initiate the packet data synchronization on PDCP and PDCP Split Bearer function block.
- 11. Start the data transfer procedure in the PDCP and PDCP Split Bearer function block. PDCP Split Bearer transfers data to the RLC of connected APs.



Figure 13-2: Message sequence chart for the inter-RAT link selection process

14 Data-layer and control-layer design for multiconnectivity

14.1 Motivation and problem statement

The notion of multi-connectivity is a well-known concept and is included already in the LTE standards, in the form of dual-connectivity. Dual connectivity implies that the User Equipments (UEs) are able to connect to two evolved Node Bs (eNBs), simultaneously, thereby exploiting such larger connectivity potential by aggregating the individual connections into a faster, unified one. Dual connectivity is thus able to provide considerably larger throughputs in LTE, particularly concerning HetNet deployment scenarios where small cells are deployed within the coverage area of macro cells.



Figure 14-1 The LTE RAN Architecture

The main points pertaining to the RAN architecture considered in LTE standards are as follows:

- a. The LTE RAN architecture consists of eNBs which, from the UE perspective, comprise the data layer and control layer network termination points. That is, the network functions of the core network (known as Evolved Packet Core (EPC) Network) are not visible to the UEs.
- b. The eNBs may be interconnected with one another via a special inter-node interface, known as the X2 interface. The X2 interface is set such that inter-eNB functions are supported, including inter-cell interference coordination (ICIC), enhanced mobility, as well as several functions which belong to the broad family of Self Organized Network (SON) functions. These basic architectural elements of the LTE RAN are depicted in Figure 14-1.
- c. The eNBs support the full protocol stack, consisting of the following protocol layers: Packet Data Convergence Protocol (PDCP); Radio Link Control (RLC); Medium Access Control (MAC); Physical Layer (PHY). On top of these protocol layers the Radio Resource Control (RRC) layer is placed; its functionality is related exclusively to control layer messages. This architecture view is depicted in the upper part of Figure 14-2.
- d. The interface between the eNBs and the EPC network is defined as the S1 interface.
- e. Overall, it can be argued that *the LTE RAN architecture is fully decentralized*. That is, the eNBs are *standalone entities* which, albeit interconnected, they are independently connected to the core network.

Nevertheless, as also reported in [5GN-D41], the dual-connectivity approach is associated with some drawbacks, which are particularly related to an increased signalling overhead due to the frequent mobility events such as channel measurements and handovers between the small cells [5GN-D41].

14.2 Related work

To this end, a rather centralized architecture has been proposed in [5GN-D41], where the main differentiation is that the RRC (control) and the PDCP layer are common for multiple eNBs and located in a centralized location known as "edge cloud" (c.f. Figure 14-2). The main advantages that stem from such architecture are two-fold: First, the core network is completely isolated from any mobility events in the RAN, since no path switching occurs each time a handover takes place between the small cells. Second, this architecture facilitates special techniques dedicated to achieve higher reliability levels. This is because a centralized PDCP layer facilitates the coordination of data duplication techniques, thereby attaining the beneficial effects of diversity which ultimately result in considerably lower error rates and latencies. Such approach is presented in detail in [5GN-D41], while security aspects are elaborated in Section 5.3.2 of [5GN-D32].



Figure 14-2 The 5G RAN Architecture considered in [5GN-D41] vs the existing LTE RAN Architecture

A similar deployment approach for HetNets is also provided in [NGMN-SBH]. Although the architecture is still oriented towards a pure LTE network, it contains the fundamental ideas for providing transport connectivity to small cells, and identifies two major options: a) Backhauling via the macro site, where the aggregation connectivity point of small cells is a nearby macro cell; b) Backhauling via a separate site, where such aggregation point is located separately from the macro cells. Such options are summarized in Figure 14-3.



Figure 14-3 Transport connectivity options to small cells [NGMN-SBH]

14.3 Architectural approaches



Dedicated aggregation for small cells

Centralized scenario

Figure 14-4: The small-cell aggregation and centralized scenario architectural approach

Taking the above considerations into account, one may argue that as far as multi-connectivity is concerned, the related 5G architecture should incorporate the available topological flexibility as well as the available backhaul capabilities. In this regard, and on the basis of the RAN architecture discussed in [5GN-D41], there are two major architectural variations that can be applied. In particular, the two proposed approaches are as shown in Figure 14-4 and described below:

• Approach 1- Dedicated aggregation for small cells: A data centre is used between the core cloud and the (macro and small) cells. The small cells are further organized into clusters such that small cells of the same geographical region belong to the same cluster. For each cluster, the data layer is connected to dedicated aggregation points, while the

control layer remains in the macro cell. The connection between the macro cell and the small cells for both the data and control layer is established via the aggregation point (for each cluster of small cells) and the data centre.

• *Approach 2 – Centralized scenario:* The macro and small cells use a common aggregation point. Such topology corresponds to a generalized version of the cloud-RAN scenario, where both the data layer and the control layer are centralized.

14.4 Throughput evaluation of multi-connectivity

The above two approaches are associated to different performance features due to the different topologies involved. Here, we provide an evaluation of the two considered approaches in terms of the throughput that is achievable by the corresponding multi-connectivity scenario.

Simulation scenario: Simulations have been conducted which involve a HetNet deployment consisting of 21 macro cells and 84 small cells; the small cells are organized in clusters of four, such that the within the coverage area of each macro cell four small cells are used with overlapping coverage. The macro cells are assumed to be typical wide-area 5G macro cells, operating at a central frequency of 5.9 GHz with a 20 MHz bandwidth. The small cells are assumed to operate in the millimetre-wave band, with central frequency of 28 GHz and 100 MHz bandwidth. Each UE within the overlapping area of a small cell and a macro cell is connected to both access points, in a dual-connectivity fashion.

Traffic flow control considerations: A traffic flow control mechanism is used, which is assumed to be located a) at the macro cell for approach 1; b) at the aggregation point for approach 2. This flow control mechanism is used to ensure that enough traffic is sent to both links, depending on the buffer status of both types of cells and the channel quality of the links involved.

The throughput performance of the two considered approaches is depicted in Figure 14-5 and Figure 14-6, obtained from the simulation scenario described above. The fronthaul delay values correspond to the delay between the small cell clusters to the data centre (which is assumed colocated with the aggregation point in approach 2), as well as the delay between the macro cell and the data centre. The throughput is shown in the cumulative distribution function (CDF) form, such that the curves correspond to the probability that the throughputs achieved are at maximum the projection values to the x-axis.

With such illustration, it is easy to identify the percentage of time that the network performs above a given throughput requirement. For instance, we notice that, for approach 1 and for the ideal case of 0 ms delay, 50 % of the time the achievable application-layer throughput is above 300 Mbps, whereas for 1ms fronthaul delay such percentage drops to approximately 20 %. Interestingly, approach 2 seems more robust to fronthaul delay as can be deduced by comparing Figure 14-5 and Figure 14-6. For the example discussed above, the percentage of time that the throughput is larger than 300 Mbps is again approximately 50 % for (the ideal case of) 0 ms delay, however for 1 ms delay the corresponding value is approximately 40 %, i.e., twice as much as for approach 1.

The above analysis reveals that a centralized architectural approach offers a relative robustness to multi-connectivity against delays that may occur in the link between the access points, where RAN is realized, and the network aggregation points. Such architectural considerations are important when assessing the performance of multi-connectivity, since, together with the design of the RAN, they represent a crucial limitation factor to the effectiveness of the network to meet its requirements.



Figure 14-5: Throughput performance (at application layer) of approach 1



Figure 14-6: Throughput performance (at application layer) of approach 2

14.5 References

[NGMN-SBH] Next Generation Mobile Networks (NGMN) Alliance, Small Cell Backhaul Requirements, White Paper, June 2012.

15 Architectural approaches for multiconnectivity of mm-wave APs and 5G low band

15.1 Motivation and problem statement

One of the key technology components envisioned for the delivery of extreme mobile broad band in 5G networks is the radio access using millimetre waves. However, using this technology, it is challenging to provide reliable and seamless data transmission to the users. Due to the propagation characteristics of millimetre waves, high interruption times may occur if line of sight is obstructed between the transmitter and the receiver. Therefore, we propose that cellular deployments of systems using millimetre waves should be supported by the low-band (below 6 GHz) coverage to enable a high reliable control layer connection to the mm-wave terminal.

Due to the high bitrate, limitations with respect to deployment, e.g. backhaul, processing power etc. have to be taken into account when defining an architecture supporting a mm-wave transmission.

15.2 Major results

Different architectural approaches have been analysed and evaluated, to enable a deployment with dedicated hardware, i.e. a standalone 5G LB and standalone mm-wave APs (mmAPs) taking into account limitations with respect to e.g. processing of mm-wave data within a 5G LB or e.g. backhaul limitations between a 5G LB and mmAPs.

Other alternatives like a cloud based implementation for the PDCP processing of mm-wave data or a specific hardware for distribution of mm-wave PDCP data have also been investigated. The different alternatives have been analysed by counting the required number of signalling messages for a setup and a reconfiguration of a mm-wave cluster.

Within [ABA+16] the architectural approaches and results are published.

15.3 Architectural approaches for provisioning of 5GmmAPs

Different architectural approaches taking into account backhaul limitations, processing power, the placement of RRC and the split of data flows have been investigated and analysed with respect to the required messages for a given scenario. For all approaches it is assumed that the RRC messages towards the UE are transmitted via the 5G low band (5G LB) to guaranty a reliability inside the defined UE cluster consisting of the 5G LB and mmAPs The location and processing of the RRC stacks differs between the investigated approaches.

Approach-1a is the straight forward approach based on LTE dual connectivity where the 5G eNB is hosting the RRC and also manages the PDCP split, see Figure 15-1. This approach requires a backhaul which can carry mm-wave data and a powerful baseband processing inside the 5G eNB. In case of sporadic mm-wave transmissions the processing power might be overprovisioned.



Figure 15-1: PDCP based MC with two variants for location of RRC protocol stack

A potential variation, i.e. Approach-1b is depicted in Figure 15-1 which has relaxed requirements to the backhaul and towards the 5G LB AP and the PDCP processing capability of the LB AP. A 5G Radio Network Controller (5G-RNC) acts as a RRC-Host and PDCP manager. The RRC connection with UE is maintained by 5G-RNC via 5G-LB AP. The 5G-RNC could also be considered as a function inside the NORMA defined edge cloud. All 5G-mmAPs of the preconfigured cluster are directly connected to the 5G RNC which has an interface to core network. The protocol stack for both approach is also depicted in Figure 15-1. Approach-1 offers a fall-back solution for the transfer of mm-wave data – with reduced bitrate – to an UE in case it is outside the coverage of mmAPs.

In case of backhaul limitation to the 5G LB node it should still be possible to couple a cluster of 5G mmAP to a 5G LB node. Approach-2 introduces a further split in data layer by introduction of a new entity or a function in edge cloud which is called a Radio Network Data Distributor (RNDD) as illustrated in Figure 15-2. The data layer is splitted in high data rate and low data rate flows in accordance with the definition of a PDCP-H and a PDCP-L respectively for high and low data rates. The 5G-LB Node hosts the PDCP-L and RRC whereas the RNDD will host the PDCP-H. The 5G-LB Node which is the RRC-Host will manage the MC-cluster via RNDD. Being the RRC-host, it will also control the traffic steering in RNDD which is hosting PDCP-H. This will require new and standardized interface between 5G-LB node and RNDD.



Figure 15-2: Architecture approach -2: split in PDCP_H and PDCP_L and related protocol stack

Using this interface AMMC and LMMC functions can also be applied. In this example, AMMC function assigns mmAP1 and mmAP2 as serving APs whereas mmAP3 is prepared with only the context of the UE. The protocol stack for approach-2 is also shown in Figure 15-2.

The core network already performs traffic steering decisions based on the service requirements and routes the traffic with e.g. low data rate requirements to the 5G-LB node while traffic with high data rate requirements is routed towards the RNDD, which hosts the PDCP-H. Inside the PDCP-H, the high data rate traffic is split between 5GmmAP1 and 5G-mmAP2. The traffic steering within PDCP-H is done by the RRM entity hosted inside the 5G-LB node.

The approach-3a and approach 3b are based on common MAC based Multi Connectivity provided by the lower radio layers, see Figure 15-3. A 5G mmW Central Access Point (5G-mmCAP) is shown which is connected to multiple mmRRHs. The 5G-mmCAP splits the data flow on MAC level. The traffic steering is controlled by a RRC. PDCP and RRC could be hosted by either 5G-LB Node (approach-3a) or a cloud entity (approach-3b). In principle, MC-cluster can contain more than one 5G-mmCAPs but in Figure 15-3 only a single 5G-mmCAP is depicted. The interface between 5G-mmCAP and mmRRH can be I/Q samples (e.g. CPRI) or preferably above PHY (MAC/PHY split as CPRI is not cost efficient due to high data rate requirements). The 5G-LB node and 5G-mmCAP perform independent scheduling. A centralized scheduler in 5GmmCAP is assumed for all mmRRHs within the cluster. Using the AMMC, LMMC functions, the RRC manages the configuration and management of 5G-mmCAPs and their respective mmRRHs. In Figure 15-3 the AMMC selects mmRRH1 and mmRRH2 associated with 5GmmCAP as servers. A link switch to mmRRH3 can be easily executed if LMMC detects mmRRH3 with better quality. In case of link switch, AMMC can update the MCcluster info.



Approach-3a Approach-3b

Figure 15-3: Architecture approach-3: MAC based MC and related protocol stack

The related protocol stack in downlink shows an example of hierarchical MC, where at first stage MC using common PDCP is applied and in the second stage MC using common MAC is applied. In case of 5G-mmAP outage, 5GLB node can still be used as a fall-back solution but with low data rates.

15.4 Evaluation



Figure 15-4: Signalling messages flow between the node and the UE for approach-1a

The assessment of the different approaches was focused on the qualitative estimation of reconfiguration time with the help of signalling messages exchange required for setup and reconfiguration, e.g. due to handover, by each approach. The preparation and update phase of the MC-cluster configuration is performed without core network involvement, i.e. messages are directly exchanged between the nodes in 5G-RAN. An example of the required messages for approach-1a is shown in Figure 15-4:

- 1: The 5G-LB node configures the UE measurement procedures according to the related information e.g. list of mmAPs.
- 2: A measurement report (indicating suitable mmAPs) is triggered and sent to the 5G-LB node by the UE.

- 3–5: With the help of measurement report and using AMMC functionality, 5G-LB node defines the set of 5G mmAPs in the MC-Cluster and their respective activity levels. The 5G-LB then issues request for addition in MCcluster to each 5G-mmAP passing necessary information.
- 6–8: Admission control may be performed by each 5GmmAP depending on the received activity level request, if the resources can be granted by the 5G-mmAP. Here we assume that all the 5G-mmAPs that have been asked by the 5G-LB node accept the request with the required activity level. Each 5G-mmAPs sends the acknowledgement along with the required related configurations. In this example, 5G-mmAP1 and 5G-mmAP2 will be involved in transmission so they send random access preambles.
- 9: The MC-cluster configuration is passed to the UE by the 5G-LB e.g. using RRC.
- 10: The UE receives the configurations and acknowledges the reception.
- 11–12: Using the allocated preambles, the UE performs the random access to 5G-mmAP1 and 5G-mmAP2.
- 13–14: 5G-mmAP1 and 5G-mmAP2 provide random access response to the UE which contains scheduling information.

The UE is receiving data through 5G-mmAP1 and 5G mmAP2. Using the LMMC functionality, the 5G-mmAP1 achieves an early detection of link degradation, which triggers the following procedure:

- 15: The 5G-mmAP1 sends link quality degradation indication to 5G-LB node.
- 16: The 5G-LB node requests the UE for proactive measurement report.
- 17: The UE sends the measurement report to the 5G-LB node.
- 18–19: With the help of measurement report and using AMMC functionality, 5G-LB node updates the MCcluster and issues accordingly the modification requests to the influences mmAPs. Here we assume that after modifications, mmAP3 would activate the transmission.
- 20–21: The requested mmAPs respond to the 5G-LB along with the related information. For example, 5GmmAP3 will provide the random access preamble.
- 22: The MC-cluster configuration update is passed to the UE by the 5G-LB.
- 23: The UE receives the configurations and acknowledges the reception.
- 24: Using the allocated preamble, the UE performs the random access to 5G-mmAP3.
- 25: 5G-mmAP3 provides random access response.

At the end of the procedure, 5G-mmAP1 suspends data transmission but 5G-mmAP2 will continue to transmit to guarantee seamless data transmission. Due to modifications, 5G mmAP3 is also set as a server and now the UE receives data from 5G mmAP2 and 5G-mmAP3. During the reconfiguration procedure, a message can be send to 5G mmAP2 optionally by 5G-LB in order to regulate the percentage of data flow.

The comparison of messages required for each approach is shown by the graphs in Figure 15-5. The highest number of messages is required by approach-2 which considers the deployment aspect with respect to the backhaul limitations and provides MC with the help of RNDD. The lowest number of messages are required by approach-3a which assumes ideal backhaul capacities between the nodes and overlooks the complexity. Slightly less messages are required by approach-1b than approach-2. This approach provides a good compromise in terms of the realization of deployments of mmAPs with required backhaul and their coupling with 5G-LB AP for reliable control.



Figure 15-5: Comparison of number of messages between different approaches

15.5 Security considerations

All approaches are feasible security-wise and can be supported by the AS security concept described in [5GN-D32] and [5GN-D33].

Security-wise, approaches were the PDCP is not located on a physically exposed entity like an 5G AP or LTE eNB have a slight advantage. In the other variants, solutions must be implemented that make it hard for attackers with physical access to the equipment to get access to the cleartext traffic. Good solutions for this may be expensive, however.

15.6 References

[AGA+16] D. Aziz, J. Gebert, A. Ambrosy, H. Bakker and H. Halbauer, "Architecture Approaches for 5G Millimetre Wave Access Assisted by 5G Low-Band using Multi-Connectivity", Globecom Workshop Washington, DC, USA, pp 1–6, December 2016.

16 Virtual cells and multi-cell coordination

16.1 Security

The proposed approach for formation of virtual cell is mainly dependant on the reports and measurement of the UEs. There may be some issues when terminals send wrong reports. If many terminals behave this way, they may be able to prevent the virtual cell from working properly, leading to a reduced service quality for other terminals in this area. However, the issue of wrong reports sent by terminals in not a new issue, and it is already addressed in [33.401].

16.2 Numeric results

In the first part of this document, the concept of virtual cell is presented. The edge-cell threshold is the key parameter in formation of virtual cells. In addition, the numeric results showed that the results are highly dependent on the placement of MTs in the deployment area. Therefore, in the next scenarios, the placements of the MT are changed from uniform to study the effect of it.

Numeric results for Scenario B:

The total network throughput for scenario B has been presented in Figure 16-1. Referring to the figure, the total network throughput has increased up to 34 % of its initial value. It can be seen throw the graph, that the optimal value for the edge-cell threshold is changing from 15 dB to 18 dB.



Figure 16-1: The network throughput increases as a function of the cell edge threshold (Scenario B)

Numeric results for the scenario C:

Moreover, the results for scenario C have been presented where distribution of terminal per eNB is non-uniform keeping the cell edge and cell centric configuration as it is. Figure 16-1 presents the increment of the total network throughput for the scenario C.



Figure 16-2: The network throughput increases as a function of the cell edge threshold (Scenario C)

The numeric results show the increased total network throughput in this scenario, up to 45 % of its initial value. The formation of the virtual cell allows the network to use the available resources in the cell with lower traffic demand to serve the cell, which was congested. Hence, the total network throughput has increased.

The maximum network throughput, comparing the Figure 16-1 with Figure 16-2, has achieved with smaller edge-cell threshold. For instance, the maximum network throughput when there are 100 MT deployed in the reference scenario is achieved by setting the threshold to 21 dB. However, the maximum throughput for the same number of MT in scenario A is achieved when the threshold is set to 18 dB. The reason behind this change is the density of the MT in the edges where the virtual cell is going to form. The higher density of MT in outer ring, the lower threshold is required.

Numeric results using the Q-learning:

As it is described in part I, the formation of virtual cells using the Q-learning technique has been studied. It worth reminding that the scenario considered has two cells. The UEs in cell 1 are all active in uplink and the UEs in cell 2 are all active in downlink direction. The changes of throughput in downlink direction of these cells are presented in Figure 16-3.





It can be seen from the graph that the formation of virtual cell increased the throughput in down link direction of cell 1. This means that the unallocated resources in the cell 1 have been allocated to the terminal in cell 2 (i.e. the terminal active in downlink). However, the formation of the virtual cell increased the interference level in cell 2 and reduced the downlink throughput in this cell. It can be seen that the throughput in downlink direction in cell 2 had been reduced about 20 Mbps, in one hand and on the other hand the throughput in the cell 1 increased up to 45 Mbps. According to the figure the total network throughput in downlink direction has increased by 25 Mbps.

Finally, the Figure 16-4 presents the throughput of the cells in uplink. As expected the uplink in cell 2 has increased up to 5 Mbps in cell 2 while the cell 1 experienced a decrement of 3 dB. Since the terminals transmission power are considered comparatively smaller than the eNBs, the improvement of throughput and the imposed interferences are also smaller comparing to the downlink.



Figure 16-4: The throughput of the cells in download direction

In conclusion, the formation of the virtual cell can improve the total network throughput with cost of increasing the interferences. The numeric results show the increment of the throughput in the downlink comparing to the uplink is bigger. It worth noting that, the same scenario has been studied with the focus on improving the network latency. The improvement of latency as the result of virtual cell formation is not considerable. However, it offers more flexibility in managing the radio resources, which enables the patterns for the cell to be chosen based on latency requirements and then compensate the loss of network throughput by forming virtual cell.

16.3 The implementation of the Virtual Cell in Demo 1

The software part demonstrates how the SDM-C can be coupled with a mobile network simulator. This enables the SDM-C to control software simulated eNBs and reconfigure them into central or edge-cloud based on the network parameters and SDM-C's decision logic. The main motivation for a software demo is to highlight the effects of 5GNORMA innovations on a larger scale, i.e. with many eNBs and UEs, and to investigate the gains on network-level, which could not be easily observed with only a few hardware eNBs.

The main part of the demo is dedicated to present the improvement of QoS by changing the placement of function network elements. As it is shown in Figure 16-5, the SDM-C received periodic reports and base on the service demands decides to put the network function elements in the edge-cloud and reduced the latency or push them to the central cloud and use centralised interference mitigation.



Figure 16-5: Network reconfiguration into edge-cloud

The overall layout of the software demo developed for the current stage is shown in Figure 16-6. The demo comprises of three main parts. The Software Defined Mobile Network Controller application is shown on the top left of the figure. This is the simplified graphical user interface (GUI) for the SDM-C application. The second part of the demo is the network simulator, which is on the right bottom of the screen. The simulator provides a 2D view of the simulated MTs, the eNBs and their respective cell layout.



Figure 16-6: Software demo with Edge-Cloud Configuration

Two centralise radio resource management algorithm is going to be demonstrated in this demo, which are:

• Beamforming with interference Mitigation: The beamforming algorithm implemented in the demo is helpful to improve the Signal-to-Interference-plus-Noise Ratio (SINR) of UEs by directing the beams towards the desired UE and reducing the interference to

other UEs. This is used when eNBs are placed in central-cloud configuration. The beamforming mechanism makes use of a linear array of antennas at the eNB with up to four antenna elements and used primarily for downlink. The division of the cells into regions is highlighted in Figure 16-7. The implemented algorithm is presented with more details in [5GN-D62].



Figure 16-7: The beam coordination method

• Formation of virtual cell: the aforementioned algorithm for formation of virtual cell is already implemented in the simulator. However, for the demonstration proposes, the related controller will be added to the GUI, in order to show the effect of forming virtual cells. The demonstration is going to have three phases regarding the virtual cell, which are: no virtual cell, manual selecting of edge-cell threshold, and the Q-learning based approach. The SDM-C GUI may be updated relatively.

16.4 References

[33.401] 3GPP TS 33.401, "3GPP System Architecture Evolution (SAE); Security architecture (Release 15)", June 2017.

17 Flexible 5G service-flow (SF) with in-SF QoS differentiation and multi-connectivity

This Chapter is a continuation of the work presented in the previous WP4 deliverable [5GN-D41] Section 7.8, where 4 contributions were presented:

- 1. facilitating in-bearer QoS differentiation;
- 2. d-layer enhancement for optimised QoS support;
- 3. c-layer for flexible radio multi-connectivity; and
- 4. MAC-level multi-connectivity for ultra-dense 5G networks.

Here, we add 2 new contributions, focusing on RAN support for network slicing in 5G:

- 5. PDCP multiplexing of SFs from different network slices; and
- 6. Coordinated dynamic scheduling and semi-persistent scheduling based allocations.

The network slicing has been considered as one of the most important key issues for 5G network systems in 5G NORMA [5GN-D31] as well as in 3GPP [23.799]. This chapter considers RAN level support for network slicing and, in particular, enabling and facilitating that a single UE may have different PDU sessions, also referred to as service flows (SF) herein, from different network slices being served by the shared RAN – same radio access node or base station, also referred to as gNB in 3GPP, on same or shared carrier(s). That is, RAN slicing Option 3 described in Section 2.1.3 is assumed.

17.1 PDCP multiplexing of SFs from different network slices

Motivation and problem statement:

The UE may be provided with only a limited number of radio bearers (RB) including signalling RBs (SRB) and data RBs (DRB) by the serving RAN for the sake of scalability and interoperability, regardless of how many network slices the UE may access simultaneously. Thus, PDCP multiplexing of SFs from different network slices which have more or less the same QoS requirements from RAN perspective into one RB is a potential option. The PDCP multiplexing also has minimum impact on lower layers of RAN protocol stack and therefore requires the least of standardization and implementation efforts for RAN.

Major outcomes:

This contribution therefore provides an efficient scheme for facilitating the PDCP multiplexing of SFs from different network slices into a RB. The scheme includes new RAN functions and procedures which are optimized in terms of protocol overhead and operation for individual multiplexed SFs. The scheme also allows for possible extension to enable slice-specific RRC functions and procedures, here referred to as RRC_Slice and implemented by the *RRC Slice* SDM-C application, to be multiplexed and signalled over a common SRB.

Related work:

In LTE, multiplexing of different SFs into a RB is not considered as RAN function or, that is, PDCP does not need to be aware of SFs.

Technical details:

The PDCP multiplexing requires that the network slice instance (NSI) specific contexts of individual SFs being multiplexed into the configured RB are applied and maintained. To enable and facilitate that, a straightforward solution is to have some identifier (ID) and sequence number (SN) per a NSI specific SF, denoted as NSI-SF ID and NSI-SF SN, included in the control header of individual PDCP PDUs being transmitted on the configured RB, in addition to SN of the

configured RB, denoted as RB SN. Note that RB ID may be omitted from the header of PDCP PDU as in LTE, assuming that there is 1:1 mapping between a RB and a logical channel (LC) of MAC and LC ID is included in the header of MAC PDU. However, the straightforward solution requires that the header of individual PDCP PDU carries two SN fields plus a NSI-SF ID, even when there is no actual multiplexing of different SFs from different NSIs, or no such multiplexing takes place on the RB, resulting in high protocol overhead. It is noted that if NSI-SFs are always independent and PDCP ciphering is always NSI specific then RB-SN may not be needed. In this case, PDCP is operating on individual NSI-SFs.

PDCP specific to NSI-SF being multiplexed	NSI specific operation (ciphering, header compression using NSI specific SN)
PDCP common to RB	RB wise operation (ciphering, in- order delivery with possible NSI prioritization using RB wise SN)

Figure 17-1: PDCP structure for PDCP multiplexing of SFs from different NSIs into one RB



Figure 17-2: PDCP multiplexing of SFs from different NSIs into one RB

Figure 17-1 and Figure 17-2 illustrate some basic PDCP structure and operation behind PDCP multiplexing of SFs from 3 different network slice instances (NSI) into one configured RB for radio packet transmission between a given UE and a serving RAN. The PDCP multiplexing function, based on pre-configured properties of individual SFs from corresponding individual NSIs coupled with pre-configured QoS scheduling rules applied on the RB across all the

multiplexed SFs, is able to sort out and maintain the queue of incoming packets from the individual SFs for transmitting on the RB. It is noted that at least SN and sufficient multiplexing control information per a PDCP PDU need to be provided to facilitate, e.g., duplication-detected-and-discarded, lossless and in-order delivery for the RB as well as de-multiplexing and routing of individual PDCP PDUs received on the RB for the belonging individual SFs and serving NSIs.

The RAN should be able to support NSI specific configurations and functions to be applied on individual SF of serving NSI. For example, QoS and security context of a SF can be specific to serving NSI. In this regard, the PDCP multiplexing needs to facilitate NSI specific configurations and functions on a corresponding SF and therefore maintain NSI specific contexts of individual SFs.

It is also highly expected in 5G that the given UE is provided with multi-connectivity wherein the multiplexed RB may be split and served by different access nodes. Hence, the receiving PDCP peer entity may receive PDCP PDUs in different orders from different serving access nodes due to, e.g., different delays thereof. The in-sequence delivery operation of the multiplexed RB based on RB wise SN (RB-SN) overall may cause unnecessary delay to some of the individual multiplexed NSI-SFs. Hence, it is beneficial that the receiving peer may be able to check and proceed with in-sequence delivery for at least some determined or prioritized NSI-SF without further delay.

Based on the above considerations, the proposed PDCP multiplexing scheme herein consists of the following:

- **NSI-SF ID configuration with zero signalling overhead.** Because there is only one NSI-SF from each NSI being multiplexed into the configured RB, NSI-SF ID is set to the specified RAN level ID of the corresponding NSI, referred to as RAN-NSI-ID. RAN-NSI-ID can be derived from the pre-configured network ID of the corresponding NSI which is unique to the given UE and the configured RB. NSI-SF ID may be included in the header of PDCP PDUs transmitted on the configured RB only when actual multiplexing of different SFs from different NSIs is taking place on the configured RB.
- **Facilitating efficient maintaining of NSI specific contexts.** SN of packets per a multiplexed SF is not to be included in each and every PDCP PDU but only some selected ones. This is based on exploring trade-offs between using explicit per-PDU signalling (control elements included in the header of each PDU) and once-off c/d-layer control signalling procedures for maintaining NSI-SF specific contexts only when needed. There can be different alternatives for realization.

Figure 17-3 illustrates a RB reconfiguration procedure to add NSI-SF#2 to be multiplexed into the existing RB carrying NSI-SF#1. The format of PDCP PDUs over the RB is changed when the actual multiplexing is taking place after the RB reconfiguration.

Figure 17-4 illustrates an inquiry of the SN mapping of a particular in-order but out-of-sequence received PDCP PDU. For instance, let us assume that the UE receives PDCP PDU of RB-SN=9 but has not received PDCP PDUs of RB-SN=6, 7, 8; and the RB is configured for delay sensitive services. It is noted that, for a high-reliable but delay tolerable RB, the UE may be configured to proceed with the configured RB wise reordering window and not need to consider advancing for individual NSI-SFs. In this case, there may be no need for initiating on-the-fly inquiry of a certain received out-of-sequence PDU.



Figure 17-3: Adding a NSI-SF into a multiplexed RB



Figure 17-4: Illustration of SN mapping inquiry procedure of PDCP



Figure 17-5: Possible extension to cover NSI specific c-layer RRC, referred to as RRC_Slice

Figure 17-5 illustrates a possible extension of the scheme to cover RRC_Slice multiplexing on a common SRB. Option (a) is to have PDCP multiplexing of RRC_Slice# on a specified SRB directly. Option (b) makes NSI specific RRC internal to RRC and hidden from PDCP. In this option, RRC is responsible for addressing NSI specific RRC information elements (IE) or messages (PDU).

Security considerations:

This contribution assumes that the shared or common RAN is in charge of setting up security for the proposed multiplexing RB. However, on top of that, NSI specific security may also be applied for the individual NSI-SFs being multiplexed into the RB. That is, 2 levels of PDCP security may be provided: higher level is specific to the individual NSI and lower level is specific to the common serving RAN.

Remarks:

This contribution has considered many practical aspects and provided an efficient scheme for PDCP multiplexing of SFs from different RAN sharing network slices for a given UE accessing those different network slices simultaneously. The scheme has certain implications on RRC in c-layer and MAC/PHY in d-layer of the serving RAN being shared by different network slices. The scheme can be incorporated in 5G RAN systems being standardized in 3GPP. Thus, the scheme may have potential for standards contributions to 3GPP for examples.

17.2 Coordinated dynamic scheduling and semipersistent scheduling based allocations

Motivation and problem statement:

Though RAN may be commonly shared for multiple network slices, certain resource isolation from one to another network slice may still be fully or partially required in order to meet the requirements of individual network slices and services. Thus, on the one hand, dynamic scheduling based resource allocation over the shared RAN resources is expectably the most spectral efficient and flexible option. On the other hand, some network slice may require to have guaranteed resource allocation in order to meet its user and service requirements. In another example, the tenant that wants to have full vertical control of the network slice may request a dedicated resource allocation for serving its users.

Major outcomes:

This contribution provides an efficient scheme for coordinating the semi-persistent scheduling and dynamic scheduling based resource allocations for the same UE accessing multiple network slices simultaneously. In particular, when dynamic resource allocation is also scheduled for the same TTI in which a slice specific SPS resource is allocated, the dynamic resource allocation is used for the UE in a coordinated way which allows for differentiating the slice with SPS allocation while maximizing utility of the scheduled resources for corresponding slices, instead of overriding the scheduled SPS allocation as currently specified in LTE. The scheme can also be used to facilitate resource reservation and separation among different network slices for the purpose of, e.g., ensuring that some network slices have certain guaranteed resources for their services and at the same time to allow flexible on-demand resource usage among multiple network slices.

Related work:

The dynamic scheduling (DS) based resource allocation refers to the mode in which the serving RAN allocates resources and transport formats to the UE for radio transmissions/receptions dynamically per each scheduled TTI. The semi-persistent scheduling (SPS) based resource allocation refers to the mode in which the serving RAN allocates at least a part of resources and transport formats to the UE semi-statically over a certain time interval consisting of a number of TTIs. In LTE, the dynamic scheduling and SPS based resource allocations are exclusive [36.300]. E-UTRAN can allocate semi-persistent downlink resources for the first HARQ transmissions to UEs:

- RRC defines the periodicity of the semi-persistent downlink grant;
- PDCCH indicates whether the downlink grant is a semi-persistent one i.e. whether it can be implicitly reused in the following TTIs according to the periodicity defined by RRC.

HARQ retransmissions when needed are explicitly signalled via the PDCCH(s). In the sub-frames where the UE has semi-persistent downlink resource, if the UE cannot find its C-RNTI on the PDCCH(s), a downlink transmission according to the semi-persistent allocation that the UE has been assigned in the TTI is assumed. Otherwise, in the sub-frames where the UE has semi-persistent downlink resource, if the UE finds its C-RNTI on the PDCCH(s), the PDCCH allocation overrides the semi-persistent allocation for that TTI and the UE does not decode the semi-persistent resources.

Technical details:

For supporting the UE accessing multiple network slices simultaneously in which the UE is configured with UE specific SPS or NSI specific SPS configuration for the given NSI and given service thereof and, at the same time, dynamic resource allocation scheduled per a transmission for other services of the same or different NSIs, the following coordination between the SPS allocation and the dynamic-scheduling based allocation is considered.

- The DRB(s) for the SF(s) of the given NSI are configured so as to differentiate the SPS and dynamic resource allocations for the given NSI vs. other NSI(s). For examples, when the SPS allocation related to the given network slice is activated for the UE via a downlink scheduling control channel (similar to PDCCH in LTE for instance), the data of the DRB(s) configured for the given NSI has the higher priority in using the SPS allocation. If the allocated SPS resources are not fully filled by the NSI specific DRB(s) of the given NSI, other DRBs that are configured to support multiplexing of SFs from other network slices or one of the other NSIs can use the allocated SPS resources.
- In case the dynamic resource allocation is also scheduled at the same TTI the SPS resource is allocated, the dynamic resource allocation is used for the UE with the following coordination with the SPS allocation, instead of overriding that as currently specified in LTE. The serving RAN now controls the UE via, e.g., an indication whether and to what extent the dynamic allocation may override the SPS allocation in the same TTI. In one option, MCS allocated in the dynamic scheduling may override that in the SPS allocation in the current and following allocated TTIs of the SPS allocation. In another option, the dynamic scheduling does not override the SPS allocation, neither MCS in the following TTIs of the allocated SPS. In yet another option, the dynamic scheduling may not include MCS setting, implying that the MCS indicated in the SPS allocation will be used for both the dynamic scheduled resources and the SPS resource.

Then, if the dynamic scheduling does not override the SPS and only one transport block (TB) can be transmitted per TTI for the given UE, the transport block size (TBS) shall be calculated based on the aggregated physical resource blocks (PRBs) allocated in both the dynamic scheduling and the SPS and the common MCS. The amount of data corresponding to the SPS allocation should be filled from the DRBs of the given NSI and service thereof related to the SPS allocation and the rest are filled from the DRBs for the given NSI or other slices based on e.g. DRBs' priority, as illustrated in the following Figure 17-6.



Figure 17-6: TB formed based on the coordination between SPS and dynamic scheduling

In case the dynamic scheduling does not override the SPS allocation and multiple TBs can be transmitted per TTI for the given UE, multiple TBs may be formed of which one is for the SPS

associated DRBs based on the SPS allocation and another one is for other DRBs based on the dynamic scheduling allocation.

Figure 17-7 illustrates high level network controlled procedures for the coordinated dynamic scheduling and SPS allocations. There are two options for configuring network slice related SPS shown in Figure 17-7, but these 2 options may not be exclusive. In option 1, UE may be configured with NSI specific SPS for one slice and UE specific SPS for another slice. In option 2, UE may be even configured with both UE specific SPS configuration and, in addition, NSI specific SPS for the same slice for different bearer services. For option 1, the network slice specific SPS configuration may be signalled via common control signalling and, by default, associated to every UE accessing to the network slice. For option 2, the association to the NSI specific SPS configuration may be indicated to the UE via dedicated signalling.



Figure 17-7: Coordinated dynamic scheduling and SPS resource allocations

By incorporating and controlling the possibility of non-overriding the SPS allocation by the dynamic scheduling based allocation, the UE is allowed to continue with the ongoing SPS allocation either with the existing MCS or updated MCS indicated in the dynamic scheduling allocation. This, on one hand, saves the signalling and processing overhead due to the SPS. On the other hand, this provides the continuous guarantee as well as separation of the resource usage between different network slices in RAN sharing scenarios. Furthermore, as the "overriding" and "non-overriding" options can be controlled in an optimized adaptive fashion as described above, both the SPS and dynamic scheduling allocations can be fully utilized for all the involved slices. For examples, unused portion of the SPS allocation by the given NSI may be made usable for other slices at any scheduled TTI. The given NSI may of course size up a portion of the dynamic scheduling allocation.

Remarks:

This contribution may be applied for many industrial use cases in which machine devices with programmable and high resilient machine operation may prefer SPS allocation. However, at least a portion of those machine devices, master or special purpose ones, may need to have access to different network slices for on-demand communications and thus dynamic scheduling based allocation is preferable to serve such the traffic dynamically.

The scheme for coordinating the SPS and dynamic scheduling based allocations in this contribution has certain implications on RRC in c-layer and MAC/PHY in d-layer of the serving RAN being shared by different network slices. The scheme can be incorporated in 5G RAN systems being standardized in 3GPP. Thus, the scheme may have potential for standards contributions to 3GPP for examples.

17.3 References

- [23.799] 3GPP TR 23.799, "Study on Architecture for Next Generation System", December 2016.
- [36.300] 3GPP TS 36.300, "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2 (Release 14)", June 2017.

18 User-centric connection area

18.1 Motivation and problem statement

LTE networks were mainly designed for the transmission of broadband packet payloads, i.e., a large amount of data is sent over the air interface, followed by a long time without activity. Once a UE is attached and registered within the LTE network, its connection to CN is managed with the help of the RRC protocol, which is terminated in the eNB. With respect to the radio activity of the UE, there are two RRC states for the UE namely RRC-Connected and RRC-Idle. The UE remains in the RRC-Connected state during the active data flow in UL or DL. The location of the UE is known on cell level.

In LTE, the connection management is controlled with the help of a RRC timer, which detects the radio in-activity. LTE allows a wide range of timer values. The study in [HQG+12] suggests timer values around 10 s based on the measurements in a practical LTE network.

Today's smartphone applications are developed such that they maintain their connection with the server using keep-alive messages. Additionally, there are also notifications from the server side. This is generally known as background traffic [ZZG+13]. Our motivation for signalling minimisation for 5G stems from this type of traffic. The background traffic mostly carries the packet payload in the order of a few bytes. However, due to the connection oriented nature of the LTE network, the setting up of the connection involves more than 20 signalling messages to transmit few IP packets.

Consider a scenario where a user is engaged in a "WhatsApp" based messaging (with packet interarrival time of around 10 s [36.822]) and the timer is set to a value, which is below 10 secs, then there will be a frequent transition between the RRC-Connected and the RRC-Idle states. Following the expiry of the RRC timer, the network will switch the UE to RRC-Idle state. Since the user is still engaged in a WhatsApp session, it may generate an uplink packet right after it has been switched to the RRC-Idle state. For this purpose, a completely new random access procedure has to be carried out and a new radio bearer has to be setup.

On the other hand, the network may have a downlink packet for the user right after it has been switched to the RRC-Idle state. In this case, a paging procedure followed by a random access procedure is required. Typically, the paging is carried out within the complete tracking area where many base stations are involved. This problem has been identified as "signalling storm", see [Won12]. Hence, we expect that in such cases, the given solution with RRC timer is not suitable enough to reduce the huge amount of signalling messages for small packet payloads on the air interface and towards the MME.

18.2 Major results

A dedicated simulation study has been performed to show the gains of the UCA concept compared to LTE. The saving of signalling messages towards the MME and for paging inside the RAN has been evaluated and substantial gains have been shown.

Based on contributions submitted to 3GPP RAN2 and RAN3 by 5G NORMA partners the proposed main proposals of the UCA concept have been adapted by 3GPP for the New Radio Access Technology [36.804].

Within [ABA+16] the UCA concept and results are published.

18.3 Related work

The UCA concept was adapted by 3GPP for the New Radio Access Technology [36.822]. The main concepts of the UCA will be standardized by RAN2 in the upcoming work item phase.

18.4 UCA concept

A UCA is defined as a coverage area where the UE-context is known in advance to all 5G access points. For each UE with small or sporadic data, an individual UCA is configured by the RAN, e.g. an anchor node, based on SON functionality (neighbour relation table) and with some assistance of the UE (measurement of surrounding access points). A UE specific UCA is depicted in Figure 18-1.

A UE, even if no data has to be transmitted, will remain in the RRC_CONNECTED state but will be transferred to a sub-state defined as UCA_enabled. Within this sub-state, the UE will not report any channel measurements (CQI) or other measurements, i.e. it will perform cell reselection without control of the access network (forward-handover).

With the help of open loop synchronization and efficient access protocols for 5G as described in [SWS15], the UE is able to perform both contention based and contention free UL transmissions in any cell (measured as best server by the UE) within the UCA. This UL transmission may carry the notification of the best server or the small uplink user data.

The best serving node will forward these UL packets to the anchor node. In case of DL transmission, the anchor node (e.g. with the help of best server updates) will forward the packets to the currently best serving node of the user. If the serving node is not known, the anchor node will trigger a paging message within the area of the UCA to identify the best serving node of the UE.



Figure 18-1: (a) Assignment and (b) update of UCA for a UE traversing the path shown by the dashed line

If the UE leaves the coverage of the UCA, a dedicated signalling (compared to the LTE handover procedure) will define a new user specific UCA, see Figure 18-1. The UCA framework minimizes

the radio and core network signalling overhead related to connection management (idle/active transitions) and to mobility (paging, handover).

With the introduction of the new RRC state "UCA_Enabled" the signalling load towards the core network is drastically reduced compared to state of the art in LTE, especially for short data packet (SDP). The transmission of SDP generates several cycles from RRC_CONNECTED to RRC_IDLE and back to RRC_CONNECTED transitions which are avoided with the new UCA state, see Figure 18-2.



Figure 18-2: Saving of signalling messages towards the core with UCA compared to 4G LTE

There are four main events in 4G-LTE that contribute to the signalling load at the MME (in terms of messages entering and leaving MME):

- UE-originated idle-active transition: service request triggered by UL DATA results in 6 messages;
- S-GW-originated idle-active transition: service request triggered by DL DATA results in (8+Paging) messages;
- eNB-terminated active-idle: eNB-initiated S1 release due to UE inactivity results in 5 messages; and
- intra-MME X2 handover: results in 4 messages.

For the evaluation, it was assumed that a UCA update event requires CN (MME) related signalling as the procedure is similar to the intra-MME LTE handover. Thus, the same amount of signalling for UCA update was assumed.

18.5 Simulation framework

A simulation framework has been setup to compare the saving of signalling messages based on the UCA concept against the LTE signalling messages for SDP transmissions based on the following modelling:

• Background traffic modelling

The arrival of SDPs is modelled as negative exponential distribution with mean inter arrival time (IAT) τ [36.814]. The size of the packet is fixed. Every packet is successfully transmitted in the same transmission time. The type of packet (whether UL or DL packet) is controlled with the help of a parameter μ which is the UL packet probability.

• RRC timer and handover modelling for 4G-LTE

After the successful transmission of SDP, the RRC timer is started. The UE remains in RRC_Connected state until the timer expires. The timer is reset upon the arrival of new packets or A3 measurement reports [36.331] in active state as depicted in Figure 18-3.



Figure 18-3: Idle/active modelling for simulation study based on SDP arrivals and UE mobility

• Definition of number or radio cells belonging to a UCA

Within the simulation environment the definition of cells belonging to a UCA is based on UE measurements, i.e. the Reference Signal Received Power (RSRP). The anchor node defines a UCA for the UE by applying a window size in dB, this is depicted in Figure 18-4. There are also other means to define a UCA, e.g. based on distance, known neighbours, mobility anticipation etc. which should be applied in commercial networks.



Figure 18-4: Definition of a UCA based on applying a window size to RSRP based UE measurements

• Simulation assumptions The simulation assumptions and parameters are summarized in Table 18-1.

Parameters	Values
Cellular layout	Hexagonal grid, 91 sites, wraparound
Inter-site distance (ISD)	500 m
Speed	3, 30, 120 km/h
Simulation duration	1 hr
Subframe duration	1 ms
Traffic type	UL+DL
Handover hysteresis	1 dB
Channel model	Urban Macro [36.814]
Mobility	Random walk
Tracking area	91 sites

Table 18-1: Simulation parameter for UCA performance assessment

The simulation environment comprises an urban cellular area with warp-around hexagon network structure. The area is covered by 91 sites equipped with 3-sectorized antennas where each sector has a unique cell ID. A mobile user moves in the simulated area with a given speed. We have used random walk mobility model with different mean path lengths. The simulation results shown here are mainly with linear mobility as they have highest impact on UCA update. Path loss, shadowing and antenna models are used as given in [36.814] Fast fading is averaged out in RSRP and handover measurements. Event triggered handover measurement report is used based on A3 event [36.331]. It is assumed that the UE is attached to the same CN/MME/S-GW and moving inside the same TA consists of 91 sites controlled by the centralized single MME. The gain of the UCA framework as compared to LTE is calculated as follows:

$$\Omega = 100 * (1 - \frac{\alpha}{\beta})$$

where, Ω is the percentage gain over LTE, α is the MME signalling due to UCA updates and β is the total MME signalling in LTE due to the events: UE originated Idle-Active transition, S-GW originated Idle-Active transition, eNB terminated active-idle and intra-MME X2 handover.

18.6 Simulation results

UCA size and UCA updates

As expected, the mean UCA size increases with the increase of the window size, see Figure 18-5. There is almost no impact of the mobility model. Three different mobility types were modelled with the help of random walk with a given Mean Path Length (MPL). Linear stands for a straight line walk with random dropping on border events. Figure 18-5 also shows that for the increase of the window size the number of UCA updates decreases; there is even a large impact on the update rate with random work compared to linear mobility.



Figure 18-5: UCA size and UCA updates with respect to UCA window size

Impact of speed

Figure 18-6 illustrates the gain in signalling reduction achieved with the UCA approach for different speeds with linear mobility and the following simulation parameters: RRC timer = 10 s, mean inter arrival time $\tau = 10$ s and UL packet probability $\mu = 0.5$. With lower mobility, there is almost a 100 % gain implying that the signalling generated with the UCA framework is negligible as compared to LTE regardless of the UCA size. Moreover, signalling related to state transitions in 4G-LTE dominates the mobility related signalling. However, with 120 km/h, as there are considerable updates in UCA, the gains decrease to 80 % for lower UCA sizes.



Figure 18-6: Gain over LTE for different UE speeds with linear mobility

There are even higher gains for none linear mobility schemes, as the probability that the UE remains inside the UCA, resulting in a lower UCA update rate. This is depicted in Figure 18-7, with RRC timer = 10 s, $\tau = 10$ s, $\mu = 0.5$ as simulation parameter.



Figure 18-7: Gain over LTE for different UE speeds and different mobility schemes

Impact of inter-arrival time (IAT)

As long as a UE stays within the UCA_Enabled RRC mode, UCA related signalling is independent of packet arrival rate whereas 4G-LTE signalling depends on it. Figure 18-8 illustrates the gains with different IAT for the speed of 30 km/h, RRC timer = 10 s, $\tau = 10$ s and $\mu = 0.5$.

There are two effects associated with IAT, the probability of packets arrival within RRC timer and the total number of packets. The RRC timer value here is 10s. Initially, with IAT increasing from $\tau = 5$ s to $\tau = 10$ s, the probability of packets arrival after RRC timer and triggering state transition increases and therefore the gain increases. Later, as we further increase the IAT (from 10s to 60s), although the probability that a packet arrives after RRC timer increases, the total number of arrived packets decreases, which leads to a decrease of the gains achieved by our method. Nevertheless, the UCA concept enables high gains even with a very high arrival time.



Figure 18-8: Gain over LTE for different packet inter-arrival times

Impact of RRC timer value

In LTE the RRC timer defines the time period a UE will remain in RRC_ACTIVE after data has been sent in uplink or received in downlink. Consequently, an increase of the RRC timer for a given mean arrival rate will decrease the active/inactive transitions and the gain of the UCA concept over LTE will decrease. The impact of RRC timer on higher speed has even more impact as there is a higher number of handovers, which also prevent the change towards RRC_IDLE and which further reduce the gain over LTE. This is shown in Figure 18-9 for the 3 speed values of 3 km/h, 30 km/h and 120 km/h, a linear mobility and $\mu = 0.5$. However, large RRC timer values increase the air-interface load due to uplink-feedback reporting of the UE. This is not covered by the simulation study.



Figure 18-9: Gain over LTE for different RRC timer values and different UE speeds

Impact of paging area

The preceding simulation results have considered a full tracking area as paging area, i.e. 90 sides with 270 radio cells, which causes a lot of paging related signalling due to the arrival of downlink packet in the RRC_Idle state. 3GPP has already addressed for 4G-LTE a paging optimization mechanisms which is implemented in commercial networks, e.g. a paging is triggered within a selected area taking into account the last known serving eNB. If this paging is not successful, the area will be increase until finally a paging within the complete tracking area will be initiated. Figure 18-10 compares the signalling with different paging area sizes with the simulation setup for an UL packet probability of $\mu = 0.5$, RRC timer = 10 s, linear mobility and with $\tau = 10$ s.



Figure 18-10: Gain over LTE for different paging are sizes

For high speed the gains are less compared to low and medium speed as in this case the mobility related signalling dominates the state transitions and paging events are reduced. Nevertheless, the UCA concept enables considerable gains, even if the LTE paging optimization, e.g. only a limited number of eNBs are triggered by the MME for paging, is applied.

The simulation study has shown, that the UCA concept can be realized with sustainable gains with the features of 5G air-interface and efficient access protocols for SDP transmissions. Compared to a traditional cell-centric architecture in 4G-LTE, this proposal allows a flexible handling of data traffic and associated signalling in an efficient manner.

18.7 Security considerations

The PDCP protocol is terminated within the anchor node (gNB), i.e. there is no need to provide key to other gNBs inside the UCA cluster. Consequently, there is no security problem with PDCP.

The same applies to the RRC protocol, i.e. also no security problem.

The UE identifier will be used as long the UE remains in the UCA, which might be quite long, depending on the UE mobility. To make tracking of UEs via this identifier harder, during the final standardization of the UCA concept a RRC controlled change of the UE identifier, e.g. timer based, should be investigated.

As the UE supports the definition of its UCA area by measurements of neighbour cell, there might be a threat potential if this information is counterfeited. But in this case it will injure itself as a none optimal UCA will be allocated based on the wrong measurement values.

18.8 References

- [36.822] 3GPP TR 36.822, "LTE Radio Access Network (RAN) enhancements for diverse data applications (Release 11)", September 2012.
- [36.331] 3GPP TS 36.331, "EUTRA Radio Resource Control (RRC) protocol specification (Release 12)", December 2013.
- [36.814] 3GPP TR 36.814, "Further advancements for E-UTRA physical layer aspects (Release 9)", March 2010.
- [38.804] 3GPP TR 38.804, Study on New Radio Access Technology; Radio Interface Protocol Aspects (Release 14), March 2017.

[ABA+16]	D. Aziz, H. Bakker, A. Ambrosy, and Q. Liao, "Signalling minimisation
	framework for short data packet transmission in 5G", IEEE Vehicular
	Technology Conference, Montreal, Canada, September 2016.
[SWS15]	S. Saur, A. Weber, and G. Schreiber, "Radio access protocols and preamble
	design for machine-type communications in 5G," in Signals, Systems and
	Computers, 2015 49th Asilomar Conference on, pp. 3–7, November 2015.
[HQG+12]	J. Huang, F. Qian, A. Gerber, Z. M. Mao, S. Sen, and O. Spatscheck, "A close
	examination of performance and power characteristics of 4G LTE networks", in
	Proceedings of the 10th International Conference on Mobile Systems,
	Applications, and Services, ser. MobiSys '12. New York, NY, USA, ACM, pp.
	225–238, 2012. Available at: http://doi.acm.org/10.1145/2307636.2307658
[ZZG+13]	S. Zhang, Z. Zhao, H. Guan, D. Miao, and H. Yang, "Statistics of RRC state
	transition caused by the background traffic in LTE networks", in 2013 IEEE
	Wireless Communications and Networking Conference (WCNC), pp. 912–916,

 [Won12] April 2013.
[Won12] Dan Wonak, "A Storm is Brewing – an LTE Signaling Storm", White Paper, Diametriq LLC, September 2012. Available at: http://www.diametriq.com/wpcontent/uploads/downloads/2013/01/A-Storm-is-Brewing.pdf

19 Massive machine-type communication RAN congestion control

Many different solutions have been designed to control RAN congestions for mMTC. In recent years, a batch of new methods are proposed, with the common idea of clustering MTCDs into groups, selecting one group coordinator (GC) for each group, and let them relay data for the other group members (GMs) with help of D2D links. This kind of method has the advantages of ultrahigh energy efficiency, low access time for duty-cycle-critical applications and simultaneously an efficient reduction in RA collision. However, all these methods are designed based on the assumption that D2D links are available and reliable, which is not always practically feasible. Study on the impact of D2D unreliability and corresponding solutions are required.

19.1 Impact analysis of D2D link exception

It is clear, that increasing the group size leads to a decrease in RA request density, and hence reduces the RAN congestion rate. However, it should be considered that the intra-group D2D links are not always available and reliable. The reliability even drops when the group size grows. And the link failures will lead to extra signalling for reports and exception handling, this will compromise the gain if these messages are transmitted through extra RA processes. Assuming an exponentially decreasing average D2D link lifetime w.r.t. the group size, unreliable D2D links can even lead to a loss in case of large groups, as Figure 19-1 shows.



Figure 19-1: The simulated impact of D2D link failure reports through extra RA processes

19.2 Enhanced grouping processes

First, three different processes are designed in the protocol:

- global group update,
- group joining,
- group leaving.

GGU tries to update all groups in the local cell. Group joining will lead to a similar process, but involves only one group, which the joining MTCD will be added to. The message flow is shown in Figure 19-2.



Figure 19-2: The global group updating process

Depending on the process trigger, the group leaving can be executed in different methods, as shown in Figure 19-3 to Figure 19-6.



Figure 19-3: The group leaving process, triggered by detaching events







Figure 19-5: The group leaving process, triggered by a GM losing D2D and macro-cell links


(b) A GC moves towards another cell.

Figure 19-6: The group leaving process, triggered by handover events

19.3 Enhanced transmission frame structure



Figure 19-7: Frame structure of the grouped mMTC transmission based on D2D; the terms R/R, UDT, A/C and DDT denote request/report, uplink data, acknowledgment/command, respectively

A new frame structure of D2D-based grouped mMTC transmission was designed, to embed most reports, acknowledgements, requests and commands into the d-layer, as shown in Figure 19-7.

19.4 Geolocation database integration

Finally, we briefly investigate the potential of the GDB concept assisting mMTC congestion control and GDB. The GDB is applicable to a multiplicity of purposes, and these have been mentioned in [5GN-D41] and later in Section 6.7 of this document. This is related to the mMTC solution using geolocation and advanced context information (e.g., propagation mapping) in order to manage the grouping of MTCDs through centralized management and decision making. The use of a GDB in this context, along with appropriate messaging built into the proposed protocol, serves a number of purposes and benefits. Particularly, it facilitates:

- The optimal and more efficient grouping of MTCDs based on the channel conditions between them, taking into account locations, propagation, the requirements in terms of communication characteristics and rate, among others. It also achieves a system-level viewpoint considering, e.g., interactions among the transmissions of different groups.
- The collection of information from the MTCDs to enhance the performance of the GDB in serving the above purpose, as well as others.

A generic GDB signalling example closely derived from a regulatory TVWS GDB is as in Figure 19-8. The procedure is as follows. First, a request from the element at the top level in the hierarchy is made to the GDB, where it is assumed that that element already has a means of connecting back to that GDB (e.g., an Internet connection). This top-level element is generally known as the MASTER device in GDB-related spectrum sharing contexts, and is the GC in the context here. The request from the MASTER/GC can be for a range of purposes; but let's say, typically, it is a resource request in spectrum sharing scenarios, and will include context information about the MASTER/GC such as its location, requirements (e.g., QoS), and technical characteristics. The response from the GDB will convey the allocated resource to the MASTER/GC, or alternatively in other usage scenarios will convey management instructions derived from geolocation, as well as other information such as the technical characteristics and requirements of the GC. The devices at the next level in the hierarchy are known as SLAVE devices, or GMs in the context here. In a generalized case, it is assumed that the SLAVEs (GMs) might not already have preexisting connectivity back to the MASTER (GC)-noting that in some D2D contexts applicable to mMTC such an assumption can be valid. The exchange in Figure 19-8, where the MASTER makes a generic (worst case) query on behalf of the SLAVE, solves this issue. Upon receiving the generic (worst case) management instructions or resources via the MASTER as calculated by the GDB, the SLAVE is then operational and connected to the MASTER, and can then make a specific request again via the MASTER in order to get more precise parameters. The specific request will lead to far better conveyed or allocated resources, or otherwise a more optimal configuration from management for the intended purpose. In some cases, where the GC does not have an initial Internet/network connection, the BS might be seen as a MASTER, the GC a slave, and GMs might be seen as a third level in the hierarchy—also depicted in Figure 19-8.

GDB MAC-layer messaging supporting the mMTC protocol

GDBs are generally assumed to interact with users through messages exchanged at the network/transport layers. However, in scenarios where GDB functions are closely integrated within particularly systems, messaging might take place at the MAC layer. This is indeed the case for the use of GDBs to support grouped mMTC in the context of this paper. Referring back to Figure 19-7, and comparing with Figure 19-8, the DA timeslot represents the communication of the SLAVE (GM) information with the MASTER (GC), which then integrates all messages into one communication to the GDB in the AUT timeslot. This reflects GDB contexts such as TVWS, whereby the MASTER communicates messages on behalf of the SLAVE. The ADT timeslot then reflects the combined response from the GDB via the MASTER (GC), subsets of which are then forwarded by the MASTER to the respective SLAVEs (GMs) in the DD timeslot. Again

considering Figure 19-7, it is noted that the R/R and A/C messages will be updated to support the aforementioned GDB-related signalling.



Figure 19-8: GDB signalling (generic case)

19.5 Security considerations

The proposed D2D-based grouped RA approaches lead to a security risk, that not only the control layer data, but also the data layer messages are relayed by the GCs. If an unreliable device are chosen as GC, a leakage may be caused. Hence, when applied in complex scenarios where unknown devices may join the local network, the GC selection algorithm must take the trust value into account. And the UP messages should be encrypted before assembled at the GC.

Furthermore, when considering the information exchange between the mMTC congestion control module and the GDB, it shall be noted that privacy information such as UE locations, UE identities and DC of the devices shall not be shared to the third-party GDB. Instead of forwarding

raw channel measurements to the GDB, the mMTC congestion control module should generate and submit an abstract report about the D2D link availability in the local area, which does not contain privacy information of UEs.

19.6 Summary

It has been derived that D2D link exceptions caused by device mobility and channel fading can seriously damage the performance of D2D-based grouped RA methods, if the exception handling generates an extra RA process.

An enhanced procedure has been proposed, taking the D2D link exceptions into account, which includes three different processes (global group updating, group joining and group leaving) and a novel group transmission frame structure.

An optimal resource allocation approach has been developed to minimize the overall congestion rate, while providing a capability of guaranteeing specialized performance to particular device groups with preference.

A concept has been proposed to connect the mMTC group management and the geolocation database (GDB), which can benefit them both.

20 Geolocation databases, use of geolocation information and associated opportunities

This section describes the innovation related to the use of a Geolocation Database (GDB) functionality within the 5G NORMA architecture.

20.1 Motivation and problem statement

Geolocation and associated management capabilities, administered through a GDB or hierarchy of GDBs operating in government/regulatory and operator contexts, will be highly relevant in 5G communication contexts. This is for a number of reasons linked to: (i) spectrum sharing for 5G and various forms of spectrum databases/GDBs being the solution to administer that sharing, (ii) heterogeneity of 5G communications and issues such as rendezvous (based on location) linked to that, and (iii) the need to consider geolocation *per-se* in 5G communication systems, e.g., in order to minimise the propagation path and achieve a true optimization of network provisioning with such things in mind.

One particular case where geolocation information per-se is useful in 5G communication contexts is the realisation of a virtualised optimised propagation path for the Tactile Internet as a key 5G and beyond application, supported by signalling over the mobile network for control purposes. Here, as well as GDBs offering optimal location of virtualised elements, the GDBs themselves may often need to be virtualised and optimally located in order to minimise latency in signalling with the GDB to make/obtain geolocation-based decisions. An example of surgeon using the Tactile Internet based on our scheme/proposal to do remote operation is given in Figure 20-1, whereby the lower-most part (more direct path) is achieved by the virtualised network elements each consulting with a GDB. More detail on such aspects is provided in Part I Section 6.7 of this deliverable.



Figure 20-1: Virtualized edge network slices achieving a more direct path compared with (fixed) network elements in a Tactile Internet remote surgical operation example

In addition to the GDB applicability to a multiplicity of purposes covered in Section X of this Deliverable, particular work has been done in the context of the application to Massive Machine-Type Communication (mMTC), and the gearing of the associated concepts to the protocol for mMTC indicated in Section X. Key applications include the mMTC solution using geolocation and advanced context information (e.g., propagation mapping) to manage the grouping of MTCDs through centralized management and decision making. Benefits of this include the optimal and more efficient grouping of MTCDs based on a number of factors, and better consideration of a system-level viewpoint, among others. More detail is provided in Section X of this Deliverable.

20.2 Major results

Figure 20-2, Figure 20-3, Table 20-1 and Table 20-2 demonstrate the TV spectrum (470–790 MHz) that can be realised through the use of such a GDB concept in a TV White Space (TVWS) context, in terms of the number of channels available, for a pessimistic case for Southern England (availability is less there than in, e.g., the North/Scotland). Here, the transmitter is 30m above ground level, with \geq 30 dBm allowed EIRP being necessary in order for the given channel to be counted as viable. This demonstrates that the use of such a scheme could realise some significant additional spectrum close to (or, until the spectrum is auctioned in the UK) the vital 700 MHz low-band, assisting aspects such as coverage/reliability and control of 5G systems.



Figure 20-2: Map of channel availability for a Class 3 white space device, for a large area of England, for a 30m height above ground level transmitter, ≥30 dBm allowed power



Figure 20-3: CCDFs of channel availability for different classes of white space device, for a large area of England, for a 30 m height above ground level transmitter,≥30 dBm allowed power scenario

Scenario	Database calc.	Location	Ave. No. chan.	Std. no. chan.	CoV no. chan.	% loc. ≥ 1 chan.	% loc. ≥ 3 chan.
MBD	past	wide area	8.6	7.2	0.83	98.0	81.6
		London M25	15.2	8.5	0.56	99.5	97.1
	present	wide area	4.0	4.9	1.21	74.1	43.1
		London M25	4.7	3.5	0.76	90.8	63.1
	future (WRC)	London M25	1.6	1.4	0.88	82.5	21.0
IBP/ MBU	past	wide area	26.5	6.1	0.23	100.0	100.0
		London M25	25.5	3.6	0.14	100.0	99.9
	present	wide area	23.5	7.1	0.30	99.9	99.7
		London M25	24.8	4.7	0.19	100.0	99.9
	future (WRC)	London M25	14.4	3.6	0.25	99.5	99.0

Table 20-1: Statistical results on number of allowed channels

Table 20-2: Scenarios (MBD: Mobile Broadband Downlink; IBP/MBU: Internal Broadband Provisioning/Mobile Broadband Uplink)

Scenario	Transmitter height (m)	Required EIRP (dBm)
MBD	30	at least 30
IBP/MBU	1	at least 20

Figure 20-4 makes clear the saving that can be achieved by selection of computational resources to optimize the communication path, using GDB functionality. This is for a 100 km*100 km simulation area, with 10 nodes present, based on a Monte-Carlo simulation approach. The non-optimal case is for communication with a fixed "anchor" what represents a next element up in the communication hierarchy or a vital network element, whereby node positions (start point and end point of the communication link) and anchor positions are randomly chosen in each simulation iteration. The optimal case is where all nodes are still randomly-placed, but the choice of best node for the anchor is based on the communication path is chosen.

Results here show that the average latency for the non-optimised case is 0.23 ms; and for the optimised case is 0.16 ms. This represents a 30 % saving in latency.



Figure 20-4: Example of latency CDFs for optimised (using the GDB) and non-optimised Network Function placement.

20.3 Related work

The geolocation database concept has been a key tenet of spectrum sharing regimes since the FCC announced their use for TV white space in 2008, and drastically increased their importance in 2010. Other countries have followed the US approach, whereas the UK then South Africa have chosen a more complex approach that allows variable white space device transmission power and 5 classes of devices.

Much work has been undertaken in this context. In the US, Harrison *et al.* have analysed white space availability [HMS10], [MDKHS15]. In Europe, van der Beek and Fitch have attempted to predict TVWS white space availability [BRAM12], [FNKBM11], and Holland *et al.* has measured availability and capacity using real TVWS GDBs, in both the trial and commercially operational phases [Hol15], [Hol16a], [Hol16b]. Such latter work has led to some fundamental conclusions on the design and deployment of white space devices.

Regarding the signalling procedures argued for the GDB in Section 20.4, it is noted that, in the context that the cognitive pilot channel (CPC) must be backed by context awareness and in some contexts decision-making processes linked to that, these signalling approaches can in some cases be seen as having parallels to the CPC concept—particularly in the case of rendezvous covered in Figure 20-6. It is noted, however, that this is particularly in the case of the so-called on-demand

version of the CPC [PSAG07], which defines the resource sharing on a signal rendezvous channel (i.e., the CPC) for that context.

20.4 Signalling procedures

The concept of GDBs has been proven in the context where regulatory approval is needed for decisions on spectrum sharing. This is particularly in the case of TVWS, where sharing is between services even without agreement of incumbent in a given band.

The signalling procedures for such contexts is as follows:

- 1. The device sends its characteristics and resource request (perhaps with other context information, e.g., its traffic requirements) to the database.
- 2. The database responds with information on available resources, given the above information from the device.
- 3. The device chooses the resources it will use, and sends information on its chosen resources and reconfirms its ID to the database along with this message.
- 4. Database logs this information, and sends a response to confirm/allow the choice of resources, at which point the devices can begin transmitting.
- 5. Periodic updates or reconfirmations are required by the device/database. In the TVWS context, the resources have to be updated or reconfirmed every 15 minutes.

It is noted that there is also the option of broadcast by the Master (i.e., the directly databaseconnected element), as a form of "pilot", device generic (worst-case) parameters anywhere within the coverage area of the Master to allow, Slave devices (e.g., the end-user devices, who don't have any pre-existing means of Internet connection hence no way of access the GDB), to allow them to connect back to Master and send their parameters to the Master, to be authorised through the Master forwarding the message to the GDB, based on the Slave's specific (i.e., non-worst case) parameters.

It is suggested that we (broadly) maintain support for such procedures to facilitate the realisation of cases where the regulator needs to be (directly/indirectly) involved in the process. In line with such objectives, this GDB approach can be generalised to any case of multi-level resource request/allocation; notably (in the context of 5G virtualisation) computational resources to realise the virtualised network functions. The generic procedure, as illustrated in Figure 20-5, might apply as follows:

- 1. The UE assesses the situation (e.g., traffic requirements) and sends request based on that, in conjunction with its geolocation information.
- 2. Other network elements/functions higher in the chain combine these requests, and may add to them, e.g., with information on their own specific constraints/requirements/observations, etc.
- 3. The GDB calculates allowed resources, or even manages resources in some scenarios.
- 4. The GDB may forward information to SDM-O which may play a part in managing those resources which have to be dealt with within the scope of the SDM (e.g., computational resources, handled based on the geolocation of their availabilities and packaged/forwarded by the GDB).
- 5. Resource responses are sent back to the UE and perhaps in some scenarios other network elements/functions.
- 6. The implementation of the decision by the UE (or in some cases other network elements) is ACKnowledged/confirmed with the GDB, and in some cases the SDM; it may be the case that the UE implements its own decision within the constraints of the response from the GDB/SDM, and feeds back information on its decision in the "ACK". For example, it may choose a lower transmission power than the upper allowed limit returned in the response from the GDB/SDM.



Figure 20-5: Signalling—resource management (generalisation)

This GDB signalling can be adapted to other purposes. Figure 20-6 shows how it might apply to rendezvous in a heterogeneous networking scenario, reducing the need for, e.g., the devices to keep all radios active to detect connection opportunities among the heterogeneity of possible links, thereby saving energy and enhancing battery life. Here, the information/requests can include location information, and detail on supported radio access technologies and associated configurations (e.g., supported bands), and the response indicates possible connection opportunities—while taking into account considerations such as privacy of the devices and elements that might provide those opportunities—perhaps in a multi-hope fashion.



Figure 20-6: Signalling—rendezvous

Figure 20-7 shows the option of Master devices/elements carrying initial information on behalf of unconnected Slave devices/elements, discussed above. This is entirely compliant with how the process is managed in the scope of TVWS, for example.



Figure 20-7: Signalling—indirect access

Finally, through the replication of the MAC for mMTC in Figure 20-8, and the reflection in the context of Master-Slave signalling procedures in Figure 20-9—although in this case applied at the MAC layer, these Figures illustrate the mapping to the proposed mMTC protocol detailed in Section X of this Deliverable on to the GDB concept.



Figure 20-8: Signalling—Link to mMTC congestion control, MAC structuring



Figure 20-9: Signalling—Link to mMTC congestion control reflected in GDB signalling (partial)

20.5 Security considerations

A number of security considerations apply in the context of the use of the GDB and associated signalling mechanisms to access the GDB. First, the signalling itself has to be secure against eavesdropping, and particularly malicious modification through man-in-the-middle attacks. This can be realised through setting up secure sockets between the GDB and it clients. Second, the authenticity of GDBs and their users has to be ensured. In the context of TVWS and other similar mechanisms, this is through the definition of a server providing information on trusted GDBs, as well as registration of trusted users of the GDBs with those GDBs. This of course has to be twinned with a secure communication channel also with server for trusted GDBs.

There are, however, significant research questions that need to be addressed concerning how these security aspects are realised in the context of GDBs applied to 5G topics highlighted in this section and Section 6.7.

20.6 References

- [HMS10] K. Harrison, S. M. Mishra, and A. Sahai, "How Much White-Space Capacity Is There?", IEEE DySPAN 2010, Singapore, April 2010.
 [MDKHS15] V. Muthukumar, A. Daruna, V. Kamble, K. Harrison, and A. Sahai, "Whitespaces
- after the USA's TV incentive auction: A spectrum reallocation case study," IEEE ICC 2015, London, UK, June 2015.

[BRAM12]	J. van de Beek, J. Riihijarvi, A. Achtzehn, and P. Mahonen, "TV white space in Europe", IEEE Trans. Mob. Comput., Vol. 11, No. 2, February 2012, pp. 178-188.		
[FNKBM11]	M. Fitch, M. Nekovee, S. Kawade, Keith Briggs, and R. MacKenzie, "Wireless service provision in TV white space with cognitive radio technology: A telecom operator's perspective and experience", IEEE Commun. Mag., Vol. 49, No. 3, March 2011.		
[Hol15]	O. Holland <i>et al.</i> , "To White Space or Not To White Space: That it the Trial within the Ofcom TV White Spaces Pilot", IEEE DySPAN 2015, Stockholm, Sweden, September–October 2015.		
[Hol16]	O. Holland, "Some Are Born With White Space, Some Achieve White Space, and Some Have White Space Thrust Upon Them", IEEE Transactions on Cognitive Communications and Networking, Vol. 2, No. 2, June 2016, pp. 178- 193.		
[Hol16b]	O. Holland <i>et al.</i> , "Changing availability of TV white space in the UK", IET Electronics Letters, Vol. 52, No. 15, July 2016, pp. 1349-1351.		
[PSAG07]	J. Perez-Romero, O. Sallent, R. Agusti, and L. Giupponi, "A novel on-demand cognitive pilot channel enabling dynamic spectrum allocation", IEEE DySPAN 2007, Dublin, Ireland, April 2007.		