



H2020 5G-TRANSFORMER Project
Grant No. 761536

5G-TRANSFORMER final system design and Techno-Economic analysis

Abstract

This deliverable reports the final 5G-TRANSFORMER system design and a techno economic study of the implemented platform. The final system design mainly includes some further enhancements to some features described in the refined design reported previously in D1.3. The techno economic study mainly focuses on analysing the 5GT ecosystem and players business interactions with regards to the vertical use case services implemented as proof of concepts from the different verticals involved in the project. The outcome of the study is a set of recommendations to allow more sustainable business models.

Document properties

Document number	D1.4
Document title	5G-TRANSFORMER final system design and Techno-Economic analysis
Document responsible	Thouraya Toukabri (ORANGE)
Document editor	Thouraya Toukabri (ORANGE)
Editorial team	Thouraya Toukabri (ORANGE)
Target dissemination level	Public
Status of the document	Final
Version	1.0

Production properties

Reviewers	Barbara Martini (SSSA), Giovanni Rigazzi (IDG), Carlos J. Bernardos (UC3M)
------------------	----------------------------------------------------------------------------

Disclaimer

This document has been produced in the context of the 5G-Transformer Project. The research leading to these results has received funding from the European Community's H2020 Programme under grant agreement N° H2020-761536.

All information in this document is provided "as is" and no guarantee or warranty is given that the information is fit for any particular purpose. The user thereof uses the information at its sole risk and liability.

For the avoidance of all doubts, the European Commission has no liability in respect of this document, which is merely representing the authors view.

1 Table of Contents

List of Contributors	5
List of Figures	6
List of Tables	8
List of Equations	11
List of Acronyms	15
Executive Summary and key contributions.....	18
1 Introduction.....	19
2 5G-TRANSFORMER Techno-economic analysis	21
2.1 Background	21
2.2 Analytical study.....	22
2.2.1 Vertical use case modelling	23
2.2.2 Infrastructure cost modelling	35
2.2.3 Analytical pricing modelling	57
2.3 Experimental study	73
2.3.1 Methodology.....	73
2.3.2 Simulation system structure.....	74
2.3.3 Actors in the study	74
2.3.4 Tool description	79
2.3.5 Scenario setting.....	86
2.4 Summary	114
2.4.1 Lessons learnt	114
2.4.2 KPIs of the project	115
2.4.3 Business model transformation of 5G-T vertical use cases	118
3 5G-TRANSFORMER Final architecture design and refinements.....	121
3.1 Summary of final architecture design and refinements.....	121
3.2 Radio Abstraction	123
3.2.1 General Concepts.....	123
3.2.2 Radio Abstraction in the 5GT-VS	126
3.2.3 Radio Abstraction in the 5GT-SO.....	127
3.2.4 Radio Abstraction in the 5GT-MTP	128
3.3 MEC integration.....	131
3.3.1 MEC capabilities in 5G-TRANSFORMER descriptors, and VS-level operations.....	131
3.3.2 SO-level operations	133

3.3.3	MTP-level operations	134
3.3.4	Service onboarding	136
3.3.5	End-to-end instantiation workflow	136
4	Conclusions	139
5	References	140
A.1	Infrastructure cost modelling examples	142
A.1.1	Management switch pricing	142
A.1.2	Network Node Pricing	143
A.1.3	Controller node pricing	144
A.1.4	Compute node pricing	145
A.1.5	Block storage HDD pricing	146
A.1.6	Block storage SDD pricing	146
A.1.7	Bill of materials for object storage low density	147
A.1.8	Bill of materials for object storage high density	148
A.1.9	Bill of materials for object storage archiving	148
A.1.10	TCO pricing	148

List of Contributors

Partner Short Name	Contributors
UC3M	Carlos J. Bernardos, Marcelo Bagnulo
NEC	Xi Li, Andres Garcia Saavedra, Josep Xavier Salvat
TEI	Paola Iovanna, Fabio Ubaldi
ATOS	Ignacio Domínguez, Jose E. González, David Salama
NOK-N	Thomas Deiß
TID	Alberto Solano, Luis M. Contreras
ORANGE	Thouraya Toukabri, Philippe Bertin, David Mathieu
BCOM	Farouk Messaoudi, Cao-Thanh Phan
NXW	Giada Landi, Juan Brenes
CTTC	Josep Mangues, Jordi Baranda, Ricardo Martínez, Ramon Casellas
POLITO	Marco Ajmone Marsan, Carla Fabiana Chiasserini
EURECOM	Pantelis Frangoudis, Adlen Ksentini
SSSA	Luca Valcarenghi, Nicola Sambo, and Piero Castoldi
CRF	Aleksandra Stojanovic, Marina Giordanino

List of Figures

Figure 1: Gartner’s hype cycle for cloud computing [10].....	22
Figure 2: Entertainment use case topology	24
Figure 3: Extended entertainment use case topology.....	24
Figure 4: MVNO use case topology.....	26
Figure 5: EVS Automotive UC topology	28
Figure 6: Non-emergency eHealth use case topology	30
Figure 7: Emergency eHealth case topology	31
Figure 8: Topology and components of the EIndustry cloud robotics demonstrator	33
Figure 9: Latency requirements for eIndustry use case	34
Figure 10: eIndustry Service Topology.....	34
Figure 11: Traditional three-tier datacentre topology	36
Figure 12: Standard Spine-and-Leaf Topology [15].....	36
Figure 13: A simple rack with several nodes	38
Figure 14: The TCO for three sizes of cloud size based on the previous examples of hardware nodes	56
Figure 15: Cost per Month VM reservation on Amazon AWS Europe.....	57
Figure 16: Pessimistic scenario	64
Figure 17: Realistic scenario.....	65
Figure 18: Optimistic scenario.....	66
Figure 19: Ideal scenario.....	68
Figure 20: Overflow scenario	69
Figure 21: Followed methodology	73
Figure 22: System Structure	74
Figure 23: Service-graphs.....	76
Figure 24: Deployment event flowchart.....	80
Figure 25: Dynamic cost variation	82
Figure 26: Scale event flowchart.....	83
Figure 27: End event flowchart.....	84
Figure 28: Pessimistic Service ID histogram	87
Figure 29: Pessimistic Service life-time histogram	88
Figure 30: Pessimistic revenue	92
Figure 31: Pessimistic profit.....	92
Figure 32: Realistic Service ID histogram	93
Figure 33: Realistic Service life-time histogram.....	93
Figure 34: Realistic revenue	97
Figure 35: Realistic profit	97
Figure 36: Optimistic Service ID histogram	98
Figure 37: Optimistic Service life-time histogram	98
Figure 38: Optimistic revenue	102
Figure 39: Optimistic profit	102
Figure 40: Ideal Service ID histogram	103
Figure 41: Ideal Service life-time histogram	103
Figure 42: Ideal revenue	107
Figure 43: Ideal profit.....	107
Figure 44: Overflow Service ID histogram.....	108
Figure 45: Overflow Service life-time histogram	108

Figure 46: Overflow blocked services.....	111
Figure 47: Overflow not-deployed services	111
Figure 48: Overflow revenue	113
Figure 49: Overflow profit.....	113
Figure 50: 5G-TRANSFORMER system architecture.....	121
Figure 51: RAN and CN entities involved in the deployment of a vertical service	124
Figure 52: RAN Abstraction approaches	125
Figure 53: Exemplary NSD for a vertical service including RAN	127
Figure 54: 5GT-SO view of the RAN and the coverage area	129
Figure 55: 5GT-MTP view.....	130
Figure 56: Relationship between AppD, SAPD, and NSI with a focus on location constraints	132
Figure 57: Descriptor extensions for MEC support.....	133
Figure 58: Workflow of deploying a network service instance that includes MEC applications.....	138
Figure 59: Cost breakdown for a small private cloud (in %).....	150
Figure 60: Cost breakdown for Small private cloud (in \$)	151
Figure 61: Cost breakdown for medium private cloud (in %)	152
Figure 62: Cost breakdown for medium private cloud (in \$).....	153
Figure 63: Cost breakdown for large private cloud (in %)	154
Figure 64: Cost breakdown for large private cloud (in \$)	155
Figure 65 : Comparison of the cost breakdowns of cloud sizes	155
Figure 66: Evolution of the hardware costs from 2013 to 2019	156
Figure 67: Evolution of the hardware costs from 2013 to 2019	156
Figure 68: Evolution the TCO of a Small DC from 2013 to 2019.....	157
Figure 69: Evolution the TCO of a Medium DC from 2013 to 2019.....	157
Figure 70: Evolution the TCO of a Large DC from 2013 to 2019.....	158

List of Tables

Table 1: Entertainment use case VNFs' resources.....	24
Table 2: Entertainment use case VNFS' Links assets	25
Table 3: MVNO use case VNFs' resources	26
Table 4: MVNO use case VNFS' Links assets.....	27
Table 5: Automotive VNFs' resources	28
Table 6: Automotive VNFS' Links assets	29
Table 7: Non-emergency eHealth VNFs' resources	30
Table 8: Non-emergency eHealth VNFs' Links assets.....	30
Table 9: Emergency eHealth VNFs' resources.....	31
Table 10: Emergency eHealth VNFs' Links assets.....	32
Table 11: eIndustry VNFs' resources	35
Table 12: eIndustry VNFs' Links assets	35
Table 13: Entry variables that influence the costs	38
Table 14: Entries Influencing the spine switch ports price	39
Table 15: Cost Equations for Spine Switches.....	39
Table 16: Entries Influencing the Leaf switch ports price	39
Table 17: Cost Equations for Leaf Switches.....	40
Table 18: Entries Influencing the Management switch ports price	40
Table 19: Cost Equations for Management Switches	40
Table 20: Entries Influencing The Network node Costs per month	41
Table 21: Cost Equations for Network Nodes.....	41
Table 22: Entries influencing the controller node costs per month.....	42
Table 23: How to compute Monthly Costs FOR Controller Node	42
Table 24: Entries Influencing The Compute Node Costs Per Month.....	42
Table 25: How to compute monthly costs for a compute node.....	43
Table 26: Entries influencing the block storage hdd costs per month	43
Table 27: How to compute monthly costs for the block storage hdd	44
Table 28: Entries influencing the block storage ssd costs per month.....	44
Table 29: How To Compute Monthly Costs For The Block Storage SSD.....	45
Table 30: Entries influencing the object storage low density costs per month.....	45
Table 31: How to compute monthly costs for the object storage low density	46
Table 32: Entries influencing the object storage high density costs per month	47
Table 33: How to compute monthly costs for the object storage high density	47
Table 34: Entries influencing the object storage archiving costs per month.....	48
Table 35: How to compute monthly costs for the object storage archiving.....	48
Table 36: Pessimistic scenario.....	64
Table 37: Realistic scenario	65
Table 38: Optimistic scenario	67
Table 39: Total amount of resources needed to run a service	70
Table 40: Maximum number of services by CPU.....	70
Table 41: Maximum number of services by RAM	70
Table 42: Maximum number of services by Disk	70
Table 43: Maximum number of simultaneous services by datacentre type	70
Table 44: Service Life-Time	71
Table 45: Total number of services per datacentre type.....	71
Table 46: Service demand distribution	72

Table 47: Estimated demand for the Pessimistic Scenario	72
Table 48: Estimated demand for the Realistic Scenario	72
Table 49: Estimated demand for the Optimistic Scenario	72
Table 50: Estimated demand for the Ideal Scenario	72
Table 51: Estimated demand for the Federation Scenario	73
Table 52: VNF definitions	76
Table 53: Links definitions	77
Table 54: Datacentres definition according to 2.2.2	79
Table 55: Datacentres distribution across domains	86
Table 56: Total CAPEX+OPEX	86
Table 57: Prices for licenses, transactions and special requirements	87
Table 58: Pessimistic Datacentres state	88
Table 59: Pessimistic interchange bandwidth among domains	90
Table 60: Pessimistic Scenario values	90
Table 61: Maximum local demand pessimistic scenario	91
Table 62: Real local demand pessimistic scenario	91
Table 63: Pessimistic Scenario Prices	91
Table 64: Pessimistic Cost	91
Table 65: Realistic Datacentres state	94
Table 66: Realistic interchange bandwidth among domains	95
Table 67: Realistic Scenario values	96
Table 68: Maximum local demand realistic scenario	96
Table 69: Real local demand realistic scenario	96
Table 70: Realistic Scenario Prices	96
Table 71: Realistic Cost	96
Table 72: Optimistic Datacentres state	99
Table 73: Optimistic interchange bandwidth among domains	100
Table 74: Optimistic Scenario values	101
Table 75: Maximum local demand optimistic scenario	101
Table 76: Real local demand optimistic scenario	101
Table 77: Optimistic Scenario Prices	101
Table 78: Optimistic Cost	101
Table 79: Ideal Datacentres state	104
Table 80: ideal interchange bandwidth among domains	105
Table 81: Ideal Scenario values	106
Table 82: Demand ideal scenario	106
Table 83: Ideal Scenario Prices	106
Table 84: Ideal Cost	106
Table 85: Overflow Datacentres state	109
Table 86: Overflow interchange bandwidth among domains	110
Table 87: Federation Scenario values	112
Table 88: Demand federation scenario	112
Table 89: Overflow Scenario Prices	112
Table 90: Overflow Cost	112
Table 91: KPIs of the project	115
Table 92: Small datacentre analysis for KPI objectives	117
Table 93: Large datacentre analysis for KPI objectives	118
Table 94: Summary of features in the 5G-TRANSFORMER Design	122

Table 95: Abstraction of the radio coverage area	130
Table 96: Costs for two arctica 3200xlp Commercial Spine Switch	142
Table 97: Costs for two arctica 4806xp Commercial Leaf Switch	143
Table 98: Costs for two arctica 4804ip Management Switches.....	143
Table 99: Bill of materials for network node.....	143
Table 100: Costs Estimated for the chosen Network Node defined in Table 13.....	144
Table 101: Controller node's material bill	144
Table 102: Compute Node's Bill of Materials.....	145
Table 103: Costs estimated for the compute node	145
Table 104: Block storage hdd's bill of materials.....	146
Table 105: Block sotrage ssd's bill of materials	146
Table 106: Bill of materials for object storage low density	147
Table 107: Bill of materials for object storage with high density.....	148
Table 108: Bill of materials for object storage archiving.....	148
Table 109: Small private cloud (8 x 5) cost.....	149
Table 110: The capacity provided by small private cloud	149
Table 111: Costs proportion by regard to the tco.....	149
Table 112: Medium private cloud (8 x 5) cost	151
Table 113: The capacity provided by medium private cloud	152
Table 114: Large private cloud (24 x 7) cost.....	153
Table 115: The capacity provided by large private cloud	154

List of Equations

Equation 1: Number of spine ports	39
Equation 2: Cost for the spine switches per month	39
Equation 3: Cost for the rack units per month	39
Equation 4: Cost for the power per month	39
Equation 5: Total cost per month	39
Equation 6: Cost per core port per month	39
Equation 7: Max number of leaf switches that can be connected to the spine	40
Equation 8: Max number of leaf ports in this configuration	40
Equation 9: Number of leaf ports	40
Equation 10: Cost for the leaf switches per month	40
Equation 11: Cost for the rack units per month	40
Equation 12: Cost for the power per month	40
Equation 13: Cost for the spine ports per month	40
Equation 14: Total cost per month	40
Equation 15: Cost per core port per month	40
Equation 16: Number of management ports	40
Equation 17: Cost for the management switches per month	40
Equation 18: Cost for the rack units per month	40
Equation 19: Cost for the power per month	40
Equation 20: Cost for the leaf ports per month	41
Equation 21: Total cost per month	41
Equation 22: Cost per management port per month	41
Equation 23: Cost per month for the hardware	41
Equation 24: Cost per month for the rack space	41
Equation 25: Cost per month for the power	41
Equation 26: Cost per month for the leaf network ports	41
Equation 27: Cost per month for the management network ports	41
Equation 28: Cost per month per controller node	41
Equation 29: Cost per month for the hardware	42
Equation 30: Cost per month for the rack space	42
Equation 31: Cost per month for the power	42
Equation 32: Cost per month for the leaf network ports	42
Equation 33: Cost per month for the management network ports	42
Equation 34: Cost per month per controller node	42
Equation 35: Cost per month for the hardware	43
Equation 36: Cost per month for the rack space	43
Equation 37: Cost per month for the leaf network ports	43
Equation 38: Cost per month for the management network ports	43
Equation 39: Cost per month for the power	43
Equation 40: Cost per month per compute node	43
Equation 41: vCPUs per node	43
Equation 42: RAM (in GB) per node	43
Equation 43: Cost per vCPU per month	43
Equation 44: Cost per vCPU per hour	43
Equation 45: Cost per GB of RAM per month	43
Equation 46: Cost per GB of RAM per hour	43

Equation 47: Cost per month for the hardware	44
Equation 48: Cost per month for the rack space	44
Equation 49: Cost per month for the leaf network ports	44
Equation 50: Cost per month for the management network ports	44
Equation 51: Cost per month for the power	44
Equation 52: Cost per month for the storage cluster	44
Equation 53: OSDs in the cluster	44
Equation 54: Total raw capacity of the cluster (in GB)	44
Equation 55: Usable capacity ratio	44
Equation 56: Usable capacity of the cluster (in GB)	44
Equation 57: Overcommitted usable capacity of the cluster (in GB)	44
Equation 58: Maximum read IOPS for the cluster (after cache is exhausted)	44
Equation 59: Maximum write IOPS for the cluster (after cache is exhausted)	44
Equation 60: Cost per GB per month	44
Equation 61: Cost per GB per hour	44
Equation 62: Cost per month for the hardware	45
Equation 63: Cost per month for the rack space	45
Equation 64: Cost per month for the leaf network ports	45
Equation 65: Cost per month for the management network ports	45
Equation 66: Cost per month for the power	45
Equation 67: Cost per month for the storage cluster	45
Equation 68: OSDs in the cluster	45
Equation 69: Total raw capacity of the cluster (in GB)	45
Equation 70: Usable capacity ratio	45
Equation 71: Usable capacity of the cluster (in GB)	45
Equation 72: Overcommitted usable capacity of the cluster (in GB)	45
Equation 73: Maximum read IOPS for the cluster (after cache is exhausted)	45
Equation 74: Maximum write IOPS for the cluster (after cache is exhausted)	45
Equation 75: Cost per GB per month	45
Equation 76: Cost per GB per hour	45
Equation 77: Cost per month for the hardware	46
Equation 78: Cost per month for the rack space	46
Equation 79: Cost per month for the leaf network ports	46
Equation 80: Cost per month for the management network ports	46
Equation 81: Cost per month for the power	46
Equation 82: Cost per month for the storage cluster	46
Equation 83: OSDs in the cluster	46
Equation 84: Total raw capacity of the cluster (in GB)	46
Equation 85: Usable capacity ratio	46
Equation 86: Usable capacity of the cluster (in GB)	46
Equation 87: Overcommitted usable capacity of the cluster (in GB)	46
Equation 88: Maximum read IOPS for the cluster (after cache is exhausted)	46
Equation 89: Maximum write IOPS for the cluster (after cache is exhausted)	46
Equation 90: Cost per GB per month	46
Equation 91: Cost per GB per hour	46
Equation 92: Cost per month for the hardware	47
Equation 93: Cost per month for the rack space	47
Equation 94: Cost per month for the leaf network ports	47

Equation 95: Cost per month for the management network ports	47
Equation 96: Cost per month for the power	47
Equation 97: Cost per month for the storage cluster.....	47
Equation 98: OSDs in the cluster	47
Equation 99: Total raw capacity of the cluster (in GB).....	47
Equation 100: Usable capacity ratio.....	47
Equation 101: Usable capacity of the cluster (in GB).....	47
Equation 102: Overcommitted usable capacity of the cluster (in GB)	47
Equation 103: Maximum read IOPS for the cluster (after cache is exhausted)	48
Equation 104: Maximum write IOPS for the cluster (after cache is exhausted).....	48
Equation 105: Cost per GB per month.....	48
Equation 106: Cost per GB per hour	48
Equation 107: Cost per month for the hardware	48
Equation 108: Cost per month for the rack space.....	48
Equation 109: Cost per month for the leaf network ports	48
Equation 110: Cost per month for the management network ports	48
Equation 111: Cost per month for the power	48
Equation 112: Cost per month for the storage cluster.....	48
Equation 113: Disks in the cluster	48
Equation 114: Total raw capacity of the cluster (in GB)	48
Equation 115: Usable capacity ratio.....	48
Equation 116: Usable capacity of the cluster (in GB).....	49
Equation 117: Overcommitted usable capacity of the cluster (in GB)	49
Equation 118: Maximum read IOPS for the cluster (after cache is exhausted)	49
Equation 119: Maximum write IOPS for the cluster (after cache is exhausted).....	49
Equation 120: Cost per GB per month.....	49
Equation 121: Cost per GB per hour	49
Equation 122: Hardware cost per month	49
Equation 123: Total number of rack units.....	49
Equation 124: Rack units cost.....	50
Equation 125: Power consumption.....	50
Equation 126: Power cost	50
Equation 127: Number of 10Gbps ports	50
Equation 128: Number of 1Gbps ports	50
Equation 129: Cost per month.....	50
Equation 130: Cost per 3 years.....	50
Equation 131: Staff cost.....	51
Equation 132: VMWare as a license	51
Equation 133: TCO	51
Equation 134: Number of ports offering 40Gbps	51
Equation 135: Number of ports offering 10Gbps	51
Equation 136: Number of ports offering 1Gbps	51
Equation 137: Number of vCPUs	51
Equation 138: RAM size.....	52
Equation 139: Block storage HDD.....	52
Equation 140: Block storage SSD	52
Equation 141: Object storage Low Density.....	52
Equation 142: Object storage High Density.....	53

Equation 143: Cost for a vCPU	53
Equation 144: Cost for a RAM.....	53
Equation 145: Cost for a block storage HDD	53
Equation 146: Cost for a block storage SDD	53
Equation 147: Cost for an object storage low density	53
Equation 148: Cost for an object storage high density.....	53
Equation 149: Cost Stuff vCPU	54
Equation 150: Cost Stuff RAM	54
Equation 151: Cost Stuff Storage.....	54
Equation 152: Cost power vCPU.....	54
Equation 153: Cost power RAM	54
Equation 154: Cost Power Block Storage HDD	54
Equation 155: Cost power Block Storage SDD.....	54
Equation 156: Cost power Object storage low density.....	55
Equation 157: Cost power Object storage high density	55
Equation 158: Cost License VMWare vCPU	55
Equation 159: Final price	55
Equation 160: Final cost	55
Equation 161: Breakeven price	59
Equation 162: Price.....	60
Equation 163: Profit margin.....	60
Equation 164: Direct revenue items	60
Equation 165: Federation margin	61
Equation 166: Real utilization of the local infrastructure	61
Equation 167: Revenue	61
Equation 168: Profit	61
Equation 169: Revenue-Profit	62
Equation 170: Revenue-Price-Profit.....	62
Equation 171: Revenue-Price-Breakeven-Profit.....	62
Equation 172: Simplified Revenue-Price-Breakeven-Profit	63
Equation 173: Local-Simplified Revenue-Price-Breakeven-Profit.....	63
Equation 174: Local Profit-Cost ratio.....	63
Equation 175: Ideal Revenue-Price-Breakeven-Profit	67
Equation 176: Ideal Profit-Cost ratio	67
Equation 177: Federated-Simplified Revenue-Price-Breakeven-Profit	68
Equation 178: Federated Profit-Cost ratio	68
Equation 179: Price variability in the overflow scenario.....	69
Equation 180: Service demand distribution	71
Equation 181: Event arrival distribution	80
Equation 182: Life-time distribution	83
Equation 183: Available resources.....	84
Equation 184: Number of scaled events distribution.....	85
Equation 185: Scaled resources distribution	85

List of Acronyms

Acronym	Description
5GT-MTP	Mobile Transport and Computing Platform
5GT-SO	Service Orchestrator
5GT-VS	Vertical Slicer
ABNO	Applications-Based Network Operations
AD	Administrative Domain
AM	Abstraction Manager
API	Application Programming Interface
AppD	Application Descriptor
AS/PCE	Active Stateful Path Computation Element
BSS	Business Support System
CIM	Cooperative Information Manager
CN	Core Network
COP	Control Orchestration Protocol
CQI	Channel Quality Indicator
CSAR	Cloud Service Archive
CSMF	Communication Service Management Function
DCSP	Data Centre Service Provider
EM	Element Manager
EVS	Extended Virtual Sensing
E/WBI	Eastbound/Westbound Interface
E2E	End-to-end
FTE	Full-Time-Equivalent staff.
GMPLS	Generalized Multi-Protocol Label Switching
GTP	GPRS Tunnelling Protocol
HMI	Human Machine Interface
HSS	Home Subscriber Server
HTTP	HyperText Transfer Protocol
IM	Information Model
JSON	JavaScript Object Notation
LCM	LifeCycle Management
LL	Logical Link
LSA	Link Selection Algorithm
MANO	Management and Orchestration
MEA	Multi-access edge application
MEAO	Multi-access edge application orchestrator
MEC	Multi-access edge computing
MEO	Multi-access edge orchestrator
MEP	Multi-access edge platform
MEPM	Multi-access edge platform manager
MEPM-V	Multi-access edge platform manager - NFV
MES	Multi-access edge service
MLPOC	Multiple Logical Point of Contact
MME	Mobility Management Element
MVNO	Mobile Virtual Network Operator
NBI	Northbound Interface
NF	Network Function
NF FG	NF Forwarding Graph
NFP	Network Forwarding Path

NFV	Network Function Virtualization
NFVlaaS	NFVI as a Service
NFV-NS	Network Service
NFV-NSI	Network Service Instance
NFV-NSO	Network Service Orchestrator
NFVI	Network Functions Virtualisation Infrastructure
NFVI-PoP	Network Functions Virtualisation Infrastructure Point-of-Presence
NFVlaaS	NFVI as a Service
NFVO	NFV Orchestrator
NFVO-RO	Resource Orchestrator
NNI	Network to Network Interface
NS	Network Slice
NS-OE	NS Orchestration Engine
NSaaS	Network Slice as a Service
NSD	Network Service Descriptor
NSI	Network Slice Instance
NSMF	Network Slice Management Function
NSSI	Network Slice Subnet Instance
NSSMF	Network Slice Subnet Management Function
NST	Network Slice Template
OEM	Original Equipment Manufacturer
OF	OpenFlow
OSS	Operating Support System
PA	Physical Application
PGW-C	Packet Gateway Control Plan
PGW-U	Packet Gateway User Plan
PNF	Physical Network Function
PNFD	Physical Network Function Descriptor
QoS	Quality of Service
RAN	Radio Access Network
REST	Representational State Transfer
RMA	Resource Management Application
RNIS	Radio Network Information
ROI	Return On Investment
ROOE	Resource Orchestration Engine
SBI	Southbound Interface
SDK	Software Development Kit
SDN	Software-Defined Networking
SGW-C	Serving Gateway Control Plan
SGW-U	Serving Gateway User Plan
SLA	Service Level Agreement
SLO	Service Level Objective
SLPOC	Single Logical Point of Contact
SM	Service Manager
S-NSSAI	Single Network Slice Selector Assistance Information
SOE	Service Orchestration Engine
TCO	Total Cost of Ownership
TD	Technology Domain
TMOP	5G-TRANSFORMER Mobile Transport and Computing Platform Operator

TMVS	5G-TRANSFORMER Managed Vertical Service
TOSCA	Topology and Orchestration Specification for Cloud Applications
TOR	Top of the Rack
TRF	5G-TRANSFORMER Resource Federation
TS	5G-TRANSFORMER Service
TSC	5G-TRANSFORMER Service consumer
TSF	5G-TRANSFORMER Service Federation
TSP	5G-TRANSFORMER Service Provider
TUVS	5G-TRANSFORMER Unmanaged Vertical Service
VA	Virtual Application
VA FG	VA Forwarding Graph
vCDN	virtual Content Delivery Network
VIM	Virtual Infrastructure Manager
VISP	Virtualization Infrastructure Service Provider
VL	Virtual Link
VLD	Virtual Link Descriptor
VNF	Virtualised Network Function
VNFFG	VNF Forwarding Graph
VNFC	Virtualised Network Function Component
VNFD	Virtualised Network Function Descriptor
VNFM	Virtual Network Functions Manager
VS	Vertical Service
VSaaS	Vertical Service as a Service
VSBlueprint	Vertical Service Blueprint
VSD	Vertical Service Descriptor
VSI	Vertical Service Instance
VXLAN	Virtual Extensible LAN
WAN	Wide Area Network
WIM	Wide area network Infrastructure Manager
XCI	5G-Crosshaul Control Infrastructure
XFE	5G-Crosshaul Forwarding Element
YAML	YAML Ain't Markup Language
YANG	Yet Another Next Generation

Executive Summary and key contributions

The vision of the 5G-TRANSFORMER project is to design a platform capable of delivering on-demand tailored services to vertical industries while meeting their very specific requirements. This vision fosters the evolution of the existing business models in today's large networks toward a new ecosystem that gathers different players and stakeholders into new interactions at different layers of the value chain. Following this perspective, 5GT investigates, beyond the technical innovations, the economic impacts of the deployment of the designed platform architecture and its services on the vertical industries businesses.

The main contributions of this deliverable are as follows:

- A detailed techno-economic study that investigates the cost models of the 5GT vertical use case services selected for implementation in the final proof of concepts. It comprises both analytical and experimental studies in lab that allow through the simulation of different scenarios the estimation of the cost of deployment of 5GT services with regards to the reality of the ecosystem in terms of heterogeneous interactions between the different players and stakeholders.
- The final 5G-TRANSFORMER architecture design describing two main extensions to its building components; (i) a RAN abstraction solution; (ii) a MEC support solution. A description of the updates and additions to the main building blocks, namely the 5GT-VS, the 5GT-SO and the 5GT-MTP is provided.

For the sake of simplicity and clarity for the reader, the detailed description of the final 5GT architecture and its components is referred to the one reported in previous deliverable D1.3 [3] and only the new extensions are reported in this deliverable.

1 Introduction

The vision of the 5G-TRANSFORMER project is that Mobile Transport Networks have to transform from today's rigid interconnection solutions into an SDN/NFV-based 5G Mobile Transport and Computing Platform (MTP) able to support a wide range of networking and computing requirements that meet the specific needs of a wide variety of vertical industries.

The 5G-TRANSFORMER platform has been designed in this context to provide a proof of concept of an SDN/NFV-based solution which implements and integrates novel concepts like network slicing, resource federation over multiple domains and MEC, in order to deploy tailored vertical services addressing heterogeneous service requirements for the vertical industries. The result is a multi-layered architecture design that defines three main functional blocks described in the baseline architecture design, published in deliverable D1.2 [2], including: the Vertical Slicer (5GT-VS), the Service Orchestrator (5GT-SO) and the Mobile Transport and Computing Platform (5GT-MTP).

The first implementation release of the 5GT platform architecture and its components has been delivered during the first year of the project and reported in D1.2 [2]. It has been then enhanced during the second year to deliver the refined 5GT architecture design reported in deliverable D1.3 [3]. This refined architecture has been implemented as part of the second release of the 5GT platform and includes mainly extensions that were motivated by advanced research conducted in the project as well as the development of 5GT architecture building components.

Aside from the technical innovations brought by this layered architecture, the 5GT system promotes new mechanisms for the monetization of the new network generation and new business models for the vertical industries. Hence, 5GT offers the possibility of hosting over-the-top applications in the network, leveraging time proximity and context information exposed by the network. This represents a unique value that can be exploited for revenue generation by operators and application service providers alike, thus creating new value chains and a variety of interaction models with operators and IT actors.

During the first phase of the project, we have reported in D1.1 [1] the results of a study that aimed at analyzing the 5GT ecosystem and its stakeholders' role models while deriving the requirements of the vertical domains involved in the project, namely, the automotive, the eIndustry, the entertainment, the eHealth and the MVNO verticals. This first study has served as an input for a more detailed techno-economic study that is proposed in this deliverable.

The current deliverable underlines the final outcome of this 5GT project and exposes the project results in two main folds:

First, in section 2, we present the results of the techno-economic study of the 5GT system. The study analyses the market and business implications of the deployment of the 5GT platform for the stakeholders and actors involved in the 5GT ecosystem. It mainly identifies the business opportunities for operators and service providers, their capability to capture value and the potential appearance of new actors and roles. It also assesses the project's impact on verticals' and providers' OPEX and CAPEX, design

business cases for both verticals and providers identifying value chain configurations and related business models and market opportunities from different players' perspectives. The study is hence articulated around an analytical study based on the modeling of the 5GT vertical use cases to estimate their pricing model and an experimental study in lab based on simulations of the system through different configuration scenarios allowing the evaluation of the 5GT platform cost variations and implications on the different actors' revenues and benefits. The final outcome of the study highlights some recommendations about the economic impacts of the 5GT platform for efficient economic mechanisms that ensure sustainability for actors.

In the second part of this deliverable, in section 3, we report some updates to the final 5GT architecture design, in which we propose some further enhancements to its main building blocks with regards to the refined architecture design proposed in D1.3 [3]. Hence, we propose a way forward to support RAN abstraction at the MTP level in order to be able to expose RAN as a resource to vertical network slices, and a solution that allows the integration of MEC in the 5GT platform at the 5GT-SO, 5GT-VS and 5GT-MTP components, and how to support it at the service on-boarding and instantiation workflows.

2 5G-TRANSFORMER Techno-economic analysis

This section covers the techno-economic analysis of the 5G-TRANSFORMER system. A techno-economic analysis is a profit-cost comparison of a project performed by using different assessment methods.

The mentioned assessment methods are used for tasks with different purposes as indicated hereafter:

- Evaluate the economic feasibility of the 5G-TRANSFORMER project.
- Setting the appropriate prices to all the vertical services so that the project is profitable.
- Investigate the cash flows over a specific timeframe: In this project, one-year timeframe is considered.
- Evaluate the evolution of the demand over time.
- Reaching some lessons learnt about the economic aspects of the project and proposes some future steps.
- Etc.

In the 5G-TRANSFORMER ecosystem, multiple actors or stakeholders interact with each other to provide services to tenants.

As a result, and in order to cover the previous points, two studies are proposed: Analytical study and experimental analysis:

- *Analytical study*: The analytical study provides an extensive theoretical analysis including: Use case description, Infrastructure cost modelling, Pricing modelling and dimensioning analysis. With this analysis, the actors of the of the 5G-TRANSFORMER ecosystem are modelled as well as the flows among them in terms of resources and money
- *Experimental study*: The aim of the experimental study is to carry out some simulations that validate the analytical study and allows to reach some conclusions for the techno-economic analysis. It includes: the methodology, the system structure, the actors in the study, the scenario setting, the tool description, the service pricing and lessons learnt.

The results that are obtained in the experimental study validate the analytical study as indicated in sections 2.3.5.2 to 2.3.5.6.

Hereafter, the background of this section can be found before entering the analytical and experimental analysis.

2.1 Background

Before getting into details about the analytical analysis in 2.2 and the experimental analysis in 2.3 some background information is given in this section regarding the different approaches to pricing models that lead to the profitability of the 5G-TRANSFORMER project.

It is possible to say that cloud computing is reaching the “early majority” in the cloud adoption life-cycle and nowadays it can be considered a mainstream market of the information technologies as indicated in [10].

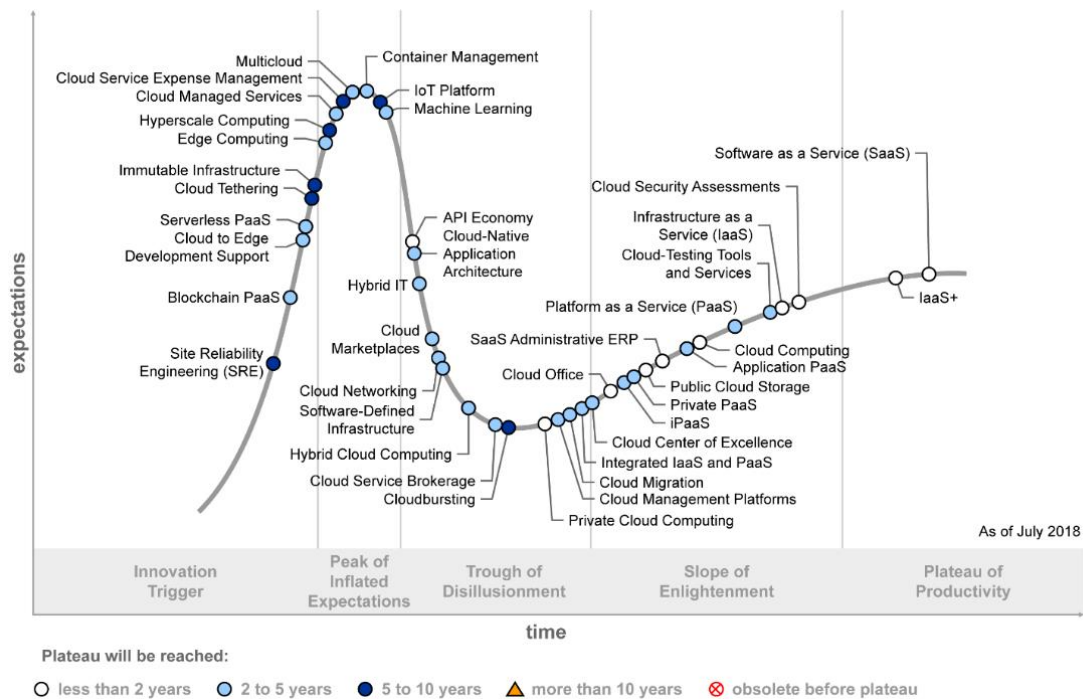


FIGURE 1: GARTNER'S HYPE CYCLE FOR CLOUD COMPUTING [10]

Accordingly, there are currently many pricing approaches for the cloud considering interdisciplinary cases providing different schemas.

In [11], a classification of more than 60 different pricing models is proposed. In all of them the weaknesses and strengths of the price model are discussed. Pricing models can be value-based, market-based or cost-based.

- *Value-based*: Value-based pricing models consider the value as the strategy for setting their approach. They integrate the customer, the experience and service pricing as the main pricing models.
- *Market-based*: Market-based pricing models consider the market as the strategy for setting their approach. They integrate the free and pay later, the retail-based, auction and online pricing as the main pricing models.
- *Cost-based*: Cost-based pricing models consider the cost as the strategy for setting their approach. They integrate the expenditure-based, the resource-based and the utility-based pricing as the main pricing models.

For the development of 5G-TRANSFORMER price modelling, the cost-based approach seems to be the most appropriate as there is no existing information on the value that users will perceive from the use of the 5G-TRANSFORMER services nor on the market value that they will obtain.

2.2 Analytical study

The aim of the analytical study is to provide an extensive analysis of the 5G-TRANSFORMER vertical use cases and explain how they interact within the 5G-TRANSFORMER ecosystem. This includes the modelling of the different actors that play a role in the study, namely:

- *The 5G-TRANSFORMER vertical use cases:* Including entertainment, MVNO, automotive, non-emergency eHealth, emergency eHealth and robotics vertical use cases. The modelling of each vertical use case includes the resources needs for the service's VNFs in terms of CPUs, RAM, Disk and licenses (understood in this project as the price that has to be paid for one user using the service) as well as the connectivity assets in terms of latency and bandwidth that must be guaranteed in order to deploy the service.
- *The Infrastructure:* The infrastructure is modelled defining the datacentres that will be considered in 5G-TRANSFORMER and the cost that the infrastructure will have, considering both the CAPEX and OPEX expenditures.
- *The pricing:* Once the vertical use cases and the infrastructure are defined, the price that has to be set for the 5G-TRANSFORMER services is modelled.
- *The demand:* Finally, a first approach to the demand modelling is done, regarding the number of services' request that might be expected for each of the services.

2.2.1 Vertical use case modelling

In this section, a description of the 5G-TRANSFORMER entertainment, MVNO, automotive, non-emergency eHealth, emergency eHealth and robotics use cases modelling can be found.

2.2.1.1 Entertainment

The Entertainment vertical is a wide domain that covers many areas of interest related to human entertainment and leisure. The current analysis will focus on a specific area of the entertainment domain, which is the Sport Events, and especially on everything related to fan interaction, known as FAN ENGAGEMENT.

The main goal of fan engagement is making the venue smart and following the fan along the fan journey. Give the fans more interaction, more engagement and make them feel like they are more part of the game than they ever could be before.

The vertical use case can be separated in two different scenarios addressing different problems: open venues and closed venues. In open venues, the key aspect is the density of users demanding a high data rate. Meanwhile, in closed venues, the key point is the broadband access in all the points with high data rate and very low latency. In both cases, there are a lot of actors that influence in the fan engagement. The main actors are:

- *Sport fans:* These are the main actors because the fans as final are the main target of the provided content.
- *Mobile Infrastructure Provider:* These operator/s can deploy extra infrastructure during certain events to support more users or provide a higher data rate.
- *Network Provider:* Sometimes it is also necessary to deploy specific communication lines to connect the venues with the data centers.
- *Main IT Integrator:* Organization that wins a tender and becomes the official integrator for a given event or set of events.

There can be also other actors like press, athletes, sponsors, ticketing provider, etc. With actual networks, the number of users supported in these events it is not enough. Here is where 5G can help to support the densely-packed environments and allow distributing immersive experiences or Ultra High Definition (UHD) content. The 5G-

TRANSFORMER entertainment use case comes with a 5G and orchestrated solution to improve the fan engagement. The use case topology can be found in Figure 2 and its components are described as follows:

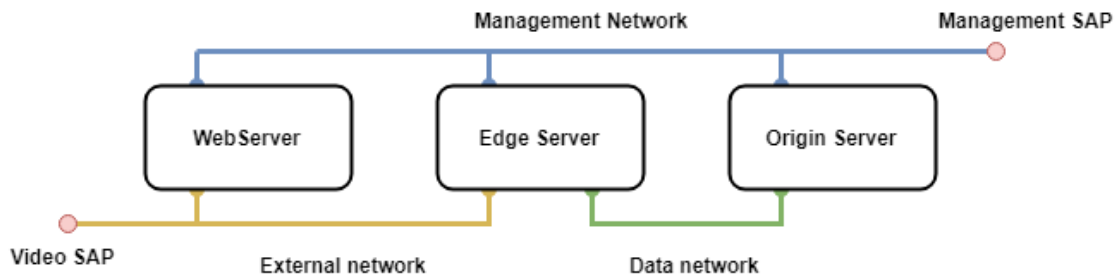


FIGURE 2: ENTERTAINMENT USE CASE TOPOLOGY

- *Origin Server (SPR1)*: This origin server receives the streams from the cameras and transcodes the content.
- *Edge Server (SPR2)*: Acts as a video cache, serving the video segments to the web server.
- *Web Server*: This Web Server provides a web player to show the videos of the system to the end-users.
- *Management Network*: It is used to manage the different VNF's and the configuration of the service.
- *External Network*: This network is accessible for the end-users and transports the video from the Edge Server to the Web Server.
- *Data Network*: The data network interconnects the Origin Server with the Edge Server transporting the transcoded video.

In this use case, it is possible to set some auto-scaling rules in the Network Service descriptor to create new instances of the Edge Server VNF. This can be configured for example, to support a higher number of users as shown in Figure 3.

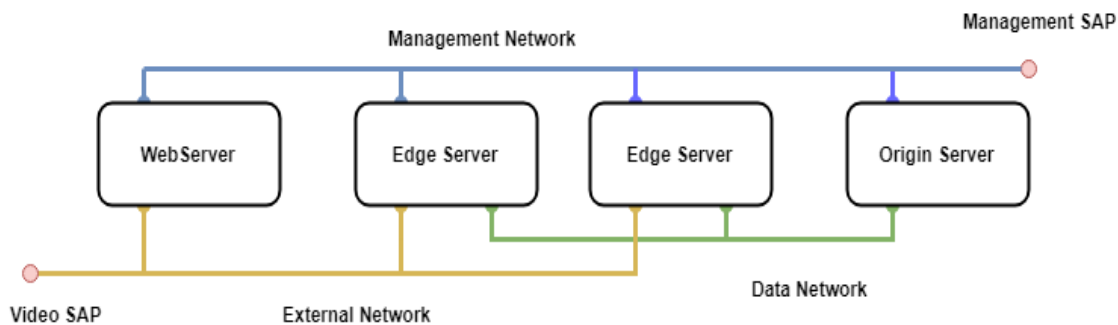


FIGURE 3: EXTENDED ENTERTAINMENT USE CASE TOPOLOGY

The use case service VNF's resources in terms of computing and networking are detailed in Table 1:

TABLE 1: ENTERTAINMENT USE CASE VNFs' RESOURCES

Service 1: Entertainment				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
SPR1	2	8	50	30
SPR2	2	2	5	30
WebServer	1	1	5	30

In Table 2, the use case service link assets can be found:

TABLE 2: ENTERTAINMENT USE CASE VNFS' LINKS ASSETS

Service 1: Entertainment		
Link	BW [Mbps]	Latency [ms]
Management Network	10	50
External Network	250	50
Data Network	10	50

2.2.1.2 MVNO

As described in deliverable D1.1 [1], the MVNO use case is built on the offering of a Network Slice as a Service (NSaaS) by a Mobile Network Operator (MNO) hosting a Mobile Virtual Network Operator (MVNO); for instance, the MNO would rely on network slicing combined to services like EPCaaS and IaaS in order to set up a virtual mobile network and provide connectivity network services to the MVNO.

In addition, verticals can be seen as customers of an MNO or an MVNO.

As explained in D1.1 [1], we chose to focus on the vEPCaaS UC. The MNO/MVNO UC aims at demonstrating the deployment and operation of a Network Slice with vEPC in “as a service” mode in order to build an MVNO service.

In deliverable D3.3 [7], a description of the service deployment requirements of the vEPCaaS use case has been provided as part of the Blueprint.

The network functions defined and required in the deployment of the use case could be divided in three types of VNFs based on their functionality:

- *Control plane functionalities:* MME, HSS, AAA, DHCP, S/PGW-C, Controller, OVS
- *Data plane functionalities are not deployed with OSM/Cloudify and that could be seen as PNF:* S/PGW-U
- *Management and monitoring capabilities:* Customer Care for subscription provisioning and management, and Dashboard to display active user sessions

The capacity deployed for these VNFs is static. VNFs DF and/or the number of VMs in VNF (in our case only the MME) are configured at instantiation time. Besides, scaling requires re-deployment.

The lifetime of the service defined as the time for which we consider the service as active, depends on the events.

For instance, a service can be active for weeks or even hours e.g. for a sport meeting event, or for years e.g. network of a company.

The specific SLA requirements for this use case are:

- *SAP.* QoS class: GBR, priority (real time, high, medium, low)
- *Service availability:* low (<95% is required)
- *Service reliability:* medium (95%-99%)
- *Traffic Density:* 1000 users/km²

The use case topology can be found in Figure 4:

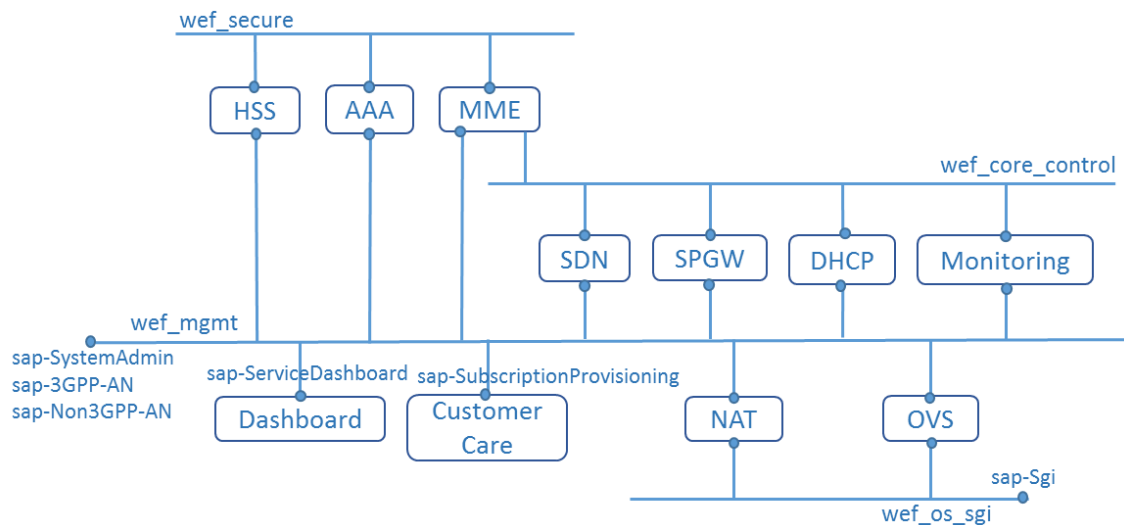


FIGURE 4: MVNO USE CASE TOPOLOGY

In order to meet the performance criteria requested by the customer, it must be ensured that sufficient virtualized resources have been allocated to the network services during the instantiation of the network slice.

A method to assess the smallest amount of resources to run a network service at a level of performance consists of obtaining from the VNF provider the template function that gives the promised performances according to the allocated resources.

To set up the template function, the VNF provider measures the upper limits of the capacity of the VNF to handle user traffic while expanding the size and the number of VNF instances. Regarding our use case, the user session creation rates handled by the MME before congestion will be measured at different scaling steps to establish the promised performance function.

Knowing the amount of virtualized resources necessary to run a network slice instance meeting the customer criteria and combined with the cost of virtualized resources whose calculation details are provided in another section, the actual infrastructure cost of a network slice instance is fully determined for the vEPCaaS provider.

Hereafter is an example of the amount of resources that could be required to run a vEPC for 1000 users, 10 sessions per second.

The use case service VNFs' resources can be found in Table 3.

TABLE 3: MVNO USE CASE VNFs' RESOURCES

Service 2: MVNO				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
MME	1	1	10	1000
HSS	1	1	10	1000
AAA	1	1	10	1000
DHCP	1	1	10	1000
S/PGW-C	1	1	10	1000
OVS	1	1	10	1000
NAT	1	1	10	1000
SDN	2	2	20	1000
Dashboard	1	1	10	1000

Customer care	1	1	10	1000
Controller	1	1	10	1000

According to [13], the number of licenses required for a high traffic density (1000 users/km²) is 1000, as indicated above.

MME is the main entry point of the core network for signaling messages issued by the access network. It runs procedures to manage sessions when a subscriber sets up or releases connection within its home or a visiting network.

For instance, during the attach procedure, the MME establishes a security context for the user containing data derived from the authentication, then it requests new IP address from the DHCP server and creates bearers to transport the user's data. As it handles many concurrent sessions and executes several procedures, it needs larger infrastructure network resources (CPU, RAM, and storage) than the other network functions to correctly absorb the traffic demand. The SDN controller needs also a large capacity as it supports features offered to the upper layer (NBI) applications to manage user traffic forwarding in the data plane.

As for the other network functions HSS, AAA, DHCP server, S/PGW-C, Dashboard, Customer Care, they need less infrastructure resources to run because they are involved in a specific phase of the processing of the user session.

The links represented in the use case topology diagram depicted in Figure 3 are virtual links and represent complete networks with specific characteristics.

We have three types of links:

1. *Wef_secure*: this network is dedicated to the transport of sensitive data such as derived keys and authentication vectors
2. *Wef_core_control*: this network supports signaling interfaces of the core network
3. *Wef_mgmt*: In this network, a floating IP address is assigned to the network functions and allows them to communicate with the external network.
4. *Wef_os_sgi*: this network is the SGi LAN

The use case service link assets can be found in Table 4.

TABLE 4: MVNO USE CASE VNFS' LINKS ASSETS

Service 2: MVNO		
Link	BW [Mbps]	Latency [ms]
wef_secure	100	50
wef_core_control	100	50
wef_mgmt	100	50
wef_os_sgi	1000	50

2.2.1.3 Automotive

The EVS (Extended Virtual Sensing) is a road safety application able to extend the view of on-board sensors to signal the presence of unseen vehicles or unexpected obstacles at intersections. The EVS has a global view of the monitored crossing, which may be exploited to provide key information to the On-Board Unit (OBU) taking decisions at the vehicles. The Extended Virtual Sensing (EVS) exploits real-time data that is collected by a 3rd party entity; the Cooperative Information Manager (CIM). The

collected information includes: the position, the heading, the speed and the acceleration for each vehicle in the monitored area. The data is provided to the EVS algorithm that estimates the probability of a collision.

In this way, vehicles can base their decisions on data fused from multiple information sources: vehicle data, on board sensors and V2I messages that act as virtual ADAS sensors. The EVS indeed extends the capability of on board sensors covering also scenarios where obstacles are hidden by buildings.

The goal of the Automotive PoC is to showcase the correct functioning of the collision avoidance service, while another service (in this case video streaming service) is active on-board.

Figure 5 represents the automotive use case deployment topology and shows the network entities involved in the use case demonstration.

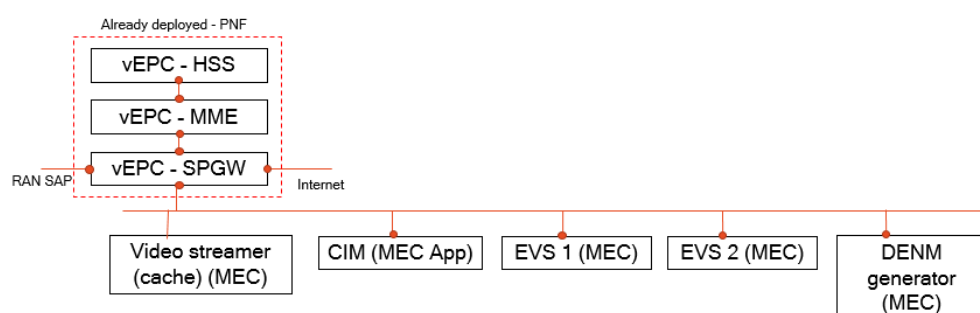


FIGURE 5: EVS AUTOMOTIVE UC TOPOLOGY

The use case service VNFs' resources can be found in Table 5:

- vCIM (Cooperative Information Manager): a VA owned by a trusted third party entity providing for each car maker a database just storing the CAM messages.
- vEVS: a vehicle collision detection algorithm. When a possible collision is detected an alert request is generated and sent to the vDENM generator.
- vDENM generator: sends DENM to the involved vehicles avoiding message conflicts.
- vEPC: receives message toward BBU and forward them to CIM.
- VS.
- VS algo.

TABLE 5: AUTOMOTIVE VNFs' RESOURCES

Service 3: Automotive				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
vCIM	2	4 GB	25 GB	4000
vEVS	1	2 GB	25 GB	4000
vDENM generator	1	2 GB	25 GB	4000
vEPC	1	2 GB	25 GB	4000
VS	1	1 GB	10 GB	4000
VS_algo	1	4 GB	10 GB	4000

According to [13], the number of licenses required for a high traffic density (4000 users/km²) is 4000, as indicated above. In addition, this source also provides information on the link characteristics.

The use case service link assets can be found in Table 6.

TABLE 6: AUTOMOTIVE VNFS' LINKS ASSETS

Service 3 Automotive		
Link	BW [Mbps]	Latency [ms]
int	50	5

2.2.1.4 eHealth

eHealth can be defined as the use of information and communication technologies (ICT) to deliver health services.

Among the possible examples of services that can be provided by eHealth systems, 5G-TRANSFORMER project focused on:

- The provision of a dedicated network slice that provides the minimum service when no emergencies are occurring (referred hereupon as Non-Emergency eHealth scenario). As examples of services are (i) a monitoring service provided to patients which, by aggregating and analysing the information from patient's wearable and portable devices, detects emergency situations; and (ii) an ambulance tracking service maintaining the location of all available ambulances.
- The dynamic adaptation of the dedicated non-emergency network slice to face emergency situations (referred hereupon as Emergency eHealth scenario) by scaling, if required, the allocated resources to the network slice and by deploying new services closer to the emergency location in order to facilitate the communication between the involved actors (e.g., patients, medical staff, ambulances and hospital) and to reduce the communication latency.

- **NON-EMERGENCY EHEALTH**

Even when no emergencies are ongoing, it is assumed that a minimum service is still available for non-emergency purposes by means of a dedicated network slice. As part of the non-emergency scenario, two services are provided: (i) a monitoring service; and (ii) an ambulance tracking service.

The former service includes the provisioning of radio connectivity to different kinds of wearable devices owned by the patients that are responsible for monitoring different health parameters. This information is then transmitted to a central server on the cloud (i.e., the "Central eServer") which is responsible for not only storing the data for historical purposes but also for processing and analysing it to detect possible emergency situations.

The latter service is also provided by the central server on the cloud and consists in tracking the location of all available ambulances, so that when an emergency occurs the nearest ambulance can be selected. Besides the required VNFs for providing radio connectivity to user devices, this scenario envisions the deployment of a VNF to act as the "Central eServer" which, since the provided services do not impose very strict latency requirements, could be deployed on the cloud. In terms of capacity required for

deployment of the required VNFs, the resources are allocated in order to provide a best effort service.

Although it is expected that the number of monitored patients remains stable over time, scaling the allocated resources of the dedicated network slice is still required in order to be able to accommodate an increase/decrease of the number of monitored patients. The use case topology can be found in Figure 6.

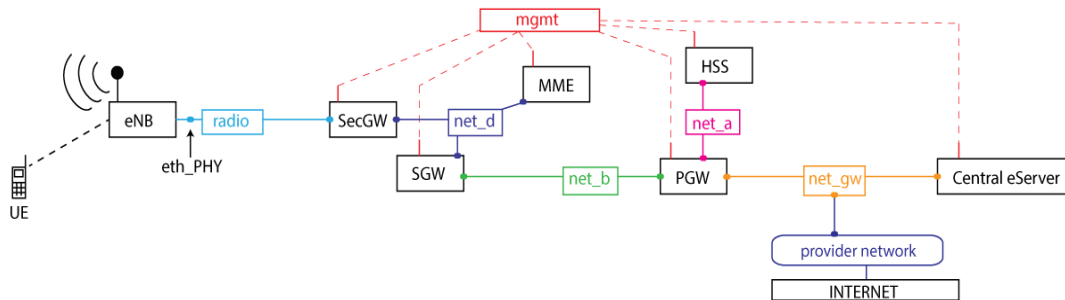


FIGURE 6: NON-EMERGENCY EHEALTH USE CASE TOPOLOGY

- *eNB*: is a physical device.
- *SecGW, SGW, PGW, MME, HSS and Central eServer* are the VNFs that build the non-Emergency service.
- *Networks (Virtual links): mgmt, radio, net_d, net_a, net_gw, provider network* are all networks used for interconnecting the deployed VNFs.

The use case service VNFs' resources can be found in Table 7.

TABLE 7: NON-EMERGENCY EHEALTH VNFs' RESOURCES

Service 4: Non-Emergency eHealth				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
SecGW	2	4	40	1
MME	2	4	20	1
HSS	1	2	20	1
SGW	2	4	20	1
PGW	2	4	20	1
Central eServer	2	4	100	1

The use case service link assets can be found in Table 8.

TABLE 8: NON-EMERGENCY EHEALTH VNFs' LINKS ASSETS

Service 4: Non-Emergency eHealth		
Link	BW [Mbps]	Latency [ms]
radio	100	5
net_d	100	10
net_b	100	10
mgmt	10	50
net_a	10	50
net_gw	10	50
provider network	10	50

The special requirements for this use case are the following:

- *Domain affinity*: non-emergency VNFs must be in same domain.
- **EMERGENCY EHEALTH**

When an emergency situation is triggered, the dedicated network slice described previously for the Non-Emergency eHealth scenario is extended to provide a localized service to the emergency. An edge service is then provided by a nearest server (i.e., the “Edge eServer”) to the emergency local. The “Edge eServer” is instantiated and deployed dynamically through a VNF, providing an emergency service to all interveners (i.e., patient, ambulances, doctors, etc) with lower latency.

In terms of capacity required for the deployment of the “Edge eServer”, the allocated resources need to scale dynamically according to the number of emergencies per area. On one hand, if an existing “Edge eServer” is already deployed near the emergency location, the same “Edge eServer” can be used to accommodate the new emergency situation, scaling the VNF resources if required. On the other hand, if the new emergency is in a different location, a new VNF to deploy the edge server is instantiated.

The use case topology can be found in Figure 7:

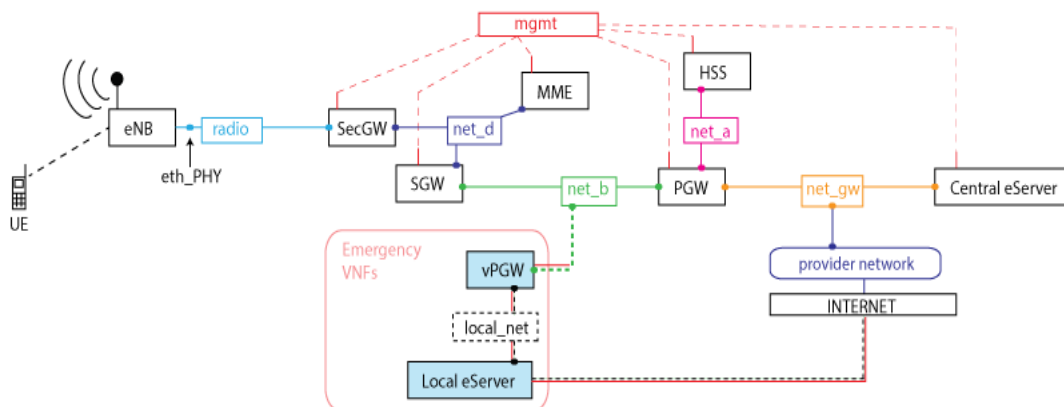


FIGURE 7: EMERGENCY EHEALTH CASE TOPOLOGY

- *eNB*: is a physical device.
- *SecGW, SGW, PGW, MME, HSS and Central eServer* are the VNFs that build the non-Emergency service.
- *Networks (Virtual links)*: *mgmt, radio, net_d, net_a, net_gw, provider network* are all networks used for interconnecting the deployed VNFs.
- *vPGW and Local eServer*: Are the two additional VNFs that build the Emergency service plus the existing ones from the non-emergency case.
- *Networks (VLs)*: *local_net* and connection to the Internet.

The use case service VNFs’ resources can be found in Table 9.

TABLE 9: EMERGENCY EHEALTH VNFs’ RESOURCES

Service 5: Emergency eHealth				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
SecGW	2	4	40	1
MME	2	4	20	1

HSS	1	2	20	1
SGW	2	4	20	1
PGW	2	4	20	1
Central eServer	2	4	100	1
vPGW	2	4	20	1
Local eServer	4	8	20	1

The use case service link assets can be found in Table 10.

TABLE 10: EMERGENCY EHEALTH VNFs' LINKS ASSETS

Service 5: Emergency eHealth		
Link	BW [Mbps]	Latency [ms]
radio	100	5
net_d	100	10
net_b	100	10
mgmt	10	50
net_a	10	50
net_gw	10	50
provider network	10	50
local_net	100	10
internet_connection	50	50

The special requirements for this use case are the following:

- *Domain affinity*: non-emergency VNFs must be in same domain. Emergency VNFs can be placed in an external domain.

2.2.1.5 eIndustry

Automatic Guided Vehicles (AGVs), as replacement of conveyor belts or as smart trolleys for local logistics, provide: i) a flexible management and reconfiguration of production line stages, something not possible with current conveyors; ii) an improved logistic as the same AGV can be used for taking product components to the line, to facilitate kitting, to bring final products to the loading bay for shipping. AGVs can also be used in other context like, for example, to move material among the different areas in a hospital.

In 5G TRANSFORMER, the eIndustry Cloud Robotics (CR) PoC simulates factory service robots and production processes that are remotely monitored and controlled in the cloud, exploiting wireless connectivity (5G) to minimize infrastructure cost, optimize processes, and implement lean manufacturing. The objective of the demonstrator is to verify the allocation of suitable resources based on the specific service requests to allow the interaction and coordination of multiple (fixed and mobile) robots controlled by remote distributed services, satisfying strict latency and bandwidth requirements.

The eIndustry use case leverages on a “factory control server” (FCS) where cloud robotics application runs, to control operations in the whole factory. An example of functionalities could be; arms and AGV tasks coordination, task planning, image recognition and so on. In terms of hardware needs, the typical requirements for processing in the edge cloud are 1CPU 3GHz, 16GB RAM. Possible scaling up is expected with the increase of number of robots in the area.

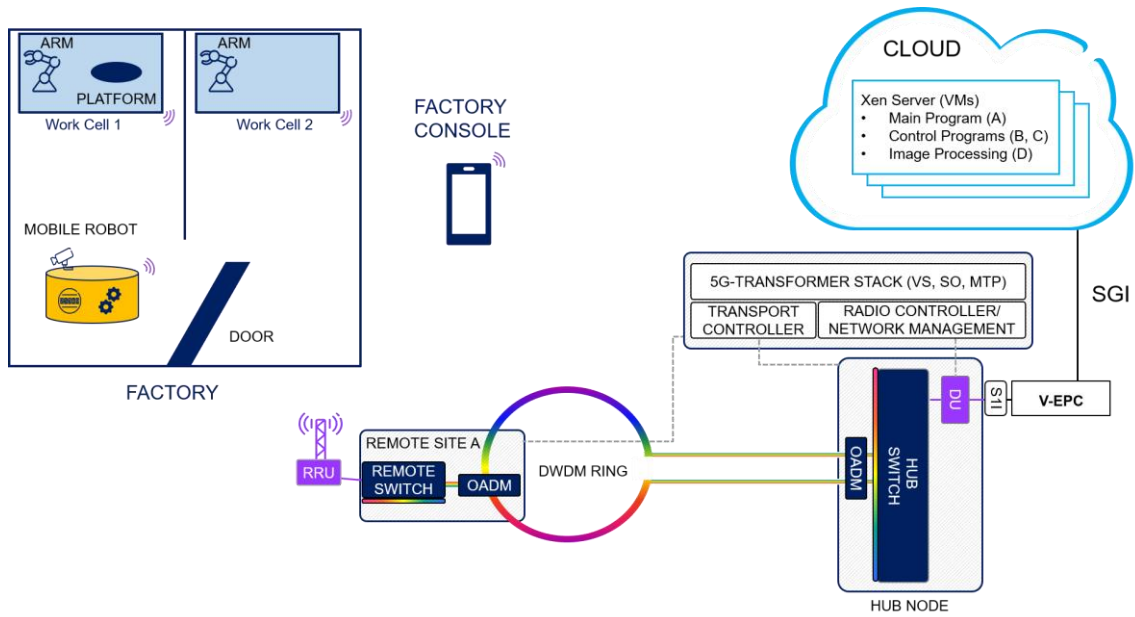


FIGURE 8: TOPOLOGY AND COMPONENTS OF THE EINDUSTRY CLOUD ROBOTICS DEMONSTRATOR

Communications among all the demo elements are ensured by a radio coverage having the EXhaul optical network infrastructure as transport layer. eXhaul serves as both backhaul and fronthaul to convey radio traffic on optical channels.

Novel photonic technologies are used to provide optical connectivity complemented by a dedicated agnostic framing, a deterministic switching module, and a flexible control paradigm. A layered concept is implemented to facilitate optimal interactions of transport and radio resources while preserving a well demarcated mutual independence. A detailed description of EXhaul can be found in [12].

With a fixed number of robots in the area, the required bandwidth capacity is static as the traffic load between robots and central server is basically stable. Dynamicity is possibly triggered by changing the number of robots. Typically, robot operations (AGVs) in an area is planned to operate for years even if, in some scenarios, shorter time windows could be considered.

The relevant service and service level is expected to be demanded for a quite long period without interruption. From a radio point of view, the eMBB and mMTC profiles are needed to transmit images (streaming) and sensor's data from the factory area to the cloud. On the opposite direction, URLLC profile is needed to send instructions from the cloud to the robot.

Requirements in term of end-to-end latency are illustrated in Figure 9 where are indicated the latencies needed to move the various functionalities (e.g. task planner) from the robot to the cloud (where the same functionality is provided by a virtual machine).

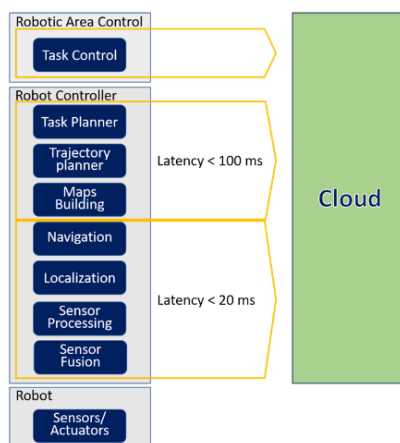


FIGURE 9: LATENCY REQUIREMENTS FOR eINDUSTRY USE CASE

The eIndustry use case topology is depicted on Figure 10. And leverage on the Ericsson vEPC. Specifically, the factory area (AGV, robot arms, etc.) is connected to the vEPC via the radio system (i.e. RU and BBU). The vEPC forwards the service traffic to an external DC site (that is usually on premises of the factory) where the factory applications run.

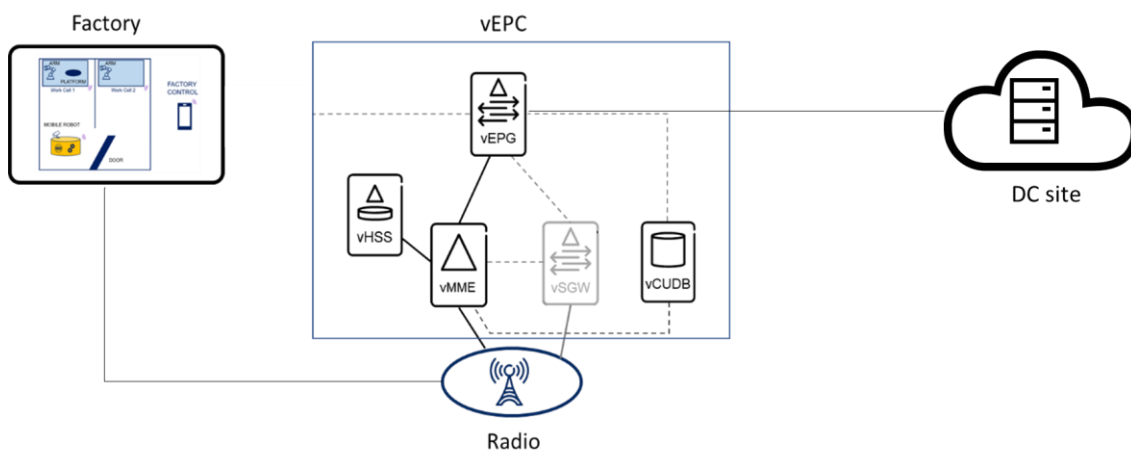


FIGURE 10: eINDUSTRY SERVICE TOPOLOGY

The EPC components virtualized as VNF in the vEPC and used for the specific use case are:

- vEPG (Evolved Packet Gateway)
- vMME (Mobility Management Entity)
- vCUDB (Centralized User Database)
- vHSS (Home Subscriber Server)

The use case service VNFs' resources used for the EPC components can be found in Table 11.

TABLE 11: eINDUSTRY VNFs' RESOURCES

Service 6: eIndustry						
Node Type	VM	No. of VMs	Virtual CPUs	Memory [GB]	Disk [GB]	Licenses
vEPG	vRP	2	2	6	40	1
	vSFO-GSC	1	12	8	40	1
	vSFO-CP/UP	3	6	10	40	1
vMME	FSB	1	2	5	160	1
	NCB	2	2	2	0.001	1
	GPB	2	6	12	0.001	1
vCUDB	SC	2	2	6	80	1
	PL	2	2	6	20	1
vHSS	SC	2	2	12	200	1
	PL	2	2	12	0.001	1

According to [12], latency and required bandwidth for this type of use case are those shown in Table 12:

TABLE 12: eINDUSTRY VNFs' LINKS ASSETS

Service 6: eIndustry		
Link	BW [Mbps]	Latency [ms]
int_a	100	10
int_b	100	10
int_c	100	10
int_d	100	10

2.2.2 Infrastructure cost modelling

The objective of this section is to provide guidance to estimate the infrastructure cost for the deployment of the 5GT vertical use cases. This guidance is especially useful for Cloud providers to estimate the budget of their platform and the benefits they may obtain from the allocation of their resources.

To better understand this infrastructure cost modelling, we introduce in the following, some background regarding the physical infrastructure architecture.

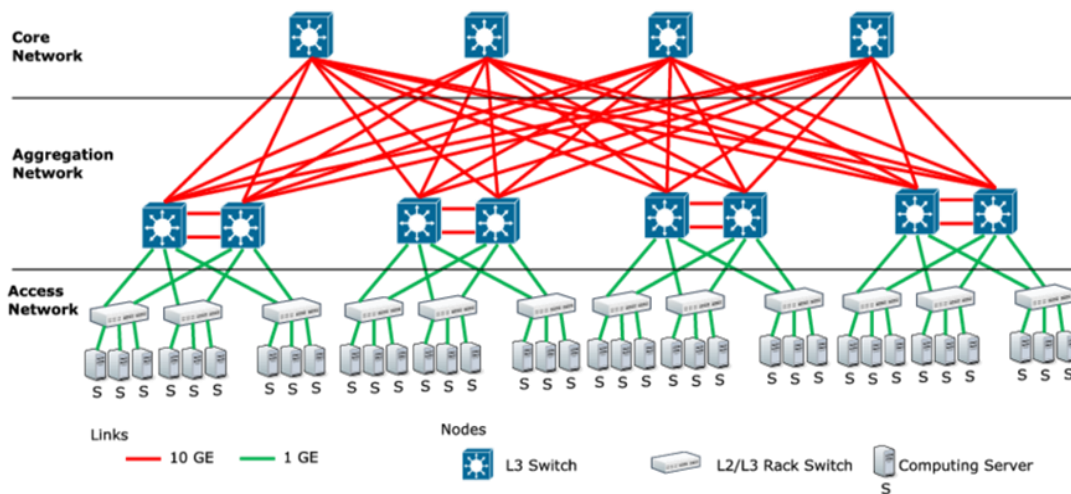


FIGURE 11: TRADITIONAL THREE-TIER DATACENTRE TOPOLOGY



FIGURE 12: STANDARD SPINE-AND-LEAF TOPOLOGY [15]

The focus will be on the interconnection, nodes, racks, and power.

- Interconnections:** The topology of datacentres has evolved over time. The traditional design was made over three-tier architecture namely; the access, aggregation, and core layer as shown in Figure 11. The servers are segmented into pods according to their location. The access layer consists of layer two switches where the servers are connected. These access switches are connected to layer 2/layer 3 switches, which are connected to layer 3 aggregation routers. A Spanning Tree Protocol (STP) is used to build a loop-free topology for the layer 2 part of network between these two layers. This architecture has been considered for many years before moving to the Spine-and-Leaf topology. Indeed, this classical architecture has reached its limitations typically from running the Spanning Tree Protocol, which uses several ports to protect against layer 2 loops, hence leave the network with half capacity. Besides that, it increases the complexity of troubleshooting. The Spine-and-Leaf architecture is two-tier architecture known as the Leaf layer and Spine layer. The Leaf layer represents the lower-tier switches, while the Spine layer consists of the top-tier switches. The Spine and Leaf switches are connected in a full-mesh topology (Figure 12). The Leaf layer represents the access switches that

are connected to devices such as servers. The Spine layer is the background network and it is responsible for interconnecting all the Leaf switches. The Spine-and-Leaf can be either layer 2 or layer 3. As a layer 2, all the connections are active/active. Instead of using the STP, we use other protocols such as Transparent Interconnection of Lots of Links (TrILL) or Shortest Path Bridging (SPB). When used as layer 3, the architecture is much simpler. We use a protocol such as Equal Cost Multi Path (ECMP) to distribute traffic across links. It is also possible to use the Open Shortest Path First (OSPF) or Border Gateway Protocol (BGP) to distribute routes across the data-center fabric. It should be noted however, that layer 3 Spine-and-Leaf deployments are not currently supported by the Openstack platform. Indeed, Openstack networking is typically layer 2 based (VLAN, VxLAN). The Leaf-Spine is a two-layer data-center network topology composed of Leaf and spine switches. The Leaf switches are connected to the servers and storage, while the Spine switches are connected to the Leaf switches. The leaf and Spine switches are connected in a mesh topology, forming the access layer that delivers network connection points for servers.

- *Compute Node*: A compute node is a hardware machine that hosts multiple instances of a virtual machine (VM). It accounts for and provides processing, memory, network, and storage that VM instances consume. In Openstack for instance, each compute node runs a hypervisor such as KVM to deploy and run VMs.
- *Network Node*: A network node is a physical server providing a network function
- *Storage node*: A storage node is usually a physical server with multiple hard-disk drives or solid-state drives. Multiple storage nodes can be clustered together to create a software-controlled storage pool.
- *Power*: It is the peak energy consumed by the components of a device or node.
- *Rack*: A rack equipment [16], also referred as a frame, is a physical device in a shape of (metallic) shelf, into which we may install all kind of nodes (including, networking, compute, and storage nodes) in a mechanical, safe, and stable manner (see Figure 13). A rack may also include additional features such as integrated cooling, integrated power, and seismic hardening. In this case, a management controller should also be included to configure, monitor, and control the cooling and power. Used to build the NFVI hardware infrastructure, a rack may hold several nodes; hence the relationship is not necessarily one-to-one. Geographically speaking, racks are physically deployed within an NFVI-PoP at fixed location in row and aisle manner. In most cases, these racks will not change their location, in order to provide a stable reference point within the NFVI-PoP. The NFVI-PoP is a set of NFVI nodes, each of which is mapped to one or more racks, each containing compute, storage, and network nodes. There are also situations where multiple NFVI nodes are mapped to a single rack of equipment. These NFVI nodes, will share the rack level infrastructure (power, cooling, and hardware management).

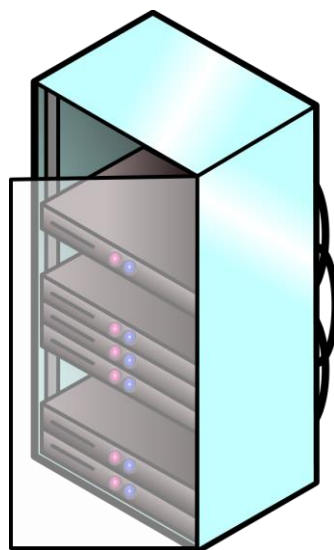


FIGURE 13: A SIMPLE RACK WITH SEVERAL NODES

We consider some entry variables that may vary according to countries, material providers, and so on. These variables are presented in Table 13.

TABLE 13: ENTRY VARIABLES THAT INFLUENCE THE COSTS

Entry	Symbol	Value
Cost per month for a rack	x_1	\$1200
Number of rack units available	x_2	42
Cost per rack unit per month	x_3	\$28.571
Raw cost for power per kW/h	x_4	\$0.100
Power usage effectiveness ratio (PUE) of the data centre	x_5	2.5
Cost for power per kW/h	x_6	\$0.250
Depreciation time for ICT hardware in months	x_7	36
Interest rates for leasing or financing hardware	x_8	6%
Hours per month	x_9	730.5
Business hours per month	x_{10}	146
Cost per hour of an OpenStack engineer	x_{11}	\$105.86
Desirable redundancy (N+?) for the clusters	x_{12}	1
Margin (for cost recovery or profit)	x_{13}	20%
Currency conversion to USD	x_{14}	1.2
VMWare licence costs/CPU sockets	x_{15}	\$1200.00
Falcon Dedicated Ratio	x_{16}	0.3
Monthly Internet Access Cost - 1Gbps	x_{17}	\$1440.00

2.2.2.1 How to compute the bill of materials

2.2.2.1.1 Network ports

Through this section, we explain how to calculate several prices for the network ports. For this, we assume a spine/ leaf architecture using MLAG, VxLAN, and ECMP.

- **SPINE SWITCHES**

For the spine switches, we consider the following entries presented in Table 14. We define symbols for the convenience of presentation.

TABLE 14: ENTRIES INFLUENCING THE SPINE SWITCH PORTS PRICE

Entry	Symbol
Cost for a spine switch	a_1
Number of ports per switch	a_2
Aggregation ports per switch	a_3
Number of ports per switch used for MLAG	a_4
Rack units per switch	a_5
Power consumption (in Watts) per switch	a_6
Number of spine switches	a_7

We are interested into several costs that we present in Table 15:

TABLE 15: COST EQUATIONS FOR SPINE SWITCHES

Entry	Formula
Equation 1: Number of spine ports	$A_1 = (a_2 - a_4) \times a_7$
Equation 2: Cost for the spine switches per month	$A_2 = -VPM \left(\frac{x_8}{12}; x_7; (a_7 \times a_1) \right)$
Equation 3: Cost for the rack units per month	$A_3 = (a_7 \times a_5) \times x_3$
Equation 4: Cost for the power per month	$A_4 = ((a_7 \times a_6)/1000) \times x_6 \times x_{10}$
Equation 5: Total cost per month	$A_5 = \sum_{i=2}^4 A_i$
Equation 6: Cost per core port per month	$A_6 = A_5 / A_1$

- **LEAF SWITCHES**

For the leaf switches, we consider the following entries in Table 16, in addition to the previous ones:

TABLE 16: ENTRIES INFLUENCING THE LEAF SWITCH PORTS PRICE

Entry	Variable
Cost for a leaf switch	b_1
Number of ports per switch	b_2
Aggregation ports per switch	b_3
Number of ports per switch used for MLAG	b_4
Number of ports per switch used to connect to spine	b_5
Rack units per switch	b_6
Power consumption (in Watts) per switch	b_7
Number of leaf switches	b_8

The costs, as the previous ones for the Spine switches are presented in Table 17 :

TABLE 17: COST EQUATIONS FOR LEAF SWITCHES

Entry	Formula
Equation 7: Max number of leaf switches that can be connected to the spine	$B_1 = A_1/b_5$
Equation 8: Max number of leaf ports in this configuration	$B_2 = B_1 \times b_2$
Equation 9: Number of leaf ports	$B_3 = b_8 \times b_2$
Equation 10: Cost for the leaf switches per month	$B_4 = -VPM\left(\frac{x_8}{12}; x_7; (b_8 \times b_1)\right)$
Equation 11: Cost for the rack units per month	$B_5 = (b_8 \times b_6) \times x_3$
Equation 12: Cost for the power per month	$B_6 = ((b_8 \times b_7)/1000) \times x_6 \times x_{10}$
Equation 13: Cost for the spine ports per month	$B_7 = b_5 \times A_6$
Equation 14: Total cost per month	$B_8 = \sum_{i=4}^7 B_i$
Equation 15: Cost per core port per month	$B_9 = B_8/B_3$

- MANAGEMENT SWITCHES

We need to define other entries for the management switches in addition to the ones presented in Table 14 and Table 16. These additional entries are defined in Table 18.

TABLE 18: ENTRIES INFLUENCING THE MANAGEMENT SWITCH PORTS PRICE

Entry	Variable
Cost for a management switch	c_1
Number of ports per switch	c_2
Number of aggregation ports per switch	c_3
Number of ports per switch used for MLAG	c_4
Number of ports per switch used to connect to spine	c_5
Rack units per switch	c_6
Power consumption (in Watts) per switch	c_7
Number of management switches	c_8

Table 19 shows how to calculate similar costs as for the spine and leaf switches.

TABLE 19: COST EQUATIONS FOR MANAGEMENT SWITCHES

Entry	Formula
Equation 16: Number of management ports	$C_1 = c_8/c_2$
Equation 17: Cost for the management switches per month	$C_2 = -VPM\left(\frac{x_8}{12}; x_7; (c_8 \times c_1)\right)$
Equation 18: Cost for the rack units per month	$C_3 = (c_8 \times c_6) \times x_3$
Equation 19: Cost for the power per month	$C_4 = ((c_8 \times c_7)/1000) \times x_6 \times x_{10}$

Equation 20: Cost for the leaf ports per month	$C_5 = c_5 \times B_9$
Equation 21: Total cost per month	$C_6 = \sum_{i=2}^5 C_i$
Equation 22: Cost per management port per month	$C_7 = C_6 / C_1$

In the annex section A.1.1, we provide examples of cost calculation for three commercial switch types: spine switch, leaf switch and management switch.

2.2.2.1.2 Network node

Like network ports, the network nodes influence the Total Cost Ownership (TCO) of the infrastructure. Here also, we introduce the elements that make different costs of the network nodes vary that in turn affect the TCO. These elements are presented in Table 20. For the convenience of presentation, we used symbols to describe these entries.

TABLE 20: ENTRIES INFLUENCING THE NETWORK NODE COSTS PER MONTH

Entry	Symbol
Cost for a network node	d_1
Rack units per network node	d_2
Power consumption (in Watts/h) per node	d_3
Number of leaf network ports per node	d_4
Number of management network ports per node	d_5

We compute costs for the network nodes, according the entries presented in Table 20. These costs, as well as the formula to estimate them are presented in Table 21.

TABLE 21: COST EQUATIONS FOR NETWORK NODES

Entry	Calculation
Equation 23: Cost per month for the hardware	$D_1 = -VPM(x_8/12; x_7; d_1)$
Equation 24: Cost per month for the rack space	$D_2 = d_2 \times x_3$
Equation 25: Cost per month for the power	$D_3 = (d_3/1000) \times x_6 \times x_9$
Equation 26: Cost per month for the leaf network ports	$D_4 = d_4 \times b_9$
Equation 27: Cost per month for the management network ports	$D_5 = d_5 \times C_7$
Equation 28: Cost per month per controller node	$D_6 = \sum_{i=1}^5 D_i$

In the annex section A.1.2, we provide an example of calculation of a network node.

2.2.2.1.3 Controller

In addition to the entries defined in the previous tables, we introduce others necessary to calculate the different costs per month for a controller node. We present these new entries in Table 22.

TABLE 22: ENTRIES INFLUENCING THE CONTROLLER NODE COSTS PER MONTH

Entry	Symbol
Cost for a network node	d_1
Rack units per network node	d_2
Power consumption (in Watts/h) per node	d_3
Number of leaf network ports per node	d_4
Number of management network ports per node	d_5

The pricing of a controller node comprises the formulas in Table 23 to estimate multiple costs per month under the form of Equation 29 to Equation 34.

TABLE 23: HOW TO COMPUTE MONTHLY COSTS FOR CONTROLLER NODE

Entry	Calculation
Equation 29: Cost per month for the hardware	$E_1 = -VPM(x_8/12; x_7; e_1)$
Equation 30: Cost per month for the rack space	$E_2 = e_2 \times x_3$
Equation 31: Cost per month for the power	$E_3 = (e_3/1000) \times x_6 \times x_9$
Equation 32: Cost per month for the leaf network ports	$E_4 = e_4 \times b_9$
Equation 33: Cost per month for the management network ports	$E_5 = e_5 \times C_7$
Equation 34: Cost per month per controller node	$E_6 = \sum_{i=1}^5 E_i$

An example of numerical application of these formulas is given in the annex A.1.3.

2.2.2.1.4 Compute

We assume that we use an open source hypervisor such as KVM. We present in Table 24, the additional entries that affect the costs for the compute nodes.

TABLE 24: ENTRIES INFLUENCING THE COMPUTE NODE COSTS PER MONTH

Entry	Symbol
Cost for a compute node	f_1
Rack units per compute node	f_2
Power consumption (in Watts/hour) per node	f_3
Number of CPU sockets per node	f_4
Number of cores per CPU	f_5
Number of threads per CPU core	f_6
CPU overcommit ratio	f_7
RAM per node (in GB)	f_8
RAM reserved for the hypervisor (base OS)	f_9
RAM overcommit ratio	f_{10}
Number of leaf network ports per node	f_{11}
Number of management network ports per node	f_{12}
CPU weight for pricing	f_{13}
RAM weight for pricing	f_{14}

Table 25 presents equations for monthly costs about compute nodes.

TABLE 25: HOW TO COMPUTE MONTHLY COSTS FOR A COMPUTE NODE

Entry	Calculation
Equation 35: Cost per month for the hardware	$F_1 = -VPM(x_8/12; x_7; f_1)$
Equation 36: Cost per month for the rack space	$F_2 = f_2 \times x_3$
Equation 37: Cost per month for the leaf network ports	$F_3 = f_{11} \times B_9$
Equation 38: Cost per month for the management network ports	$F_4 = f_{12} \times C_7$
Equation 39: Cost per month for the power	$F_5 = (f_3/1000) \times x_6 \times x_9$
Equation 40: Cost per month per compute node	$F_6 = \sum_{i=1}^5 F_i$
Equation 41: vCPUs per node	$F_7 = f_4 \times f_5 \times f_6 \times f_7$
Equation 42: RAM (in GB) per node	$F_8 = f_8 - f_9$
Equation 43: Cost per vCPU per month	$F_9 = (F_6 \times f_{13})/F_7$
Equation 44: Cost per vCPU per hour	$F_{10} = F_9/x_9$
Equation 45: Cost per GB of RAM per month	$F_{11} = (F_6 \times f_{14})/F_8$
Equation 46: Cost per GB of RAM per hour	$F_{12} = F_{11}/x_9$

An example of pricing of a compute node is presented in the appendix A.1.4.

2.2.2.1.5 Block storage HDD

We assume the use of a distributed storage system such as Ceph. Table 26 presents the entries impacting the block storage.

TABLE 26: ENTRIES INFLUENCING THE BLOCK STORAGE HDD COSTS PER MONTH

Entry	Symbol
Cost for a block storage node	g_1
Rack units per block storage node	g_2
Power consumption (in Watts/hour) per node	g_3
Number of leaf ports per node	g_4
Number of management ports per node	g_5
Number of OSDs per server	g_6
Size of OSD disks (in GB)	g_7
Storage overcommit ratio	g_8
IOPS per OSD disk for random reads	g_9
IOPS per OSD disk for random writes	g_{10}
IOPS efficiency of storage software for reads	g_{11}
IOPS efficiency of storage software for writes	g_{12}
Number of nodes in the storage cluster	g_{13}
Number of replicas per object	g_{14}
Number of acceptable node failures	g_{15}

Equation 47 to Equation 61 represent how to compute multiple costs per month, but also other values. We present these equations in Table 27. In the appendix A.1.5 we show a bill of material for block storage HDD.

TABLE 27: HOW TO COMPUTE MONTHLY COSTS FOR THE BLOCK STORAGE HDD

Entry	Calculation
Equation 47: Cost per month for the hardware	$G_1 = -VPM(x_8/12; x_7; g_{13} \times g_1)$
Equation 48: Cost per month for the rack space	$G_2 = g_{13} \times g_2 \times x_3$
Equation 49: Cost per month for the leaf network ports	$G_3 = g_{13} \times g_4 \times B_9$
Equation 50: Cost per month for the management network ports	$G_4 = g_{13} \times g_5 \times C_7$
Equation 51: Cost per month for the power	$G_5 = g_{13} \times (g_3/1000) \times x_6 \times x_9$
Equation 52: Cost per month for the storage cluster	$G_6 = \sum_{i=1}^5 G_i$
Equation 53: OSDs in the cluster	$G_7 = g_{13} \times g_6$
Equation 54: Total raw capacity of the cluster (in GB)	$G_8 = G_7 \times g_7$
Equation 55: Usable capacity ratio	$G_9 = \frac{(g_{13} - g_{15})}{g_{14}}$
Equation 56: Usable capacity of the cluster (in GB)	$G_{10} = G_9 \times G_8$
Equation 57: Overcommitted usable capacity of the cluster (in GB)	$G_{11} = G_{10} + (G_{10} \times g_8)$
Equation 58: Maximum read IOPS for the cluster (after cache is exhausted)	$G_{12} = g_{13} \times g_6 \times g_9 \times g_{11}$
Equation 59: Maximum write IOPS for the cluster (after cache is exhausted)	$G_{13} = g_{13} \times g_6 \times g_{10} \times g_{12}$
Equation 60: Cost per GB per month	$G_{14} = G_6 / G_{11}$
Equation 61: Cost per GB per hour	$G_{15} = G_{14} / x_9$

2.2.2.1.6 Block storage SSD

Assumes here too, the use of a distributed storage system like Ceph. Additional entries are described in Table 29, which presents Equation 62 to Equation 76 that we use to estimate the costs per month for one block storage SSD. An example of Bill material is presented in A.1.5.

TABLE 28: ENTRIES INFLUENCING THE BLOCK STORAGE SSD COSTS PER MONTH

Entry	Symbol
Cost for a block storage node	h_1
Rack units per block storage node	h_2
Power consumption (in Watts/hour) per node	h_3
Number of leaf ports per node	h_4
Number of management ports per node	h_5
Number of OSDs per server	h_6
Size of OSD disks (in GB)	h_7
Storage overcommit ratio	h_8

IOPS per OSD disk for random reads	h_9
IOPS per OSD disk for Equation 61c random writes	h_{10}
IOPS efficiency of storage software for reads	h_{11}
IOPS efficiency of storage software for writes	h_{12}
Number of nodes in the storage cluster	h_{13}
Number of replicas per object	h_{14}
Number of acceptable node failures	h_{15}

TABLE 29: HOW TO COMPUTE MONTHLY COSTS FOR THE BLOCK STORAGE SSD

Entry	Calculation
Equation 62: Cost per month for the hardware	$H_1 = -VPM(x_8/12; x_7; h_{13} \times h_1)$
Equation 63: Cost per month for the rack space	$H_2 = h_{13} \times h_2 \times x_3$
Equation 64: Cost per month for the leaf network ports	$H_3 = h_{13} \times h_4 \times B_9$
Equation 65: Cost per month for the management network ports	$H_4 = h_{13} \times h_5 \times C_7$
Equation 66: Cost per month for the power	$H_5 = h_{13} \times (h_3/1000) \times x_6 \times x_9$
Equation 67: Cost per month for the storage cluster	$H_6 = \sum_{i=1}^5 H_i$
Equation 68: OSDs in the cluster	$H_7 = h_{13} \times h_6$
Equation 69: Total raw capacity of the cluster (in GB)	$H_8 = H_7 \times h_7$
Equation 70: Usable capacity ratio	$H_9 = \frac{(h_{13} - h_{15})}{h_{14}}$
Equation 71: Usable capacity of the cluster (in GB)	$H_{10} = H_9 \times H_8$
Equation 72: Overcommitted usable capacity of the cluster (in GB)	$H_{11} = H_{10} + (H_{10} \times h_8)$
Equation 73: Maximum read IOPS for the cluster (after cache is exhausted)	$H_{12} = h_{13} \times h_6 \times h_9 \times h_{11}$
Equation 74: Maximum write IOPS for the cluster (after cache is exhausted)	$H_{13} = h_{13} \times h_6 \times h_{10} \times h_{12}$
Equation 75: Cost per GB per month	$H_{14} = H_6/H_{11}$
Equation 76: Cost per GB per hour	$H_{15} = H_{14}/x_9$

2.2.2.1.7 Object storage low density

Table 30 presents the entries that are influencing the different costs for the storage with low density. We assume the use of a distributed storage system like Ceph. The costs are computed using these entries as shown in Equation 77 to Equation 91 presented in Table 31. In appendix A.1.7, we present the prices for different components that compose the object storage with low density.

TABLE 30: ENTRIES INFLUENCING THE OBJECT STORAGE LOW DENSITY COSTS PER MONTH

Entry	Symbol
Cost for a block storage node	i_1
Rack units per block storage node	i_2

Power consumption (in Watts/hour) per node	i_3
Number of leaf ports per node	i_4
Number of management ports per node	i_5
Number of OSDs per server	i_6
Size of OSD disks (in GB)	i_7
Storage overcommit ratio	i_8
IOPS per OSD disk for random reads	i_9
IOPS per OSD disk for random writes	i_{10}
IOPS efficiency of storage software for reads	i_{11}
IOPS efficiency of storage software for writes	i_{12}
Number of nodes in the storage cluster	i_{13}
Number of replicas per object	i_{14}
Number of acceptable node failures	i_{15}

TABLE 31: HOW TO COMPUTE MONTHLY COSTS FOR THE OBJECT STORAGE LOW DENSITY

Entry	Calculation
Equation 77: Cost per month for the hardware	$I_1 = -VPM(x_8/12; x_7; i_{13} \times i_1)$
Equation 78: Cost per month for the rack space	$I_2 = i_{13} \times i_2 \times x_3$
Equation 79: Cost per month for the leaf network ports	$I_3 = i_{13} \times i_4 \times B_9$
Equation 80: Cost per month for the management network ports	$I_4 = i_{13} \times i_5 \times C_7$
Equation 81: Cost per month for the power	$I_5 = i_{13} \times (i_3/1000) \times x_6 \times x_9$
Equation 82: Cost per month for the storage cluster	$I_6 = \sum_{i=1}^5 I_i$
Equation 83: OSDs in the cluster	$I_7 = i_{13} \times i_6$
Equation 84: Total raw capacity of the cluster (in GB)	$I_8 = I_7 \times i_7$
Equation 85: Usable capacity ratio	$I_9 = \frac{(i_{13} - i_{15})}{i_{14}}$
Equation 86: Usable capacity of the cluster (in GB)	$I_{10} = I_9 \times I_8$
Equation 87: Overcommitted usable capacity of the cluster (in GB)	$I_{11} = I_{10} + (I_{10} \times i_8)$
Equation 88: Maximum read IOPS for the cluster (after cache is exhausted)	$I_{12} = i_{13} \times i_6 \times i_9 \times i_{11}$
Equation 89: Maximum write IOPS for the cluster (after cache is exhausted)	$I_{13} = i_{13} \times i_6 \times i_{10} \times i_{12}$
Equation 90: Cost per GB per month	$I_{14} = I_6 / I_{11}$
Equation 91: Cost per GB per hour	$I_{15} = I_{14} / x_9$

2.2.2.1.8 Object storage high density

Similarly as for the object storage with low storage, we present herein how to compute the costs for object storage with high density. We assume the use of a distributed storage system like Swift or Ceph. The entries for formulations are presented in Table

32, and the formulas in Table 33. In annex 0, the bill of materials for the components that devise the object storage with high density.

TABLE 32: ENTRIES INFLUENCING THE OBJECT STORAGE HIGH DENSITY COSTS PER MONTH

Entry	Symbol
Cost for an object storage node	j_1
Rack units per object storage node	j_2
Power consumption (in Watts/hour) per node	j_3
Number of leaf ports per node	j_4
Number of management ports per node	j_5
Number of OSDs per server	j_6
Size of OSD disks (in GB)	j_7
Storage overcommit ratio	j_8
IOPS per disk for random reads	j_9
IOPS per disk for random writes	j_{10}
IOPS efficiency of storage software for reads	j_{11}
IOPS efficiency of storage software for writes	j_{12}
Number of nodes in the storage cluster	j_{13}
Number of replicas per object	j_{14}
Number of acceptable node failures	j_{15}

TABLE 33: HOW TO COMPUTE MONTHLY COSTS FOR THE OBJECT STORAGE HIGH DENSITY

Entry	Calculation
Equation 92: Cost per month for the hardware	$J_1 = -VPM(x_8/12; x_7; j_{13} \times j_1)$
Equation 93: Cost per month for the rack space	$J_2 = j_{13} \times j_2 \times x_3$
Equation 94: Cost per month for the leaf network ports	$J_3 = j_{13} \times j_4 \times B_9$
Equation 95: Cost per month for the management network ports	$J_4 = j_{13} \times j_5 \times C_7$
Equation 96: Cost per month for the power	$J_5 = j_{13} \times (j_3/1000) \times x_6 \times x_9$
Equation 97: Cost per month for the storage cluster	$J_6 = \sum_{i=1}^5 J_i$
Equation 98: OSDs in the cluster	$J_7 = j_{13} \times j_6$
Equation 99: Total raw capacity of the cluster (in GB)	$J_8 = J_7 \times j_7$
Equation 100: Usable capacity ratio	$J_9 = \frac{(j_{13} - j_{15})}{j_{14}}$
Equation 101: Usable capacity of the cluster (in GB)	$J_{10} = J_9 \times J_8$
Equation 102: Overcommitted usable capacity of the cluster (in GB)	$J_{11} = J_{10} + (J_{10} \times j_8)$

Equation 103: Maximum read IOPS for the cluster (after cache is exhausted)	$J_{12} = j_{13} \times j_6 \times j_9 \times j_{11}$
Equation 104: Maximum write IOPS for the cluster (after cache is exhausted)	$J_{13} = j_{13} \times j_6 \times j_{10} \times j_{12}$
Equation 105: Cost per GB per month	$J_{14} = J_6 / J_{11}$
Equation 106: Cost per GB per hour	$J_{15} = J_{14} / x_9$

2.2.2.1.9 Object storage archiving

We assume the use of a distributed storage system like Swift or Ceph. Entries, Equations, and Bill of Materials are presented respectively in Table 34 and Table 35.

TABLE 34: ENTRIES INFLUENCING THE OBJECT STORAGE ARCHIVING COSTS PER MONTH

Entry	Symbol
Cost for an object storage node	k_1
Rack units per object storage node	k_2
Power consumption (in Watts/hour) per node	k_3
Number of leaf ports per node	k_4
Number of management ports per node	k_5
Number of OSDs per server	k_6
Size of OSD disks (in GB)	k_7
Storage overcommit ratio	k_8
IOPS per disk for random reads	k_9
IOPS per disk for random writes	k_{10}
IOPS efficiency of storage software for reads	k_{11}
IOPS efficiency of storage software for writes	k_{12}
Number of nodes in the storage cluster	k_{13}
Erasur coding overhead	k_{14}
Number of acceptable node failures	k_{15}

TABLE 35: HOW TO COMPUTE MONTHLY COSTS FOR THE OBJECT STORAGE ARCHIVING

Entry	Calculation
Equation 107: Cost per month for the hardware	$K_1 = -VPM(x_8/12; x_7; k_{13} \times k_1)$
Equation 108: Cost per month for the rack space	$K_2 = k_{13} \times k_2 \times x_3$
Equation 109: Cost per month for the leaf network ports	$K_3 = k_{13} \times k_4 \times B_9$
Equation 110: Cost per month for the management network ports	$K_4 = k_{13} \times k_5 \times C_7$
Equation 111: Cost per month for the power	$K_5 = k_{13} \times (k_3/1000) \times x_6 \times x_9$
Equation 112: Cost per month for the storage cluster	$K_6 = \sum_{i=1}^5 K_i$
Equation 113: Disks in the cluster	$K_7 = k_{13} \times k_6$
Equation 114: Total raw capacity of the cluster (in GB)	$K_8 = K_7 \times k_7$
Equation 115: Usable capacity ratio	$K_9 = \frac{(k_{13} - k_{15})}{k_{14}}$

Equation 116: Usable capacity of the cluster (in GB)	$K_{10} = K_9 \times K_8$
Equation 117: Overcommitted usable capacity of the cluster (in GB)	$K_{11} = K_{10} + (K_{10} \times k_8)$
Equation 118: Maximum read IOPS for the cluster (after cache is exhausted)	$K_{12} = k_{13} \times k_6 \times k_9 \times k_{11}$
Equation 119: Maximum write IOPS for the cluster (after cache is exhausted)	$K_{13} = k_{13} \times k_6 \times k_{10} \times k_{12}$
Equation 120: Cost per GB per month	$K_{14} = K_6 / K_{11}$
Equation 121: Cost per GB per hour	$K_{15} = K_{14} / x_9$

2.2.2.1.10 Staff

The cost per month for the Staff is defined according to how the teams are sized, the salary for an employee, and the number of working hours. For example, if a team is composed of 2 Full Time Equivalent (FTE) working for 8 x 5, then the minimum cost would be \$30 911. While, for 6 persons FTE for 24 hours over 7 days a week, the minimum cost should be \$92 733.

2.2.2.2 Total cost of ownership

In order to compute the financial profitability from the Cloud, infrastructure providers need to add the ownership cost. This cost is known as the Total Cost Ownership (TCO). The TCO is a financial estimation of direct and indirect cost of the platform which includes the material cost, software, and licences costs, consumption, and staff. The TCO is needed in order to compute the Return On Investment (ROI). In what follows we present how to estimate TCO for three types of Cloud: Small Private Cloud with 8 hours business working during 5 days a week (8*5), Medium Private Cloud (8 * 5), and Large private Cloud (24 * 7). For the hardware cost per month (HW_Cost), the formula is given by Equation 122.

$$HW_{Cost} = -VPM(x_8/12; x_7; HW_{UnitPrice}) \times Quantity_{HW}$$

EQUATION 122: HARDWARE COST PER MONTH

Where $HW_{UnitPrice}$ is the unit price of the hardware and $Quantity_{HW}$ is the quantity of used hardware. HW concerns all the nodes including the spine, leaf, and management switches, controller, network, and compute nodes, Block storage nodes (HDD and/or SDD), object storage nodes with low and/or high density. The number of rack units ($TotalNumberofRackUnits_{HW}$) used per HW or node is obtained in Equation 123. The $Number_of_Rack_Units_per_Hardware$ represents the number of rack units used for one type of HW .

$$\begin{aligned} TotalNumberofRackUnits_{HW} \\ = Quantity_{HW} \times Number_of_Rack_Units_per_Hardware \end{aligned}$$

EQUATION 123: TOTAL NUMBER OF RACK UNITS

The cost of the rack units (RU_Cost) per month per HW is estimated Equation 124. The $TotalNumberofRackUnits_{HW}$ is the result of Equation 123. For reminder, x_3 is defined in Table 13 as the cost for a rack unit (per month).

$$RU_{Cost} = TotalNumberofRackUnits_{HW} \times x_3$$

EQUATION 124: RACK UNITS COST

We provide also in the following estimations for respectively, the power consumption in Watts/hour, the power cost per month, the number of 10Gbps ports, number of 1Gbps ports, and the sum of the hardware cost, rack units cost, and power cost per month. Equation 125 represents the power consumption in Watts per hour per node. The *Power_Consumption_per_hardware* is the unit consumption per node.

$$Power \left(\frac{Watts}{hour} \right) = Quantity_{HW} \times Power_Consumption_per_hardware$$

EQUATION 125: POWER CONSUMPTION

The power cost per month is estimated by Equation 126, wherein $Power \left(\frac{Watts}{hour} \right)$ is the result of Equation 125, and x_6, x_9 describe respectively, the power cost for KW/h, and the number of hours per month. These two variables are defined in Table 13.

$$Power_{Cost} = \frac{Power \left(\frac{Watts}{hour} \right)}{1000} \times x_6 \times x_9$$

EQUATION 126: POWER COST

Equation 127 counts the number of ports with 10Gbps. The *Number_of_leaf_network_ports_per_node* is the number of leaf ports that one node may provide.

$$\begin{aligned} Number_of_10Gbps_ports \\ = Quantity_{HW} \times Number_of_leaf_network_ports_per_node \end{aligned}$$

EQUATION 127: NUMBER OF 10GBPS PORTS

Similarly, Equation 128 counts the number of ports with 1Gbps that concern the management ports. *Number_of_management_network_ports_per_node* is the number of management network ports that are provided by the node.

$$\begin{aligned} Number\ of\ 1Gbps\ ports \\ = Quantity_{HW} \\ \times Number_of_management_network_ports_per_node \end{aligned}$$

EQUATION 128: NUMBER OF 1GBPS PORTS

The cost of the hardware spent per month is obtained in Equation 129 which is the sum of costs obtained in Equation 122, Equation 124 and Equation 126. During 3 years, this cost is estimated by Equation 130.

$$Cost_{per\ month} = HW_{Cost} + RU_{Cost} + Power_{cost}$$

EQUATION 129: COST PER MONTH

$$Cost_{per\ 3\ years} = Cost_{per\ month} \times 36$$

EQUATION 130: COST PER 3 YEARS

For the Staff, the cost per month is computed, in Equation 131, as follows:

$$Staff_{Cost} = Quantity_{Staff} \times x_{10} \times x_{11}$$

EQUATION 131: STAFF COST

x_{10} and x_{11} are two variables defined in Table 13 as respectively, the business hours per month, and the cost per hour for an OpenStack engineer. For the VMWare licence, the cost is given by the Equation 132, where $Quantity_{VMWareLicence}$ is the number of licences, and x_{15} is a variable defined in Table 13, it represents the VMWare licence costs/CPU sockets:

$$VMWareLicence_{cost} = Quantity_{VMWareLicence} \times x_{15}$$

EQUATION 132: VMWARE AS A LICENSE

To sum up, the Total Cost of Ownership (TCO) is given by:

$$TCO = \sum_{per\ HW} HW_{Cost} + Staff_{Cost} + VMWareLicence_{cost}$$

EQUATION 133: TCO

Where $\sum_{per\ HW} HW_{cost}$ is the cost sum of spine switches, leaf switches, management switches, controller nodes, network nodes, compute nodes, and storage nodes. Each cost of these nodes is obtained in Equation 122. The $Staff_{Cost}$ and $VMWareLicence_{cost}$ are respectively, obtained from Equation 131 and Equation 132.

For the provided capacity, we have computed several entries as follow:

$$NumberPorts_{40Gbps} = Quantity_{Spine\ Switches} \times a_2$$

EQUATION 134: NUMBER OF PORTS OFFERING 40GBPS

Where $Quantity_{Spine\ Switches}$ is the number of used spine switches, and a_2 is defined in Table 14, as the number of ports per the spine switch.

$$NumberPorts_{10Gbps} = Quantity_{Leaf\ Switches} \times b_2$$

EQUATION 135: NUMBER OF PORTS OFFERING 10GBPS

Where $Quantity_{Leaf\ Switches}$ is the number of leaf switches that have been used, and b_2 is an entry defined in Table 16, as the number of ports per leaf switch.

$$NumberPorts_{1Gbps} = Quantity_{Management\ Switches} \times c_2$$

EQUATION 136: NUMBER OF PORTS OFFERING 1GBPS

Where $Quantity_{Management\ Switches}$ represents the number of used management switches, and c_2 is defined in Table 18 as the number of ports per management switch.

$$NumberVCPUs = (Quantity_{Compute\ Nodes} - x_{12}) \times F_7$$

EQUATION 137: NUMBER OF VCPUS

$Quantity_{Compute\ Nodes}$ represents the number of used compute nodes, x_{12} is the desirable redundancy for the clusters as defined in Table 13, and F_7 defined in Table 25, as the number of vCPUs per node, which is obtained by Equation 137.

$$Size_{RAM} = \left(Quantity_{Compute\ Nodes} - x_{12} \right) \times F_8$$

EQUATION 138: RAM SIZE

Where F_8 is the RAM (in GB) per node obtained by Equation 138 as defined in Table 25.

$$Size_{BlockStorageHDD} = \left(\frac{\left(\left(Quantity_{Block\ Storage\ Nodes\ (HDD)} - x_{12} \right) \times g_6 \times g_7 \right)}{g_{14}} \right) \times (1 + g_8)$$

EQUATION 139: BLOCK STORAGE HDD

Where $Quantity_{Block\ Storage\ Nodes\ (HDD)}$ is the number of block storage HDD nodes, g_6 , g_7 , and g_8 represent respectively, the number of OSDs per server, the size of OSD disks (in GB), and the storage overcommit ratio, while g_{14} is the number of replicas per object. These entries are defined in Table 26 for the block storage HDD.

$$Size_{BlockStorageSSD} = \left(\frac{\left(\left(Quantity_{Block\ Storage\ Nodes\ (SSD)} - x_{12} \right) \times h_6 \times h_7 \right)}{h_{14}} \right) \times (1 + h_8)$$

EQUATION 140: BLOCK STORAGE SSD

$Quantity_{Block\ Storage\ Nodes\ (SSD)}$ represents the number of nodes for block storage SSD, h_6 , h_7 , h_8 , and h_{14} are defined in Table 29 for the block storage SSD, as respectively, the number of OSDs per server, the size of OSD disks (in GB), the storage overcommit ratio, and the number of replicas per object.

$$Size_{ObjectStorageLowDensity} = \left(\frac{\left(\left(Quantity_{Object\ Storage\ Nodes\ (Low\ density)} - x_{12} \right) \times i_6 \times i_7 \right)}{i_{14}} \right) \times (1 + i_8)$$

EQUATION 141: OBJECT STORAGE LOW DENSITY

Where $Quantity_{Object\ Storage\ Nodes\ (Low\ density)}$ is the number of nodes for the objects storage with low density. The entries i_6 , i_7 , i_8 , and i_{14} are defined in Table 30 for the object storage with low density as respectively, the number of OSDs per server, the size of OSD disks (in GB), the storage overcommit ratio, and the number of replicas per object.

$$\begin{aligned}
 & \text{Size}_{\text{ObjectStorageHighDensity}} \\
 &= \left(\frac{\left(\left(\left(\text{Quantity}_{\text{Object Storage Nodes (High density)}} - x_{12} \right) \times j_6 \times j_7 \right) \right)}{j_{14}} \right) \\
 & \times (1 + j_8)
 \end{aligned}$$

EQUATION 142: OBJECT STORAGE HIGH DENSITY

Similarly as Object storage with low density and both SSD and HDD object storage size obtained in Equation 140, Equation 139 and Equation 141 the entries j_6, j_7, j_8, j_{14} are defined in Table 32.

2.2.2.3 Resource comparison

To calculate the cost for resource allocation (compute and storage), we need to know the cost of vCPU, RAM, and storage, as well as the cost for staff. We remind that such cost has been estimated in the previous tables.

$$\text{Cost}_{vCPU} = F_9$$

EQUATION 143: COST FOR A VCPU

Where F_9 is obtained by Equation 143as presented in Table 25.

$$\text{Cost}_{RAM} = F_{11}$$

EQUATION 144: COST FOR A RAM

Where F_{11} is obtained by Equation 145 as presented in Table 25.

$$\text{Cost}_{\text{BlockStorageHDD}} = G_{14}$$

EQUATION 145: COST FOR A BLOCK STORAGE HDD

G_{14} is obtained in Equation 145 as shown in Table 27.

$$\text{Cost}_{\text{BlockStorageSSD}} = H_{14}$$

EQUATION 146: COST FOR A BLOCK STORAGE SSD

H_{14} is obtained Equation 146 as shown in Table 29.

$$\text{Cost}_{\text{ObjectStorageLowDensity}} = I_{14}$$

EQUATION 147: COST FOR AN OBJECT STORAGE LOW DENSITY

Where I_{14} is estimated by Equation 147 as presented in Table 31.

$$\text{Cost}_{\text{ObjectStorageHighDensity}} = J_{14}$$

EQUATION 148: COST FOR AN OBJECT STORAGE HIGH DENSITY

Where J_{14} is estimated by Equation 148 as presented in Table 33. For the staff cost, it is given per vCPU, RAM, and storage according to the Cloud size (small, medium or large Cloud):

$$CostStuff_{vCPU} = (1 + x_{13}) \times \frac{\frac{Stuff_{Cost}}{4}}{NumberVCPU}$$

EQUATION 149: COST STUFF VCPU

Where x_{13} is a variable defined in Table 13, as the margin (for cost recovery or profit), $Stuff_{Cost}$ is the staff cost per month obtained in Equation 129, and $NumberVCPU$ is the total capacity provided in term of vCPUs, obtained by Equation 150.

$$\begin{array}{l} \text{Equation} \\ 150 \end{array} \quad CostStuff_{RAM} = (1 + x_{13}) \times \frac{\frac{Stuff_{Cost}}{4}}{Size_{RAM}}$$

EQUATION 150: COST STUFF RAM

Where, $Size_{RAM}$ is the total memory provided by the infrastructure, it is obtained by Equation 150:

$$\begin{array}{l} CostStuff_{Storage} \\ = (1 + x_{13}) \\ \times \frac{\frac{Stuff_{Cost}}{2}}{Size_{BlockStorageHDD} + Size_{BlockStorageSSD} + Size_{ObjectStorageHighDensity}} \end{array}$$

EQUATION 151: COST STUFF STORAGE

Where $Size_{BlockStorageHDD}$, $Size_{BlockStorageSSD}$, and $Size_{ObjectStorageHighDensity}$ represent respectively the total storage capacity for HDD, SSD, and Object storage with high density, which are obtained Equation 139, Equation 140 and Equation 142. The power cost is also calculated per vCPU, RAM, and storage:

$$CostPower_{vCPU} = 0.8 \times \frac{F_5}{F_7}$$

EQUATION 152: COST POWER VCPU

$$CostPower_{RAM} = 0.2 \times \frac{F_5}{F_8}$$

EQUATION 153: COST POWER RAM

$$CostPower_{BlockStorageHDD} = \frac{G_5}{G_{10}}$$

EQUATION 154: COST POWER BLOCK STORAGE HDD

$$CostPower_{BlockStorageSSD} = \frac{H_5}{H_{10}}$$

EQUATION 155: COST POWER BLOCK STORAGE SDD

$$CostPower_{ObjectStorageLowDensity} = \frac{I_5}{I_{10}}$$

EQUATION 156: COST POWER OBJECT STORAGE LOW DENSITY

$$CostPower_{ObjectStorageHighDensity} = \frac{J_5}{J_{10}}$$

EQUATION 157: COST POWER OBJECT STORAGE HIGH DENSITY

The values F_5 , F_7 , and F_8 represent respectively the cost per month for the power, the number of vCPUs per node, and the memory size (in GB) per node. These values are obtained using Equation 39, Equation 41 and Equation 42 respectively, the G_5 , G_{10} are defined in Table 27 by Equations respectively, Equation 51 and Equation 54 as the cost per month for the power, and the usable capacity of the cluster in GB. H_5 , H_{10} represent respectively, the cost per month for the power consumption for the block storage SSD and the usable capacity of the cluster in GB, represented via Equation 155 and Equation 54 as shown in Table 29. I_5 , I_{10} are estimated by Equation 62 and Equation 66 as highlighted in Table 31. Finally, J_5 , J_{10} are defined in Table 33 using Equation 92 and Equation 96. For the License, we compute it for the vCPU:

$$CostLicenceVMWare_{vCPU} = \frac{VMWareLicence_{cost}}{NumberVCPU}$$

EQUATION 158: COST LICENSE VMWARE VCPU

$VMWareLicence_{cost}$ is defined in Equation 132 and $NumberVCPU$ in Equation 137. The price for the Cloud is then computed as the sum of a resource cost (i.e., vCPU, RAM, block storage HDD/SSD, and object storage high/low density) with the people and the power cost for that resource. Here an example of these costs for large, medium and small Cloud.

$$price_r = cost_r + CostStuff_r + CostPower_r + CostLicenceVMWare_r$$

EQUATION 159: FINAL PRICE

Where r represents a resource, which can be a vCPU, RAM, block storage HDD/SSD, or object storage low/high density. In annex A.1.10, we estimate three the virtualized resources costs computed for the large, medium and small private clouds. Finally, the allocation cost of a VM with a flavour $c_m r_n$ (i.e., m vCPUs with n GB memory (RAM)) is calculated as follow:

$$cost_{c_m r_n} = m \times F_9 + n \times F_{11} + \frac{(m + n) \times \sum_{hw} HW_{Cost}}{NumberVCPU + Size_{RAM} + Size_{BlockStorageSSD}}$$

EQUATION 160: FINAL COST

Where $\sum_{hw} HW_{Cost}$ represents the sum of hardware costs, which include the spine, leaf, management switches and controller nodes. F_9 and F_{11} represent respectively the cost per vCPU per month and the cost per GB of RAM per month. These values are obtained by Equation 152 and Equation 153 from Table 25. Figure 14 depicts the costs per month (in \$) for VM reservation on three types of Cloud; private small, private

medium, and private large. These costs depend on the VM flavour $c_x r_y$ where c_x represent the number (i.e., x) of vCPUs and r_y is the memory size (i.e., y) in GB. x and y are integers that vary from 1 to respectively, 32 and 128 with a power of 2. We may remark two things:

- *The cost is proportional to the flavour:* indeed, more the flavour increases (i.e., x and/or y increase), the more will be the cost of the reservation, which is obvious because increasing the flavour means that the size of the allocated resource is increased.
- *Small Cloud costs more expensive than the medium, which in turn is more expensive than large Cloud:* we believe that an infrastructure with small amount of resources is naturally more expensive than an infrastructure with large amount of resources; this is due to the scarcity of these resources.

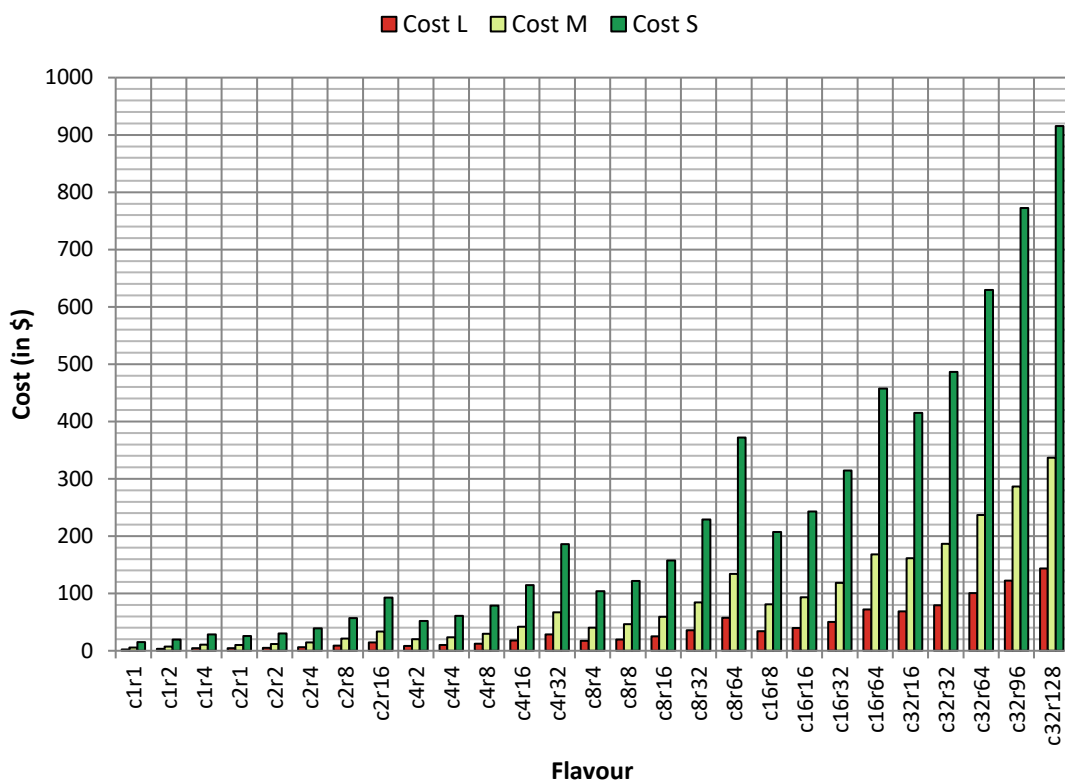


FIGURE 14: THE TCO FOR THREE SIZES OF CLOUD SIZE BASED ON THE PREVIOUS EXAMPLES OF HARDWARE NODES

Figure 15 depicts the costs per month for several flavours on Amazon AWS Europe during the period of April 5th, 2018. These costs depend on the reservation strategy that has been adopted. The *on-demand* reservation, let you reserve the exact capacity you need for any duration, in the location you need, and can keep for as long as you need it. The resources are activated as soon as they are requested, and they stay active until cancelled. Once created, the EC2 capacity is held for you regardless of whether you run the instances or not. The *spot* reservation is another strategy, where the request is specified with the number of instances, their types, the availability zone, the maximum price that the customer is willing to pay for an instance per hour. If the maximum price

exceeds the current spot price, the request will be satisfied immediately if the capacity is available, otherwise, The Amazon AWS waits until the request can be satisfied or cancelled. A spot request can be in one of the following modes; open -when the request is waiting to be fulfilled, active -the request is satisfied and has an associated Spot instance, failed -the request has one or more bad parameters, closed -the Spot instance is interrupted or terminated, and finally cancelled - the request is cancelled or expired. The Spot instance can be specified with duration. Therefore, the Spot instance will run continuously for the chosen duration without interruption. In our case, *Spot 1h*, means that, when active, the instance will run during one hour without interruption. The *reserved* instance is a billing discount applied to the use of the on-demand instance. From this figure we notice that the on-demand reservation is the most expensive in comparison with the other strategies. The cost increases with the flavour. This is obvious, as we request for larger flavour with more resources. We are surprised to see the strategy with Spot 1h is more expensive than the reserved strategy. We may explain that by the addition of the hard constraint of continuously running the instance during one hour, without considering the processing load of the node. The spot strategy is the cheapest one as it is best efforts; the reservation is done when resources are available.

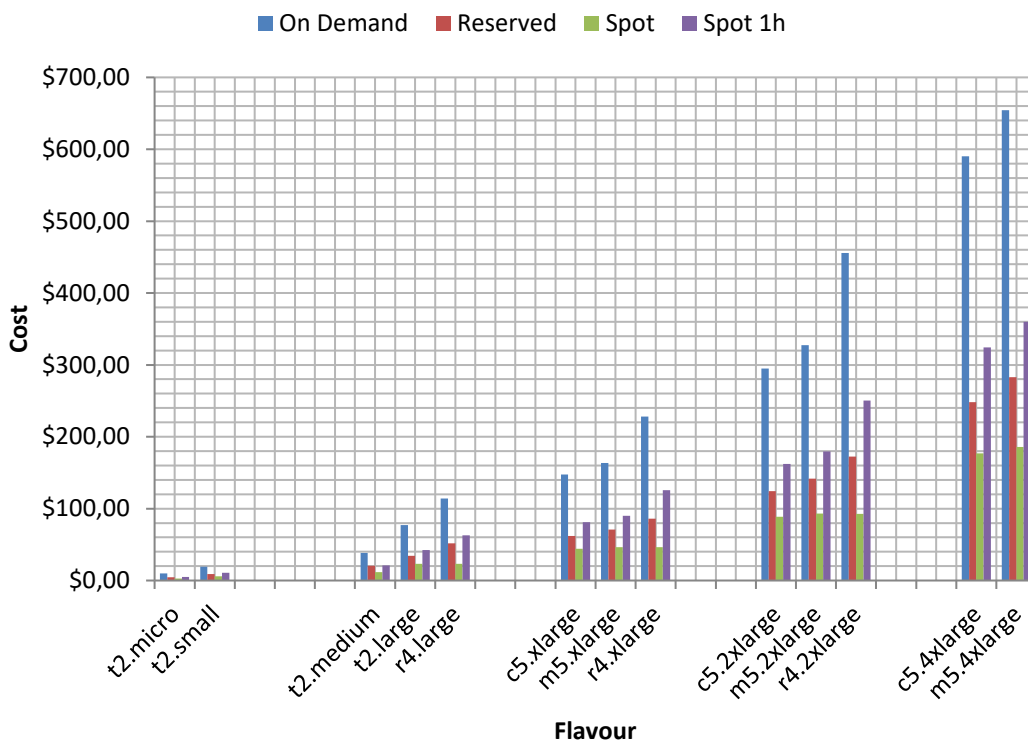


FIGURE 15: COST PER MONTH VM RESERVATION ON AMAZON AWS EUROPE

2.2.3 Analytical pricing modelling

In the 5G-TRANSFORMER ecosystem, price modelling refers to the study of the cost, revenue and profit derived from the detachment of VNFs from their services and their later allocation to different domains.

2.2.3.1 General considerations

There are different actors in the service price modelling that will determine the profitability of the 5G-TRANSFORMER business model. The actors playing a role in the service price modelling are as follows:

- *Infrastructure owners:* Infrastructure owners have datacentres of several sizes placed in different locations both close to the user and close to the edge of the network as well as links to connect them. In the latter experimental analysis, infrastructure owners can be directly mapped to domains containing datacentres as explained in 2.3.3.3. and 2.3.3.2 respectively.
- *Service providers:* Service providers are also known as verticals. They provide services by making use of infrastructure owner's resources and connectivity assets. In the latter experimental analysis, service providers can be directly mapped to the services, as explained in 2.3.3.1.
- *End users:* End users are those users calling at services provided by service providers.

Hereafter, we provide an analysis based on the cost, revenue, profit and the consequent dimensioning of the system.

2.2.3.2 Cost

The costs, understood as a monetary compensation given up in order to obtain a good or a service, are different for each of the actors as it is described hereafter:

- *Infrastructure owners:* Infrastructure owners will lead with the costs of maintaining and operating their datacentres and links as well as those of using the operator connectivity assets. In addition, when a service provider request exceeds the number of resources that the infrastructure owner can provide the federation enters into action. Infrastructure owners have to deal also with federation costs.
- *Service providers:* Service providers will assume the cost of using a network slice provided by the combination of both datacentres and connectivity assets that an infrastructure owner orchestrates.
- *End users:* End users are those paying for the cost of using the services provided by service providers.

2.2.3.3 Revenue

The revenue, understood as the income generated from sale of goods or services, or any other use of capital or assets will be coming from the price that, excluding the end users understood here as costumers, is set for each group of actors:

- *Infrastructure owners:* Infrastructure owners perceive revenue coming from the use of their infrastructure as well as connectivity assets by service providers. When federation takes place, the infrastructure owner that deploys the VNFs perceives a margin
- *Service providers:* Service providers will gain revenue from the payments of the subscription of final users to the service provider services. The process of setting prices is not a straightforward process. Apart from the revenue coming from the service price, there are other items that are also charged to the end user and that provide direct revenue to service providers:

- *Transactions*: For a given VNF, when it is already deployed, it can be scaled up or down. Each scaled, supposes a fee apart from the payment of the new consumed resources if scaled up or just the fee when scaling down. These fees are also direct revenues.
- *Licenses*: For a given VNF, when it is already deployed, several licenses are required for making it work. These licenses will change according to the number of users using the VNF, the type of VNF or the performance that it is required from it. Licenses suppose a direct revenue.
- *Special requirements*: For a given service, the special requirements such as domain affinity, datacentre affinity or protection are also charged to the user. Special requirements suppose direct revenue.

2.2.3.4 Profit

For both infrastructure owners and service providers the profit can be calculated as the subtraction of cost to the revenue. According to the previous idea, the profit for infrastructure providers and service providers will be:

- *Infrastructure owners*: the profit can be calculated as the sum of all the revenues coming from the service providers minus the cost of maintaining their infrastructure.
- *Service providers*: Their profit will come from the sum of all the revenues coming from the end user payments minus the cost of accessing the network slices necessary to operate their services inside the infrastructure owners which results in the margin profit of 2.2.3.3.

2.2.3.5 Price setting

In order to develop the price setting analysis, the view of a service provider that has some infrastructure is taken in this section. This means that in the following, the view of a new entity that acts both as infrastructure owner and service provider is assumed as the local domain. The other infrastructure owners are referred to as overflow domains.

As the services proposed in the 5G-T ecosystem are novel services, it is not possible to make an estimation of the demand that they will have. Therefore, it is complicated to set a price for them. In addition, a further degree of complexity is introduced with the use of federation, as the flows of money follow a completely different approach.

The concept of breakeven price defines the price that a service should have to ensure that cost and revenue are the same and no profit or losses are perceived. The breakeven price seems to be a good starting point for setting a price for the 5G-TRANSFORMER services. It should cover the situation in which the local infrastructure is used as its maximum capacity and also when it is underused leading to a vacancy of resources. The breakeven price can be defined as a function of the cost, the maximum demand and the forecast utilization of the infrastructure, as indicated by Equation 161:

$$B = \frac{C_l}{d_{l_{max}} \cdot u_{f_l}} \quad 0 \leq u_{f_l} \leq 1$$

EQUATION 161: BREAKEVEN PRICE

Where:

- $B \in \mathbb{R}^+$ is the breakeven price for the service.
- $d_{l,max} \in \mathbb{N}$ is the maximum local demand for the service.
- $C_l \in \mathbb{R}^+$ is the cost for the infrastructure of the local domain.
- $u_{f,l} \in \mathbb{R}$ is the forecast utilization of the local infrastructure.

Based on the previous equation, a price for a particular service can be easily derived by defining a profit margin in Equation 162:

$$P_s = B(1 + m) \quad m > 0$$

EQUATION 162: PRICE

Where:

- $P_s \in \mathbb{R}^+$ is the price for the service.
- $m \in \mathbb{R}^+$ is the profit margin set to the breakeven price.

The profit margin is set in such a way that in any situation it can ensure profit. As a consequence, the profit margin must ensure profit even in a situation in which a service is federated to an overflow domain that can be operating at a high capacity and charges a more expensive price for deploying their VNFs. This idea is presented in Equation 163:

$$m > m_f - m_{max}$$

EQUATION 163: PROFIT MARGIN

Where:

- $m_f \in \mathbb{R}^+$ is the federation margin as defined in the Equation 165.
- $m_{max} \in \mathbb{R}^+$ is the margin set for a high occupation of the infrastructure of an overflow domain participating in the federation.

The direct revenue items are defined by Equation 164. It is important to note here that Equation 164 must be understood as a general case that has to be particularized for each service instance.

$$r = r_t + \sum_{i=1}^N n_{l_i} \cdot r_{l_i} + r_s$$

EQUATION 164: DIRECT REVENUE ITEMS

Where:

- $r_t \in \mathbb{R}^+$ is the fee paid for scaling up or down the service VNFs. This term will enter into stage when a scale event arrives to the local domain.
- $N \in \mathbb{N}$ is the number of VNFs in a service.
- $n_l \in \mathbb{N}$ is the number of licenses per VNF.
- $r_l \in \mathbb{R}^+$ is the price of the license for a VNF.
- $r_s \in \mathbb{R}^+$ is the fee paid for deploying the service with any special requirement such as domain affinity, datacentre affinity or protection.

With the direct revenue items defined, the federation margin can be determined by Equation 165:

$$m_f = 1 - \frac{P_f + r}{P_s} \quad P_s > P_f + r$$

EQUATION 165: FEDERATION MARGIN

Where:

- $P_f \in \mathbb{R}^+$ is the price that an overflow domain charges the local domain for deploying a service VNF in its infrastructure. It is always lower than the market price, as it is a wholesale price that is charged only among domains and includes the connectivity cost that the federation could have.

The real utilization of the local infrastructure is defined in Equation 166 as the relationship between the real demand of the local infrastructure and the maximum demand of the real infrastructure:

$$u_{r_l} = \frac{d_{r_l}}{d_{max_l}}$$

EQUATION 166: REAL UTILIZATION OF THE LOCAL INFRASTRUCTURE

Where:

- $d_{r_l} \in \mathbb{N}$ is the real demand attended by the local domain.

The revenue that is perceived in the 5G-TRANSFORMER ecosystem can be defined as the revenue perceived by the local demand and the revenue federated demand, as indicated in Equation 167:

$$R = \begin{cases} d_{r_l}(P_s + r) & u_{r_l} < 1 \\ d_{max_l}(P_s + r) + d_{r_f}(P_s - P_f - r) & u_{r_l} \geq 1 \end{cases}$$

EQUATION 167: REVENUE

Where:

- $d_{r_f} \in \mathbb{N}$ is the real demand attended by the overflow domains. This demand is only different from zero if the local domain infrastructure is operating at its maximum capacity and cannot deploy any further VNF.

Profit can be defined as the subtraction of cost to revenue particularized for both the own and the federated demand, as indicated in Equation 168:

$$P = R - C$$

EQUATION 168: PROFIT

Equation 168 can be rewritten by substituting the revenue for both the local and federated demand, resulting in Equation 169. In this equation there is no cost for the demand attended by the overflow domains as it is not paid by the local domain.

$$P = \begin{cases} d_{rl}(P_s + r) - C_l & u_{rl} < 1 \\ d_{max_l}(P_s + r) - C_l + d_{rf}(P_s - P_f - r) & u_{rl} \geq 1 \end{cases}$$

EQUATION 169: REVENUE-PROFIT

By substituting the price in Equation 169, an equation that is also dependant on the breakeven price is obtained in Equation 170:

$$P = \begin{cases} d_{rl}(B(1 + m) + r) - C_l & u_{rl} < 1 \\ d_{max_l}(B(1 + m) + r) - C_l + d_{rf}(P_s - P_f - r) & u_{rl} \geq 1 \end{cases}$$

EQUATION 170: REVENUE-PRICE-PROFIT

Finally, by substituting the breakeven in Equation 170, an equation that is also dependant on the maximum local demand and the forecast utilization of the local infrastructure is obtained as in Equation 171:

$$P = \begin{cases} d_{rl} \left(\frac{C_l(1 + m)}{d_{lmax} \cdot u_{f_l}} + r \right) - C_l & u_{rl} < 1 \\ d_{max_l} \left(\frac{C_l(1 + m)}{d_{lmax} \cdot u_{f_l}} + r \right) - C_l + d_{rf}(P_s - P_f - r) & u_{rl} \geq 1 \end{cases}$$

EQUATION 171: REVENUE-PRICE-BREAKEVEN-PROFIT

The previous equation provides an accurate model for analysing the profitability of the 5G-TRANSFORMER ecosystem. However, it also provides a great uncertainty, as the demand and therefore the utilization of the local infrastructure is not known a priori. In order to overcome this uncertainty, a sensitivity analysis is proposed.

A sensitivity analysis determines how different values of an independent variable affect a particular dependent variable under a given set of assumptions [14].

This type of analysis allows the 5G-TRANSFORMER ecosystem to deal with different levels of uncertainty in the utilization of the local infrastructure and to create a mathematical model that contributes to a better understanding of the overall model's uncertainty.

By applying the sensitivity analysis to the utilization of the local infrastructure, five possible scenarios are defined:

- *Pessimistic scenario*: 20% utilization of the local infrastructure.
- *Realistic scenario*: 50% utilization of the local infrastructure.
- *Optimistic scenario*: 80% utilization of the local infrastructure.
- *Ideal scenario*: 100% utilization of the local infrastructure.
- *Overflow scenario*: +100% utilization of the local infrastructure.

With these five scenarios, it is possible to continue the analysis for the pricing setting by analysing each of them separately.

In order to ease the analysis, it is considered that there's no direct revenue for any of the services that are deployed, which implies that $r = 0$ and that Equation 171 can be simplified to Equation 172:

$$P = \begin{cases} d_{r_l} \frac{C_l(1+m)}{d_{l_{max}} \cdot u_{f_l}} - C_l & u_{r_l} < 1 \\ \frac{C_l(1+m)}{u_{f_l}} - C_l + d_{r_f}(P_s - P_f) & u_{r_l} \geq 1 \end{cases}$$

EQUATION 172: SIMPLIFIED REVENUE-PRICE-BREAKEVEN-PROFIT

Equation 172 will be the equation as the reference equation for the sensitivity analysis in the following subsections.

2.2.3.5.1 Pessimistic scenario

In a scenario where the real utilization of the local infrastructure raises until 0.2, the case $u_{r_l} < 1$ defined in Equation 172 is fulfilled, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

By substituting the real utilization of the local infrastructure in Equation 172 an equation that includes the real utilization of the local infrastructure is obtained as shown in Equation 173:

$$P = u_{r_l} \frac{C_l(1+m)}{u_{f_l}} - C_l$$

EQUATION 173: LOCAL-SIMPLIFIED REVENUE-PRICE-BREAKEVEN-PROFIT

Finally, by taking the cost for the infrastructure of the local domain to the other size of the equation the profit-cost ratio can be defined in Equation 174:

$$R(P/C_l) = u_{r_l} \frac{1+m}{u_{f_l}} - 1$$

EQUATION 174: LOCAL PROFIT-COST RATIO

The profit-cost ratio provides a great overview of the profitability of the 5G-TRANSFORMER ecosystem. The greater this margin is, the greater the profit will be. According to Equation 174, the profit-cost ratio depends on the real utilization of the local infrastructure, the forecast utilization of the local infrastructure and the profit margin.

As stated before, the pessimistic scenario sets the real utilization of the local infrastructure to 0.2. Therefore, there are still two degrees of freedom: the margin-profit and the forecast utilization of the local infrastructure. Their relationship between these two variables is shown in Figure 16:

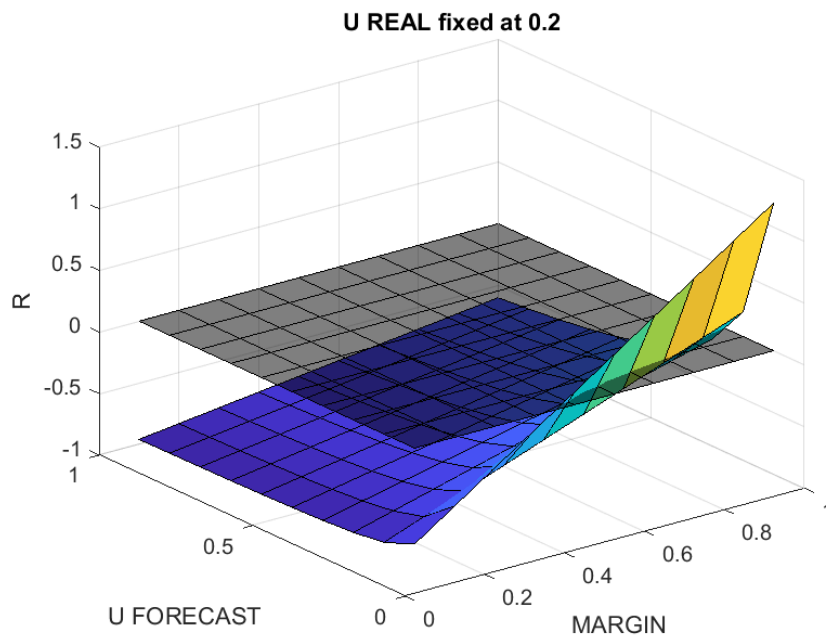


FIGURE 16: PESSIMISTIC SCENARIO

Their relationship can be more easily understood by representing in the plain projections in the figures contained in Table 36:

TABLE 36: PESSIMISTIC SCENARIO

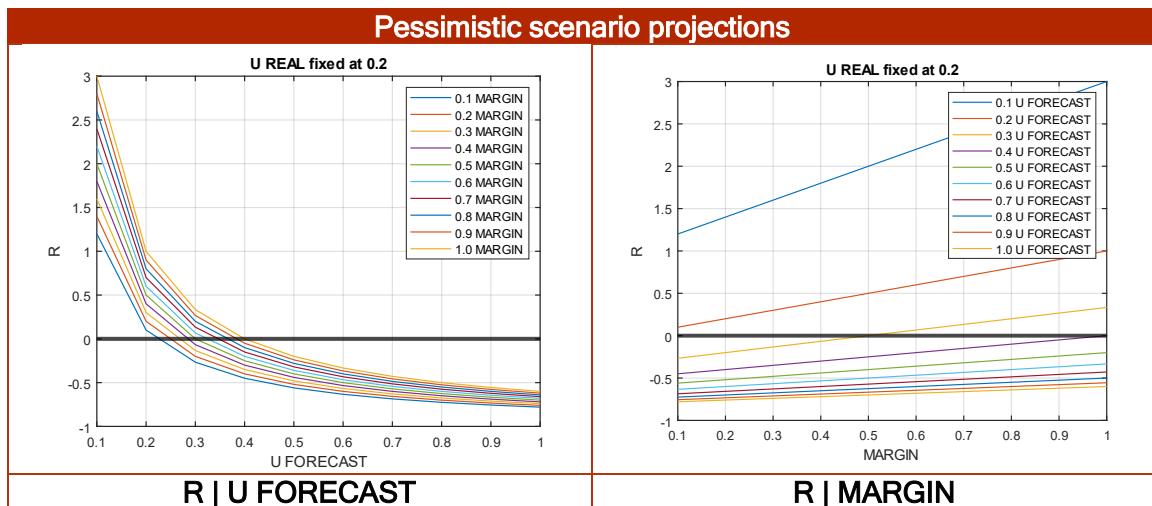


Figure 16 also represents the plane where $R=0$, that is the breakeven plane in which the ratio cost-profit is nor negative nor positive.

It can be derived that with $u_{real} = 0.2$ the forecast utilization of the local infrastructure has to be greater than 0.3 for a minimum profit margin of 0.5 (if greater than the profit margin indicated by Equation 163, otherwise it would be greater than $m_f - m_{max}$).

In addition, it can be also derived that with $u_{real} = 0.2$ with a profit margin of 0.1 the minimum forecast utilization of the local infrastructure has to be greater than 0.2.

2.2.3.5.2 Realistic scenario

In a scenario where the real utilization of the local infrastructure raises until 0.5, the case $u_{rl} < 1$ defined in Equation 172 is again fulfilled and, as there is still vacancy in the local domain and no service will be federated to any overflow domain, Equation 173 can be used to analyse the realistic scenario.

As indicated before in Equation 174, the real scenario sets the real utilization of the local infrastructure to 0.5. As before, there are still two degrees of freedom: the margin-profit and the forecast utilization of the local infrastructure. Their relationship between these two variables is shown in Figure 17:

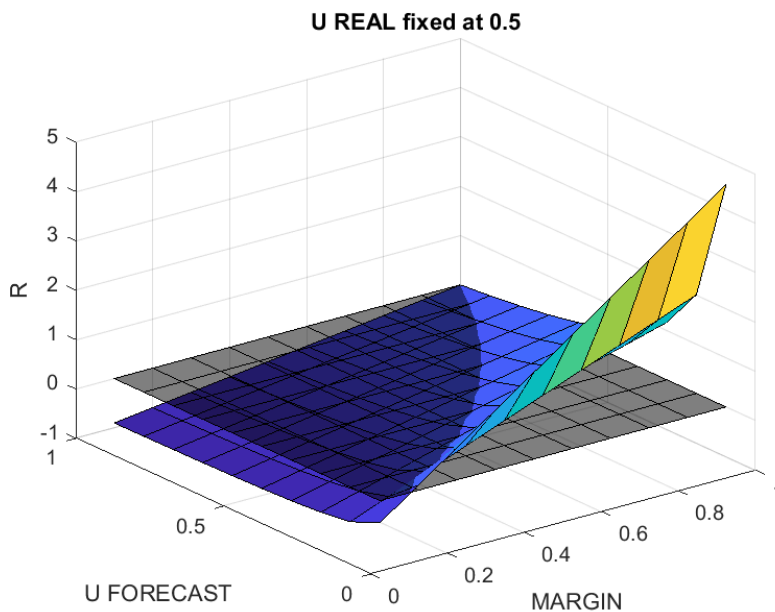


FIGURE 17: REALISTIC SCENARIO

Their relationship can be more easily understood by representing the plain projections in the figures contained in Table 37.

TABLE 37: REALISTIC SCENARIO

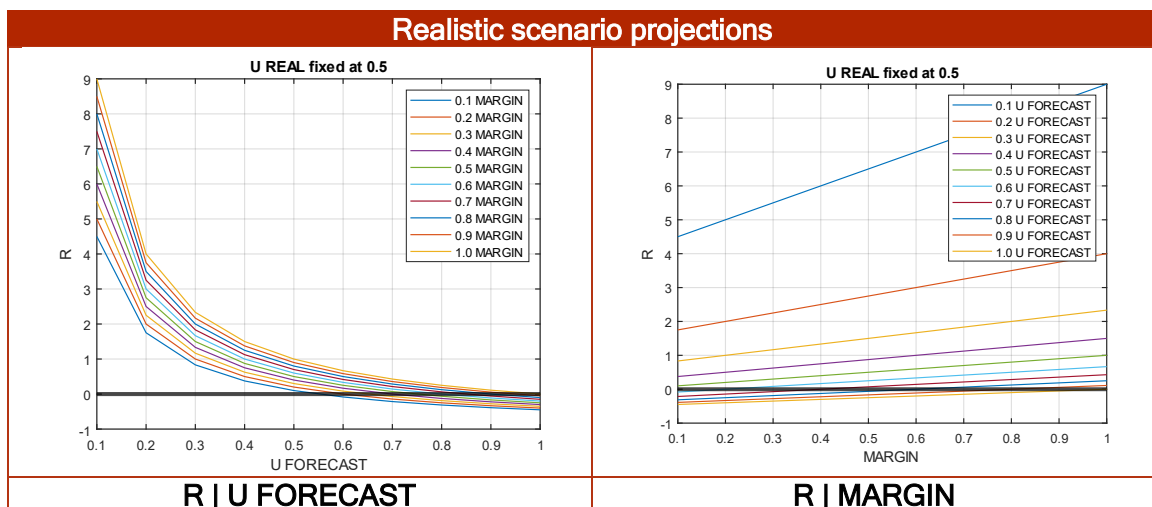


Figure 17 also represents the plane where $R=0$, that is the breakeven plane in which the ratio cost-profit is neither negative nor positive.

It can be derived that with $u_{real} = 0.5$ the forecast utilization of the local infrastructure has to be greater than 0.6 for a minimum profit margin of 0.2 (if greater than the profit margin indicated by Equation 163, otherwise it would be greater than $m_f - m_{max}$).

In addition, it can be also derived that with $u_{real} = 0.5$ with a profit margin of 0.1 the minimum forecast utilization of the local infrastructure has to be greater than 0.5.

2.2.3.5.3 Optimistic scenario

In a scenario where the real utilization of the local infrastructure raises until 0.8, the case $u_{rl} < 1$ defined in Equation 172 is again fulfilled and, as there is still vacancy in the local domain and no service will be federated to any overflow domain, Equation 173 can be used to analyse the optimistic scenario.

As indicated before in Equation 174, the real scenario sets the real utilization of the local infrastructure to 0.8. As before, there are still two degrees of freedom: the margin-profit and the forecast utilization of the local infrastructure. Their relationship between these two variables is shown in Figure 18:

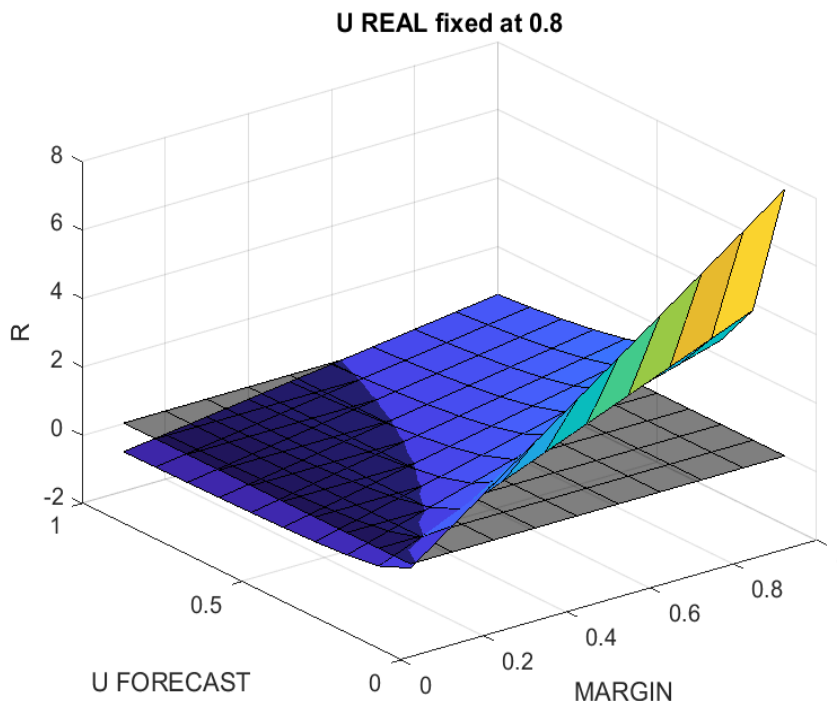


FIGURE 18: OPTIMISTIC SCENARIO

This situation can be also understood by representing the plain projections in the figures contained in Table 38.

However, this time is not so clear as the different profit margins represented on the left figure are pretty closed and it is not easy to distinguish when they cut with the 0 line.

The same happens with the forecast utilizations of the local infrastructure in the right figure. They are also pretty closed, and it is not easy to distinguish when they cut with the 0 line.

TABLE 38: OPTIMISTIC SCENARIO

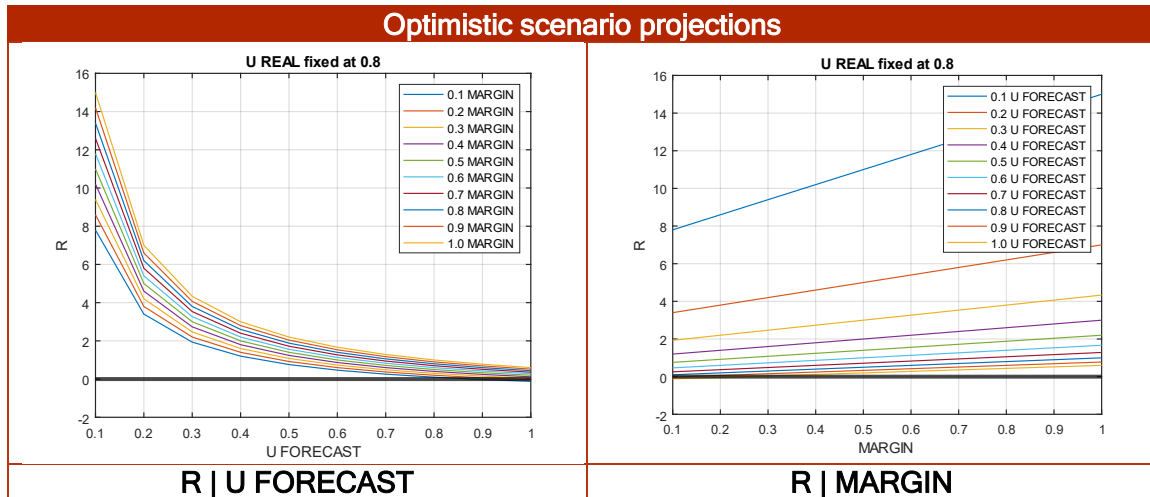


Figure 18 also represents the plane where $R=0$, that is the breakeven plane in which the ratio cost-profit neither is negative nor positive.

It can be derived that with $u_{real} = 0.8$ the forecast utilization of the local infrastructure has to be greater than 0.9 for a minimum profit margin of 0.1 (if greater than the profit margin indicated by Equation 163, otherwise it would be greater than $m_f - m_{max}$).

In addition, it can be also derived that with $u_{real} = 0.8$ with a profit margin of 0.1 the minimum forecast utilization of the local infrastructure has to be greater than 0.9.

2.2.3.5.4 Ideal scenario

The ideal scenario is slightly different from the previous scenarios, as the local infrastructure is utilized as its maximum capacity.

In a scenario where the real utilization of the local infrastructure raises until 1, the case $u_{rl} \leq 1$ defined in Equation 172 is again fulfilled and Equation 173 can be used to analyse the ideal scenario. However, as in the ideal scenario the real and forecast utilization of the local infrastructure are equal to one, Equation 173 can be simplified to Equation 175:

$$P = C_l \cdot m$$

EQUATION 175: IDEAL REVENUE-PRICE-BREAKEVEN-PROFIT

The ideal profit-cost ratio doesn't depend on the real utilization of the local infrastructure and Equation 174 can be simplified to Equation 176 which results to be the profit margin m :

$$R(P/C_l) = m$$

EQUATION 176: IDEAL PROFIT-COST RATIO

As a result, the main difference between the previous scenarios and the ideal scenario is that there only one degree of freedom, the profit margin.

The previous idea is represented in Figure 19:

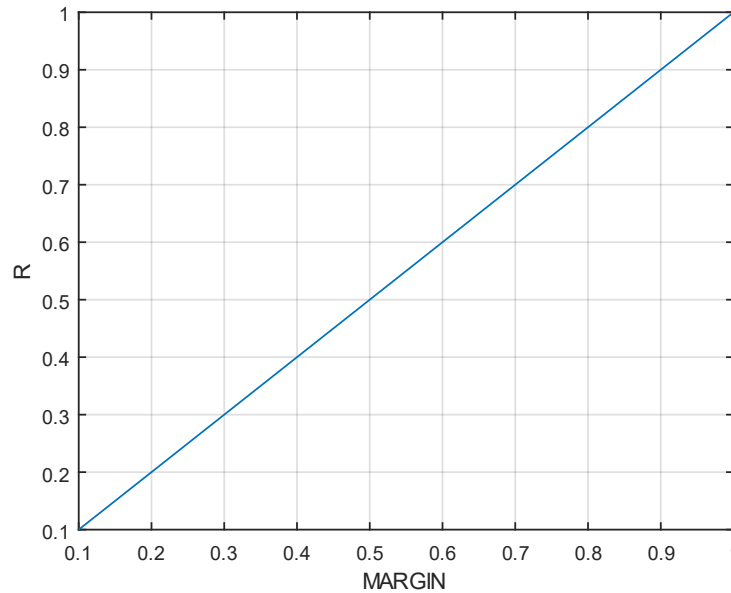


FIGURE 19: IDEAL SCENARIO

It can be concluded that there is a linear relationship between the profit margin and the profit-cost ratio.

2.2.3.5.5 Overflow scenario

In a scenario where the real utilization of the local infrastructure raises over the maximum capacity of the local infrastructure, the case $u_{r_l} = 1$ defined in Equation 172 is fulfilled and, as there is no vacancy in the local domain and every service will be federated to some overflow domain. Equation 173 cannot be used to analyse the realistic scenario.

As a result, a new equation describing the federation scenario must be derived in Equation 177 taking into consideration that the real and forecast utilization of the local infrastructure are equal to one as in the ideal scenario:

$$P = C_l \cdot m + d_{r_f}(P_s - P_f)$$

EQUATION 177: FEDERATED-SIMPLIFIED REVENUE-PRICE-BREAKEVEN-PROFIT

The overflow profit-cost ratio can be obtained by moving the cost for the infrastructure of the local domain to the other side of the equation and multiplying by the cost for the infrastructure of the local domain to the last term of Equation 177, as indicated in Equation 178:

$$R(P/C_l) = m + C_l \cdot d_{r_f}(P_s - P_f)$$

EQUATION 178: FEDERATED PROFIT-COST RATIO

The overflow scenario can be understood as an extension of the ideal scenario. The overflow scenario uses the ideal scenario model until there is no vacancy left in the local domain and then, when a new service demand arrives, federation takes places and the overflow infrastructure is used.

In Equation 165 it was stated that the $P_s > P_f$ as P_f is a wholesale price that is charged only charged among domains and includes the connectivity cost that the federation could have because of that, it can be affirmed that the term $C_l \cdot d_{rf}(P_s - P_f)$ will represent a positive value. As the profit margin is also positive by definition, the profit-cost ratio will take the form of a line with positive slope. How steep this slope is will depend on the difference between P_s and P_f . And as P_s is constant it can be concluded that the variable that determines the slope is P_f .

Accordingly, a new sensitivity analysis is applied to P_f as a function of P_s as indicated by Equation 179:

$$P_v = \frac{P_f}{P_s} \quad P_s > P_f$$

EQUATION 179: PRICE VARIABILITY IN THE OVERFLOW SCENARIO

By giving values to the price variability every 0.1, the federation scenario profit-cost ratio figure is represented in Figure 20:

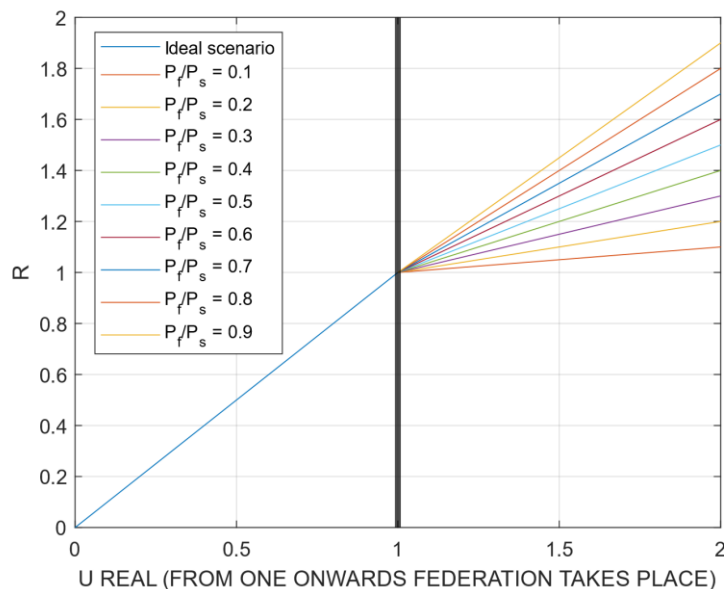


FIGURE 20: OVERFLOW SCENARIO

As it can be seen, the profit-cost ratio in the federation scenario will be always lower than in the ideal scenario. Depending on the price variability, the profit-cost ratio will be closer to that calculated for the ideal scenario.

2.2.3.6 Dimensioning

The dimensioning analysis aims at specifying the service demand definition for each of the services defined in 5G-TRANSFORMER.

2.2.3.6.1 Maximum number of services per datacentre

For each of the services, the total amount of resources that is needed to run them is calculated and then shown in Table 39:

TABLE 39: TOTAL AMOUNT OF RESOURCES NEEDED TO RUN A SERVICE

	S1	S2	S3	S4	S5	S6
CPU	4	12	15	11	17	38
RAM [GB]	4	12	29	22	34	79
Disk [GB]	60	120	125	240	280	580

With this information, it is calculated how many services can be run for each of the datacentre types by analysing each of the resources separately, CPU in Table 40, RAM in Table 41 and Disk in Table 42:

TABLE 40: MAXIMUM NUMBER OF SERVICES BY CPU

CPU	S1	S2	S3	S4	S5	S6
Small	144	48	38	52	33	15
Medium	864	288	230	314	203	90
Big	21528	7176	5740	7828	5065	2266

TABLE 41: MAXIMUM NUMBER OF SERVICES BY RAM

RAM	S1	S2	S3	S4	S5	S6
Small	378	126	100	137	88	39
Medium	2268	756	604	824	533	238
Big	56511	18837	15069	20549	13296	5948

TABLE 42: MAXIMUM NUMBER OF SERVICES BY DISK

Disk	S1	S2	S3	S4	S5	S6
Small	10800	3600	2880	3927	2541	1136
Medium	27000	9000	7200	9818	6352	2842
Big	58320	19440	15552	21207	13722	6138

In view of the previous tables, for all the services, the CPU is always the most restrictive resource. This means that CPU is the resource acting as the bottleneck and when the available CPUs are no longer available there will still be available RAM and disk storage. As a result, CPUs are going to determine the maximum number of services that can be run in a datacentre.

The maximum number of services that can be running in each type of datacentre simultaneously will be determined by Table 39 as it is indicated in Table 43:

TABLE 43: MAXIMUM NUMBER OF SIMULTANEOUS SERVICES BY DATACENTRE TYPE

	S1	S2	S3	S4	S5	S6
Small	144	48	38	52	33	15
Medium	864	288	230	314	203	90
Big	21528	7176	5740	7828	5065	2266

2.2.3.6.2 Service life-time

With the data provided in Table 44, it is possible to estimate the mean time service life-time, as the mean of the minimum life-time and the maximum life-time of each of the 5G-T services.

TABLE 44: SERVICE LIFE-TIME

Service ID	Description	Minimum life-time [Hours]	Maximum life-time [Hours]	Mean life-time [Hours]
1	Entertainment	0.5	6	3.25
2	MVNO	2160	8640	5400
3	Automotive	24	168	96
4	Non-emergency eHealth	168	672	420
5	Emergency eHealth	0.05	0.25	0.15
6	Robots	720	2160	1440

The previous table has to be understood as a first approach to the description of novel services proposed in the 5G-T ecosystem and its value would have to be adjusted with real data when available.

2.2.3.6.3 Service demand

Estimating the demand is not an easy task, especially when talking about novel services as those of 5G-TRANSFORMER.

According to the use case descriptions presented in and in order to proceed with the estimation of the demand for each of the possible scenarios, the service demand distribution can be calculated by using:

$$D_s = \frac{n_{smax}}{N_s}$$

EQUATION 180: SERVICE DEMAND DISTRIBUTION

Where:

- $D_s \in \mathbb{N}$ is the distributed demand for the service of type s , expressed as a percentage of N_s .
- $n_{smax} \in \mathbb{N}$ is the maximum number of services of type s that can be run in a particular datacentre, nominally the values of Table 43.
- $N_s \in \mathbb{N}$ is the total number of services, obtained as the sum of the maximum number of services per service type that can be run in a datacentre.

The total number of datacentre types is shown in Table 45:

TABLE 45: TOTAL NUMBER OF SERVICES PER DATACENTRE TYPE

	S1	S2	S3
Total number of services	330	1989	49603

Accordingly, the service demand distribution resulting from applying Equation 180 are shown in Table 46:

TABLE 46: SERVICE DEMAND DISTRIBUTION

	S1	S2	S3	S4	S5	S6
Demand	0,44	0,15	0,12	0,16	0,10	0,05

The previous table has to be understood as a first approach to the description of novel services proposed in the 5G-TRANSFORMER ecosystem and its value would have to be adjusted with real data when available.

The demand for the pessimistic scenario is indicated in Table 47:

TABLE 47: ESTIMATED DEMAND FOR THE PESSIMISTIC SCENARIO

UC	Description	Estimated demand		
		Small	Medium	Big
1	Entertainment	13	75	1879
2	MVNO	1	8	209
3	Automotive	1	5	132
4	Non-emergency eHealth	2	10	247
5	Emergency eHealth	1	4	101
6	Robots	0	1	21

The demand for the realistic scenario is indicated in Table 48:

TABLE 48: ESTIMATED DEMAND FOR THE REALISTIC SCENARIO

UC	Description	Estimated demand		
		Small	Medium	Big
1	Entertainment	31	189	4697
2	MVNO	3	21	522
3	Automotive	2	13	330
4	Non-emergency eHealth	4	25	617
5	Emergency eHealth	2	10	253
6	Robots	0	2	52

The demand for the optimistic scenario is indicated in Table 49:

TABLE 49: ESTIMATED DEMAND FOR THE OPTIMISTIC SCENARIO

UC	Description	Estimated demand		
		Small	Medium	Big
1	Entertainment	50	302	7515
2	MVNO	6	34	835
3	Automotive	4	21	529
4	Non-emergency eHealth	7	40	987
5	Emergency eHealth	3	16	405
6	Robots	1	3	82

The demand for the ideal scenario is indicated in Table 50. This scenario supposes the local infrastructure to be operating as its maximum capacity:

TABLE 50: ESTIMATED DEMAND FOR THE IDEAL SCENARIO

UC	Description	Estimated demand		
		Small	Medium	Big
1	Entertainment	63	377	9394
2	MVNO	7	42	1044
3	Automotive	4	26	661

4	Non-emergency eHealth	8	49	1234
5	Emergency eHealth	3	20	507
6	Robots	1	4	103

The demand for the federation scenario is indicated in Table 51. This scenario supposes the local infrastructure to be operating as its maximum capacity and some of the services are deployed by the overflow infrastructure managed by the federation. The overflow infrastructure attends an additional 300% of the local demand.

TABLE 51: ESTIMATED DEMAND FOR THE FEDERATION SCENARIO

UC	Description	Estimated demand		
		Small	Medium	Big
1	Entertainment	189	1131	28182
2	MVNO	21	126	3132
3	Automotive	12	78	1983
4	Non-emergency eHealth	24	147	3702
5	Emergency eHealth	9	60	1521
6	Robots	3	12	309

The previous tables have to be understood as a first approach to the description of novel services proposed in the 5G-TRANSFORMER ecosystem and its value would have to be adjusted with real data when available.

2.3 Experimental study

The aim of the experimental study is to carry out some simulations that validate the analytical study presented in section 2.2 and allows reaching some conclusions for the techno-economic analysis.

Hereafter, we present the methodology that we followed in our experimental study, the simulation system structure that we used, the actors involved in the study, the simulation scenarios settings and the tool description. We conclude the study with an analysis of the service pricing for each simulation scenario on the cost, the revenue and the profit of the service.

2.3.1 Methodology

The followed methodology consists of two main blocks: System and pricing, as it can be seen in Figure 21:

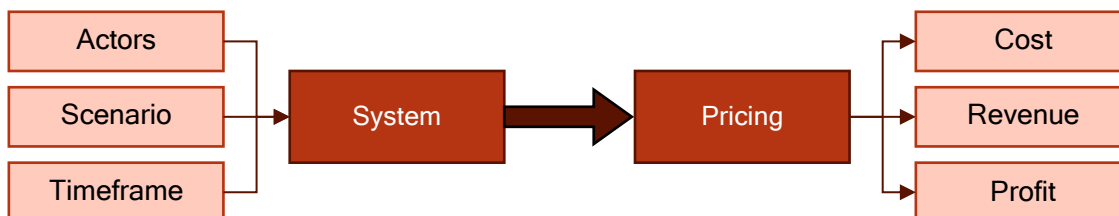


FIGURE 21: FOLLOWED METHODOLOGY

In the first block, the system structure is defined, taking into consideration the actors, the scenario setting and the timeframe of the experimental analysis.

In the second block, the pricing simulations are carried out. Taking into consideration the service pricing modelling defined in 2.2.3 and the system structure defined in the first block, the total cost for the system, revenue and profit are calculated.

2.3.2 Simulation system structure

A simulation scenario has been built up using MATLAB based on a simulation module that we have developed in order to validate the analytical modelling study through simulations. MATLAB is a multi-paradigm numerical computing development environment and proprietary programming language that allows different tasks such as matrix operation, plotting of data or implementation of algorithms.

Apart from its potential, MATLAB was selected due to its easiness to be integrated with other platforms and programming languages like Python.

The simulation system structure for the experimental study is compounded by three modules as depicted in Figure 22. These three modules oversee different functions in the simulation of the experimental study:

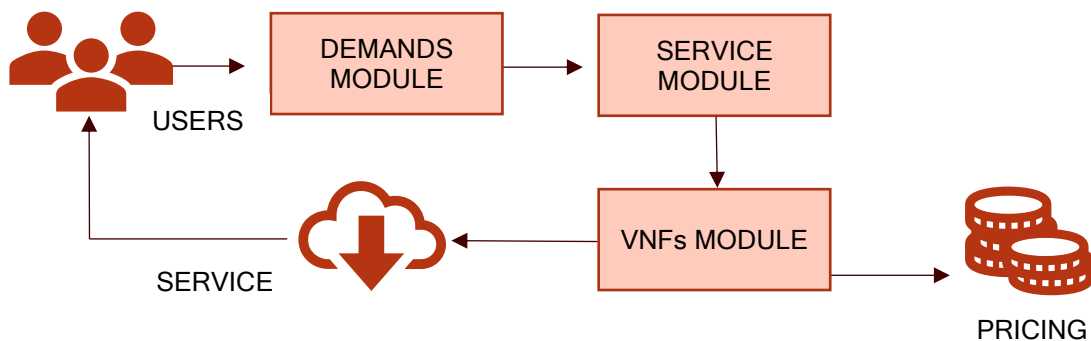


FIGURE 22: SYSTEM STRUCTURE

- *Demands module*: This module is in charge of generating service demands over time as well as the service network-graph decomposition into their VNFs and link assets.
- *Service module*: This module performs the service network-graph partitioning by associating each VNF with a domain.
- *VNFs module*: This module takes care of the mapping of resources of each of the VNFs compounding a service to the domain datacentre architecture complying with the services constrains and the notification of the remaining domain resources availability.

In 2.3.4, we provide a more-detailed explanation of each module, including the different elements they are composed of and the data they take as input.

2.3.3 Actors in the study

In the experimental study, three types of actors are distinguished: services, datacentres and domains. Those are considered actors in the sense that in all the scenarios that are considered in the study, services, datacentres and domains will be exchanging information and money.

- *Services*: Services are requested by tenants and will be understood as a series of resources that are used by those tenants for a concrete time length. In order

to do that, tenants will pay a service price, a license per each user, a fee for scaling up VNFs resources and also for special requirements such as protection, domain affinity and datacentre affinity.

- *Datacentres*: Datacentres are where the services are run. They have a cost that has to be paid, in the experimental analysis at the beginning of the year.
- *Domains*: Domains can be understood in the experimental study as the owners of the datacentres.

2.3.3.1 Service modelling

In the 5G-T ecosystem, services can be understood as a set of VNFs. At the same time, each of these VNFs can be understood as a set of resources (CPUs, RAM and disk) and connectivity assets with other VNFs (bandwidth and latency). The combination of both resources and connectivity assets among VNFs forms the service-graph of a vertical service.

In the following lines, the service-graphs, the services' VNFs, links and INPUT and OUTPUT bandwidths are specified to define the configuration of the experimental study described in 2.3. The service-graphs of the vertical use cases described and modelled in section 2.2.2 are summarized in Figure 23:

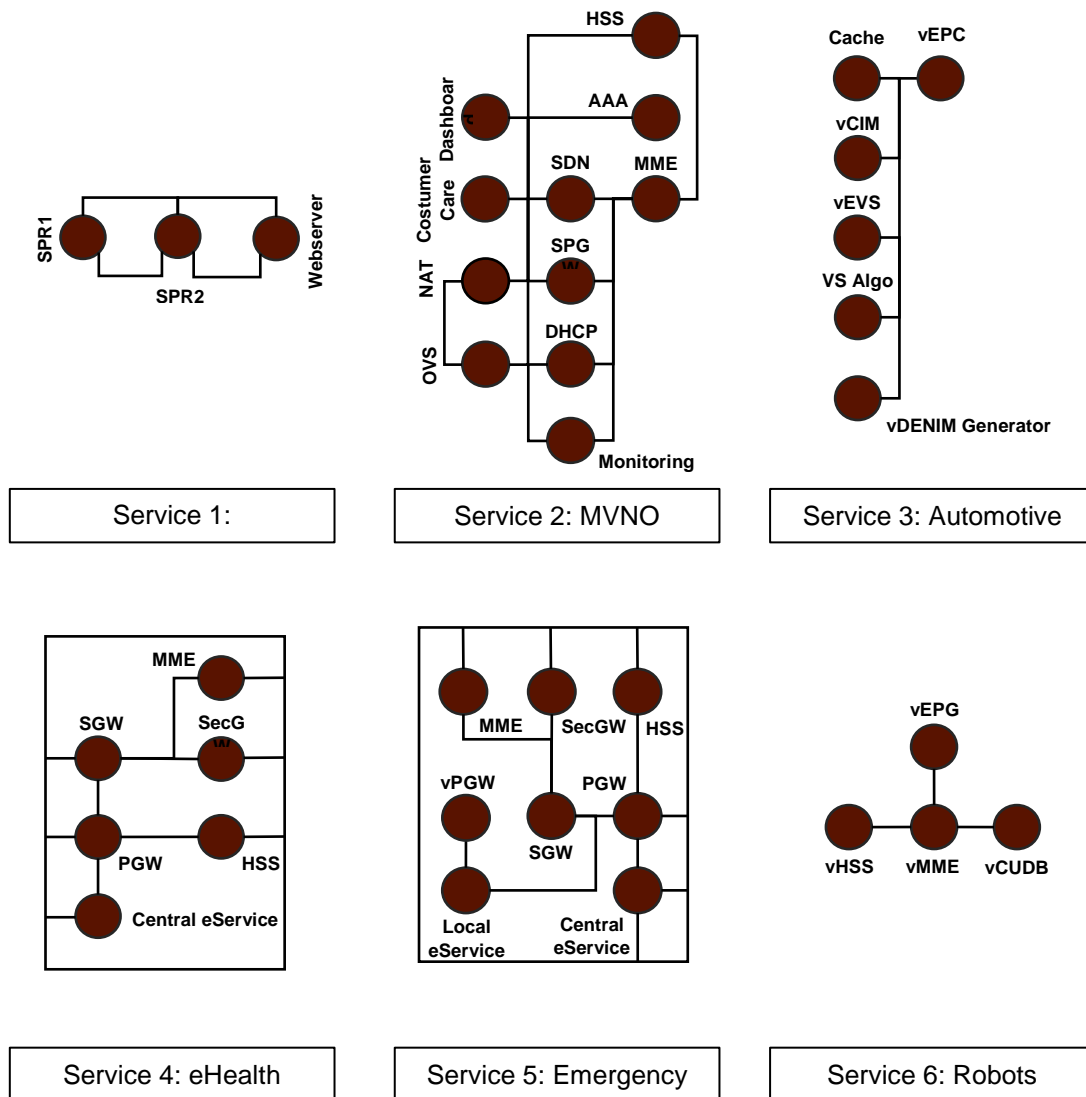


FIGURE 23: SERVICE-GRAPHS

In Table 52, the VNFs contained in each service can be found. For each VNF, the number of CPUs, required RAM, required disk and number of licenses associated to them are provided.

TABLE 52: VNF DEFINITIONS

Service 1: Entertainment				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
SPR1	2	2	50	30
SPR2	1	1	5	30
WebServer	1	1	5	30
Service 2: MVNO				
VNF	CPUs	RAM [GB]	Disk [GB]	Licenses
MME	1	1	10	1000
HSS	1	1	10	1000
AAA	1	1	10	1000
DHCP	1	1	10	1000

S/PGW-C	1	1	10	1000
OVS	1	1	10	1000
NAT	1	1	10	1000
SDN	2	2	20	1000
Dashboard	1	1	10	1000
Customer Care	1	1	10	1000
Monitoring	1	1	10	1000
Service 3: Automotive				
VNF	CPU	RAM [GB]	Disk [GB]	Licenses
vCIM	1	1	10	4000
vEVS	4	8	40	4000
vDENMgenerator	1	2	15	4000
vEPC	1	2	20	4000
Cache	4	8	30	4000
VS Algo	4	8	30	4000
Service 4: Non-Emergency eHealth				
VNF	CPU	RAM [GB]	Disk [GB]	Licenses
SecGW	2	4	40	1
MME	2	4	20	1
HSS	1	2	20	1
SGW	2	4	20	1
PGW	2	4	20	1
Central eService	2	4	100	1
Service 5: Emergency eHealth				
VNF	CPU	RAM [GB]	Disk [GB]	Licenses
SecGW	2	4	40	1
MME	2	4	20	1
HSS	1	2	20	1
SGW	2	4	20	1
PGW	2	4	20	1
Central eService	2	4	100	1
vPGW	2	4	20	1
Local eServer	4	8	20	1
Service 6: Robots				
VNF	CPU	RAM [GB]	Disk [GB]	Licenses
vEPG	20	24	120	1
vMME	10	19	160	1
vCUDB	4	12	100	1
vHSS	4	24	200	1

In Table 53, the links definitions for the connections of each VNF contained in each of the 5G-T like services can be found:

TABLE 53: LINKS DEFINITIONS

Service 1: Entertainment		
Link	BW [Mbps]	Latency [ms]
Management Network	10	50
External Network	250	50
Data Network	10	50
Service 2: MVNO		
Link	Latency [ms]	Latency [ms]

wef_secure	100	50
wef_core_control	100	50
wef_mgmt	100	50
wef_os_sgi	1000	50
Service 3: Automotive		
Link	Latency [ms]	Latency [ms]
int	50	5
Service 4: Non-Emergency eHealth		
Link	Latency [ms]	Latency [ms]
net_d	100	10
net_b	100	10
mgmt	10	50
net_a	10	50
net_gw	10	50
Service 5: Emergency eHealth		
Link	Latency [ms]	Latency [ms]
net_d	100	10
net_b	100	10
mgmt	10	50
net_a	10	50
net_gw	10	50
local_net	100	10
Service 6: Robotics		
Link	Latency [ms]	Latency [ms]
int_a	100	10
int_b	100	10
int_c	100	10
int_d	100	10

For each service the INPUT bandwidths are:

- *Service 1:* 250Mbps
- *Service 2:* 100Mbps
- *Service 3:* 1Gbps
- *Service 4:* 100Mbps
- *Service 5:* 100Mbps
- *Service 6:* 1Gbps

For each service the OUTPUT bandwidths are:

- *Service 1:* 10Mbps
- *Service 2:* 1Gbps
- *Service 3:* 1Gbps
- *Service 4:* 10Mbps
- *Service 5:* 10Mbps
- *Service 6:* 1Gbps

Finally, there are special requirements to some of the services:

- *Service 1:* This service requires protection.
- *Service 4:* This service requires datacentre affinity placing MME, HSS, SGW and PGW in one domain and Central eService in a different datacentre. All the VNFs must remain inside the same domain.

- *Service 5*: This service requires domain affinity, placing MME, HSS, SGW, PGW and Central eService in one domain and vPGW and Local eService in a different in a different domain.

2.3.3.2 Datacentre modelling

Three types of datacentres are proposed:

- *Small datacentres*: This type of datacentre is located at the edge of the network, near to the user. Because of that, the latency that it can guarantee is very low. The resources of the small datacentre, however, are more limited.
- *Medium datacentres*: This type of datacentre is located between the core and the edge of the network. The latency it can guarantee increases with respect to the small datacentre, but its resources will also increase.
- *Large datacentres*: This type of datacentre is located at core of the network. It provides the higher value of latency as well as the higher availability of resources.

In Table 54, the most relevant datacentre-parameters for the carried-out MATLAB simulation in this section can be found.

TABLE 54: DATACENTRES DEFINITION ACCORDING TO 2.2.2

	CPU	RAM [GB]	Disk [GB]	Bandwidth [Gbps]	Guaranteed Latency [ms]
Small	576	1512	40960	1008	5
Medium	3456	9072	112640	1008	10
Large	86112	226044	880640	17990	30

2.3.3.3 Domain modelling

Domains in a federation situation can be either a local domain or an overflow domain. Therefore, in 5G-TRANSFORMER there will be one local domain and several overflow domains.

When deploying a service in an overflow domain, a latency increase of 5ms is considered to account the physical latency among domains.

2.3.4 Tool description

The tool description is divided in two parts:

- The event generator which characterizes the complete life cycle of a 5G-TRANSFORMER like service.
- The pricing which analyses in detail the cost, revenue and profit of the scenario described in 2.3.

2.3.4.1 Event generator

The event generator intends to be a tool that allows the project to simulate arrivals of the different 5G-TRANSFORMER service-like events as described in 2.3.3.

In order to delimit the temporal timespan of the MATLAB simulation, a timeframe of one year with a resolution of one hour was defined. In the experimental analysis, three types of events are distinguished: Deployment event, scale event and end event. Their description is developed in the following subsections.

2.3.4.1.1 Deployment event

The deployment event refers to the arrival of a service demand. This kind of event is characterized by a Poisson probability distribution. This probability distribution gives the probability of occurrence of events $P(t; \lambda)$ in a fixed interval of time T provided that these events occur with a known constant rate λ and with independency of the time t since the last event occurred, as indicated by Equation 181.

$$P(t; \lambda) = \frac{e^{-\lambda} \cdot \lambda^t}{t!}$$

EQUATION 181: EVENT ARRIVAL DISTRIBUTION

Where:

- $\lambda \in \mathbb{R}^+$ is the service frequency depending on each particular 5G-TRANSFORMER service like.
- $t \in \mathbb{R}^+$ is the time instant within the fixed interval of time T with a resolution of one hour.
- $T \in \mathbb{R}^+$ is a fixed interval of time, set in this study to one year.

In 5G-TRANSFORMER, VNFs are considered as an atomic unit. This means that VNFs are taken as a single set of resources and connectivity assets. In order to keep that atomicity, when assigning VNFs to a domain, their set of resources cannot be split and must remain inside the same datacentre.

As a result of the previous idea, VNFs can be tracked from the demands structure to the service module and the resource module.

The deployment event flowchart is depicted in Figure 29:

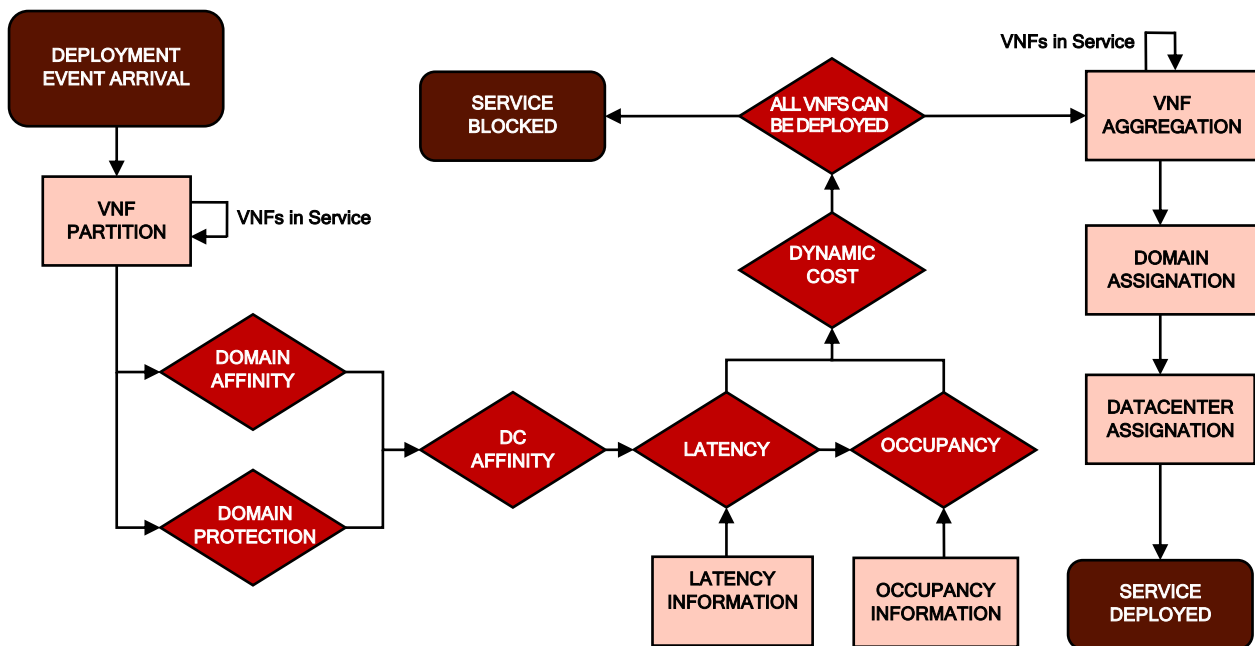


FIGURE 24: DEPLOYMENT EVENT FLOWCHART

When a deployment event arrives, the VNF partition module partitions the service in the VNFs that compose it.

Once the partition is done, it is checked whether the VNFs must comply with domain affinity of domain protection.

The next step is to check whether the VNFs have any datacentre affinity constrain before checking the compliance with the latency and resource requirements.

If all the previous tests are passed for all the service VNFs, the service is deployed by selecting the lowest cost datacentre for each of the service VNFs. Otherwise, the service deployment is blocked.

Hereafter, each of the tests that must be passed by each VNF in order to deploy a service are explained in detail.

2.3.4.1.1.1 Domain affinity

Domain affinity refers to the situation in which for the correct operation of the service, some VNFs have to be assigned to one domain while the other VNFs have to be assigned to another domain.

This affinity model is generally used to ensure for that management and data planes are placed in different domains

In the 5G-TRANSFORMER ecosystem, the Emergency eHealth use case requires domain affinity for its correct operation.

2.3.4.1.1.2 Datacentre affinity

Datacentre affinity refers to the situation in which for the correct operation of the service, some VNFs have to be assigned to one datacentre while the other VNFs have to be assigned to another datacentre.

In the 5G-TRANSFORMER ecosystem, the Non-Emergency eHealth use case requires domain affinity for its correct operation.

2.3.4.1.1.3 Service protection

Service protection implies the deployment of all the VNFs within a service in two different domains. This protection model is used to highly demanding applications in which availability must be ensure at all the time. In the 5G-TRANSFORMER ecosystem, the entertainment use case requires service protection for its correct operation.

2.3.4.1.1.4 Latency constrains

In a service, each VNF has a minimum latency for the correct operation of the service. This latency constrains must be met for all the VNFs in order to make the service deployment possible.

2.3.4.1.1.5 Occupancy constrains

When deploying a VNF, the selected datacentre has to have enough resources for all the service VNFs to allow the deployment.

2.3.4.1.1.6 Dynamic cost

In the experimental analysis, a dynamic cost function is used. This function aims to select the most convenient domain when federating a service by offering different rates according to location and occupation. The dynamic cost variation can be found in Figure 25:

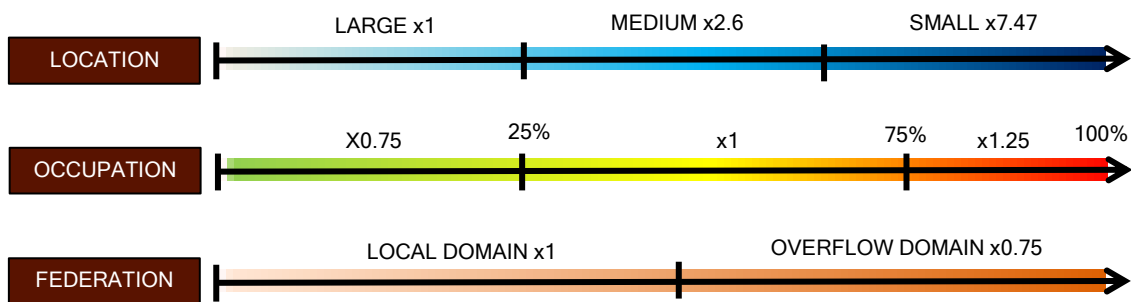


FIGURE 25: DYNAMIC COST VARIATION

The dynamic cost variation depends on:

- *Location:* When a VNF is placed in a large datacentre the cost of deployment is lower than when is placed in a medium datacentre and at the same time it is lower than when is placed in small datacentre. This is due to the fact that the larger the datacentre is, the lower the cost of operation and maintenance is. According to the datacentre prices that will be presented in Table 56, an increment of 2.6 is foreseen in the medium datacentre and of 7.47 in the small datacentre.
- *Occupation:* When federation takes place, overflow domains apply dynamic rates in order to attract the demand when the occupation is low and only accept the most profitable services when it is high. This fact is achieved by reducing prices up to 0.25 when the occupancy level is lower than 25% and increasing the prices by 0.25 when the occupancy level is higher than 75%.
- *Federation:* When a VNF is federated, the price that an overflow domain charges the local domain for using its infrastructure will be 0.75 times the market price. This is a wholesale price that is charged only charged among domains and includes the connectivity cost that the federation could have.

2.3.4.1.2 Scale event

The service scale event refers to the increase of resources that is required by a service, once it has been deployed. The scale event flowchart can be found in Figure 26:

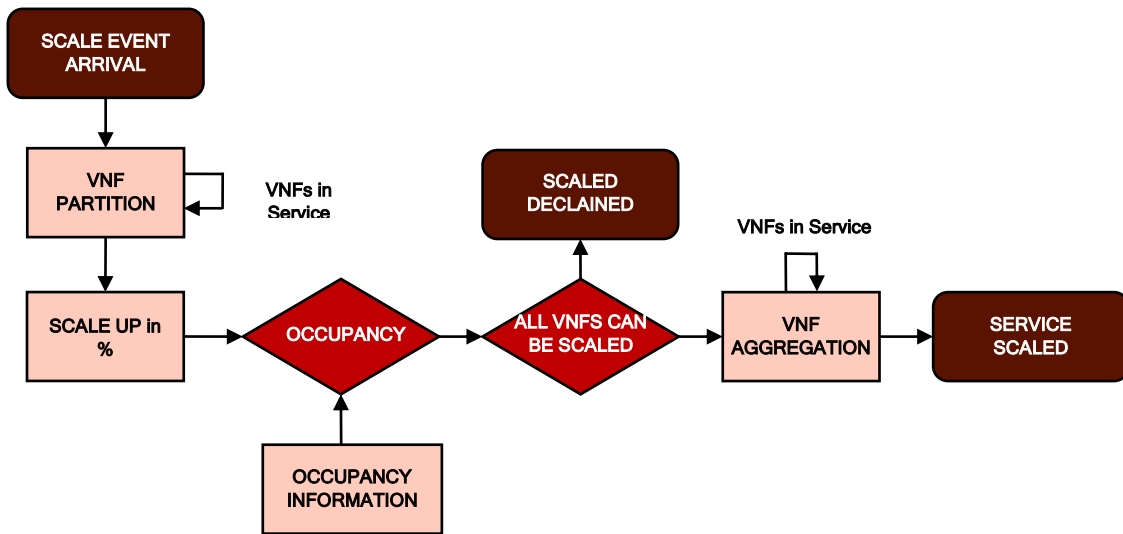


FIGURE 26: SCALE EVENT FLOWCHART

When a deployment event arrives, the VNF partition module partitions the service in the VNFs that compose it. Once the partition is done, the increase in resources for each VNF is calculated and it is checked in terms of occupancy of the current VNF datacentre whether it is possible to assign extra resources to the current VNF in the same datacentre. If the previous test passes for all the service VNFS, the service is scaled up by the required percentage. Otherwise, the service scaling is declined.

2.3.4.1.3 End event

The end event refers to the finalization of a deployed service. The service life-time indicates this finalization, and it is characterized by a truncated Gaussian distribution with mean μ and standard deviation σ^2 that lies within the interval (a, b) , with $-\infty < a \leq b < \infty$ as indicated by Equation 182:

$$P(x; \mu, \sigma, a, b) = \begin{cases} \frac{\phi\left(\frac{x-\mu}{\sigma}\right)}{\sigma\left(\Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)\right)} & a \leq x \leq b \\ 0 & \text{Otherwise} \end{cases}$$

EQUATION 182: LIFE-TIME DISTRIBUTION

Where:

- $\mu \in \mathbb{R}^+$ is the mean which corresponds to the average life-time from Table 44.
- $\sigma^2 \in \mathbb{R}^+$ is the standard deviation which is set to 1 in order to ensure a small deviation from the mean.
- $\phi(\xi) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\xi^2}$ is the probability density function of the Gaussian distribution.
- $\Phi(\xi) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\xi}{\sqrt{2}}\right)\right)$ is the probability density function of the Gaussian distribution.

The flowchart for the end event can be found in Figure 27:

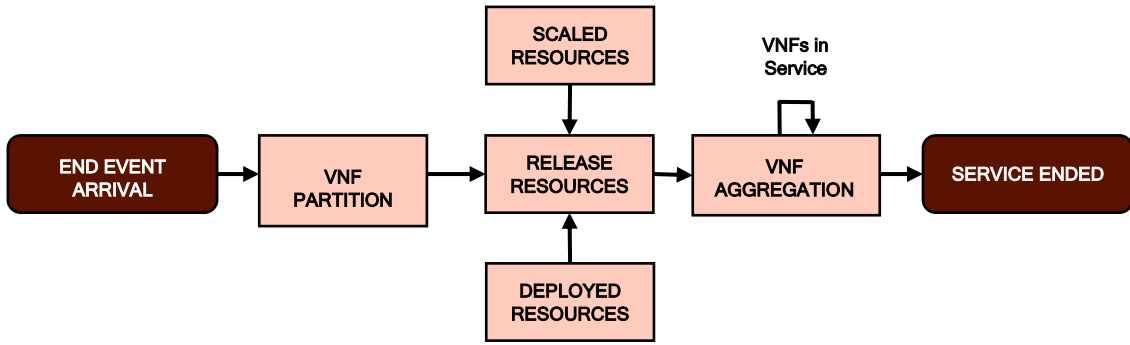


FIGURE 27: END EVENT FLOWCHART

When an end event arrives, the VNF partition module partitions the service in the VNFs that compose it. Once the partition is done, both the deployed and scaled resources are released, and their availability is notified by the domain to federation.

2.3.4.2 Blocked services

Blocked services occurred when the scarcity of resources prevent the deployment of a particular service. Each overflow domain participating in the federation makes an announcement indicating the minimum latency that it can ensure. Thanks to that announcement, when federating a service, the local domain knows which VNFs can be assigned to which of the overflow domains. This communication aims to guarantee that the local domain is able to know the current resource availability of the datacentres in each overflow domain and therefore assign that VNF to one that is able to deploy the service. Provided that the previous announcement takes place, the available resources R_t at the time t for federating a service's VNF can be calculated with Equation 183:

$$R_t = \begin{cases} \sum_{k=1}^K \sum_{n=1}^{N_k} \sum_{t=0}^T (C_{DC_n} - R_{t-1})_k, & R_{t-1} < C_{DC_n} \\ 0, & R_{t-1} \geq C_{DC_n} \end{cases}$$

EQUATION 183: AVAILABLE RESOURCES

Where:

- $R_t \in \mathbb{R}^+$ is the available resources at the time t for federating a particular service VNF.
- $k \in \mathbb{N}$ corresponds to the overflow domain of $K \in \mathbb{N}$ overflow domains that participate in the federation.
- $n \in \mathbb{N}$ corresponds to the datacentre of the $N_k \in \mathbb{N}$ datacenters that complies with the latency constraints of the service that is being deployed.
- $C_{DC} \in \mathbb{R}^+$ is the capacity of a datacentre.

When $R_t = 0$, the service cannot be deployed, and the datacenter enters into a blocking service state until some other service ends and releases resources and is unable to deploy any further VNF. The blocking service state can be reached by the occurrence of different situations as described hereafter:

- *Domain affinity*: When domain affinity is required, it could happen that there are not enough domains to comply with the domain affinity requirements or that those have not enough resources to deploy the service VNFs. In this case, the service is blocked and report as domain affinity blocking cause.
- *Domain protection*: In a similar way to domain affinity, when domain protection is required it could be the case that none of the other available domains can deploy a copy of the current service. In this case, the service is blocked and reported as a domain protection blocking cause.
- *Datacentre affinity*: When datacentre affinity is required, in a similar way to domain affinity, it could happen that there are not enough datacentres available within the domain or that those available have not enough resources to deploy the service VNFs. In this case, the service is blocked and report as datacentre affinity blocking cause.
- *Latency constrains*: When the latency constraints of the service VNFs are not met, then the service cannot be deployed, and it is blocked and reported as a latency blocking cause.
- *Occupancy constrains*: When the datacentres that comply with all the previous requirements have not enough capacity to deploy the service, then it is blocked and reported as occupancy blocking cause.

In the simulation, blocked services are stored together to the cause that blocked the service so that it can be accounted.

2.3.4.3 Not-scaled services

The number of events that are scaled up are selected by a half-gaussian distribution defined in Equation 184:

$$P(x; \sigma) = \frac{\sqrt{2}}{\sigma\sqrt{\pi}} e^{-\frac{x^2}{2\sigma^2}}$$

EQUATION 184: NUMBER OF SCALED EVENTS DISTRIBUTION

Where:

- $\sigma^2 \in \mathbb{R}^+$ is the standard deviation which is set to 1 in order to ensure a small deviation from the mean.

In addition, the number of resources that are scaled inside an event proposed to being scaled has to be defined. For such a purpose, a percentage of the resources that are required for a particular service is used. This percentage is selected by using a uniform distribution defined between a and b as indicated by Equation 185:

$$P(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{Otherwise} \end{cases}$$

EQUATION 185: SCALED RESOURCES DISTRIBUTION

Where:

- $a \in \mathbb{R}^+$ is lower limit set to 0%.
- $b \in \mathbb{R}^+$ is lower limit set to 100%.

Not-scaled services occurred when the scarcity of resources prevents the scaling-up of a particular service VNF inside the same datacentre it is already deployed. The analysis on the available resources and the reasons why a service could or could not be scaled up are those already presented in 2.3.4.2.

2.3.5 Scenario setting

The setting will vary according to the demand of the scenario that is being considered: pessimistic, realistic, optimistic, ideal or overflow scenario. A more detailed description of these scenarios and how the demand varies and affects the service prices is detailed in the following subsections. However, there is a common configuration for all of them that it is described in the following subsection.

2.3.5.1 Common configuration

The common configuration is composed by domains, datacentres and the total CAPEX and OPEX for a year timeframe per datacentre.

2.3.5.1.1 Domains

For all the scenarios three different domains are considered, one local domain (Domain 1) and two overflow domains (Domain 2 and Domain 3). With this configuration it can be guaranteed that competition among the overflow domains participating in the federation takes place.

2.3.5.1.2 Datacentres distribution across domains

The datacentres distribution across the previous domains can be found in Table 55. This configuration aims to be as close to reality as possible, by representing different configurations of datacentres.

TABLE 55: DATACENTRES DISTRIBUTION ACROSS DOMAINS

	Small	Medium	Large
Domain 1	1	1	1
Domain 2	1	1	0
Domain 3	0	1	0

2.3.5.1.3 Datacentres cost

The total CAPEX and OPEX in a year are calculated by using the data provided in the infrastructure cost modelling in Table 56:

TABLE 56: TOTAL CAPEX+OPEX

	CAPEX	OPEX	Total
Small	\$62.458,39	\$197.708,44	\$260.166,83
Medium	\$145.290,43	\$397.071,46	\$542.361,89
Large	\$2.625.579,01	\$2.578.037,67	\$5.203.616,68

2.3.5.1.4 Prices for licenses, transactions and special requirements

The prices are shown in Table 57:

TABLE 57: PRICES FOR LICENSES, TRANSACTIONS AND SPECIAL REQUIREMENTS

	License	Transaction	Special requirement		
			Domain Affinity	DC Affinity	Protection
Price	\$10	\$25	\$100	\$50	\$150

2.3.5.1.5 Number of scaled up events and grade of scaling

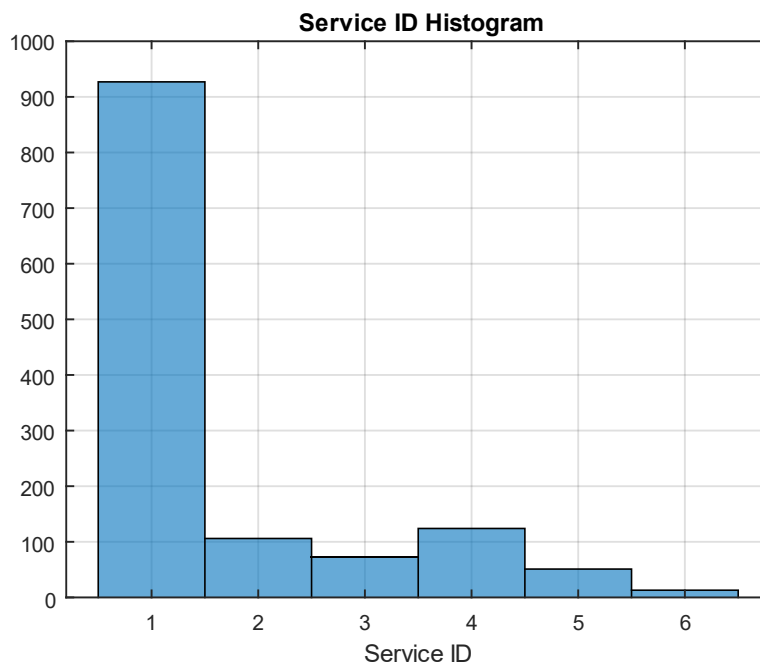
The number of scaled-up events will be the 30% and the grade of scaling will be the 50% of the total number of resources.

2.3.5.2 Pessimistic scenario

The pessimistic scenario supposes that the infrastructure provided by the local domain is utilized at 20% and no infrastructure from the overflow domains is used. Hereafter, the generated events, datacentre states, interchanged bandwidth among domains, blocked-services, not-scaled services and the price analysis for this scenario can be found.

2.3.5.2.1 Generated events

For the pessimistic scenario the number of generated events is 2710 events during a year timeframe. In Figure 28, the Service ID Histogram indicates how many events of each of the different 5G-TRANSFORMER service-like are generated in comparison with the other events can be found:

**FIGURE 28: PESSIMISTIC SERVICE ID HISTOGRAM**

In Figure 29, the Life-Time Histogram providing a general view of how much time last the services once they are deployed can be found:

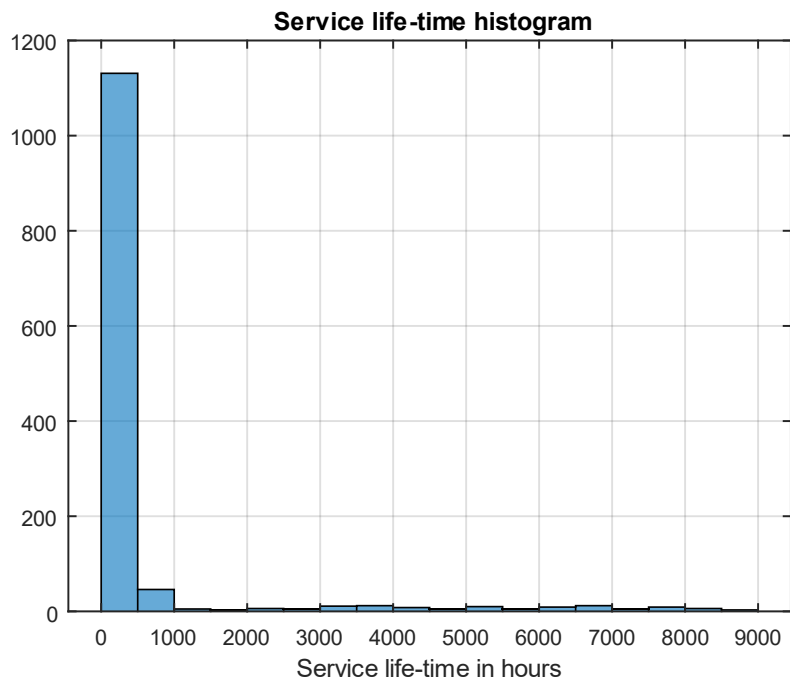


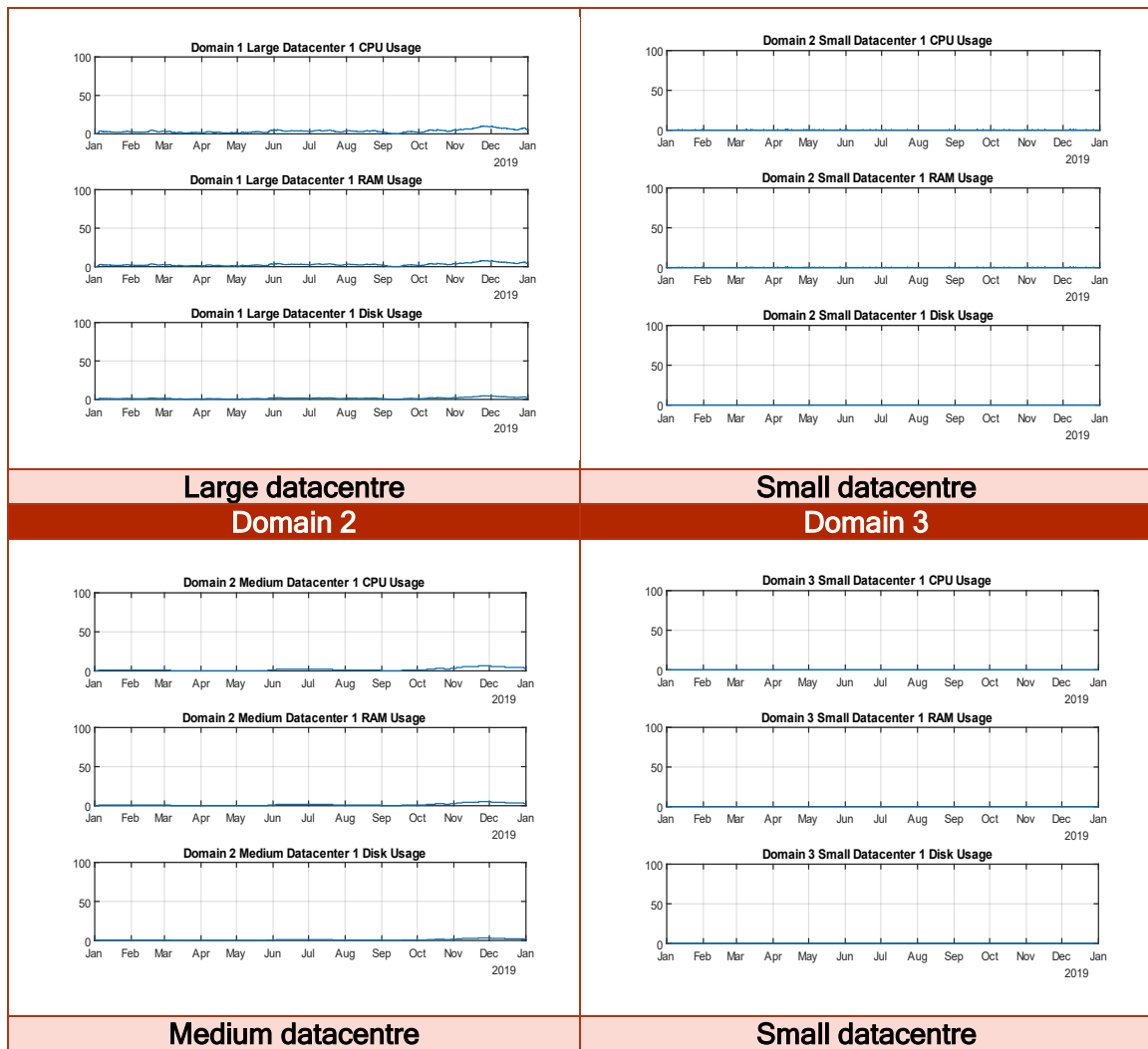
FIGURE 29: PESSIMISTIC SERVICE LIFE-TIME HISTOGRAM

2.3.5.2.2 Datacentres state

The state of the datacentres for the pessimistic scenario for the local domain can be found in Table 58:

TABLE 58: PESSIMISTIC DATACENTRES STATE

Domain 1	Domain 1
<p>Domain 1 Small Datacenter 1 CPU Usage</p>	<p>Domain 1 Medium Datacenter 1 CPU Usage</p>
<p>Domain 1 Small Datacenter 1 RAM Usage</p>	<p>Domain 1 Medium Datacenter 1 RAM Usage</p>
<p>Domain 1 Small Datacenter 1 Disk Usage</p>	<p>Domain 1 Medium Datacenter 1 Disk Usage</p>
Small datacentre Domain 1	Medium datacentre Domain 2

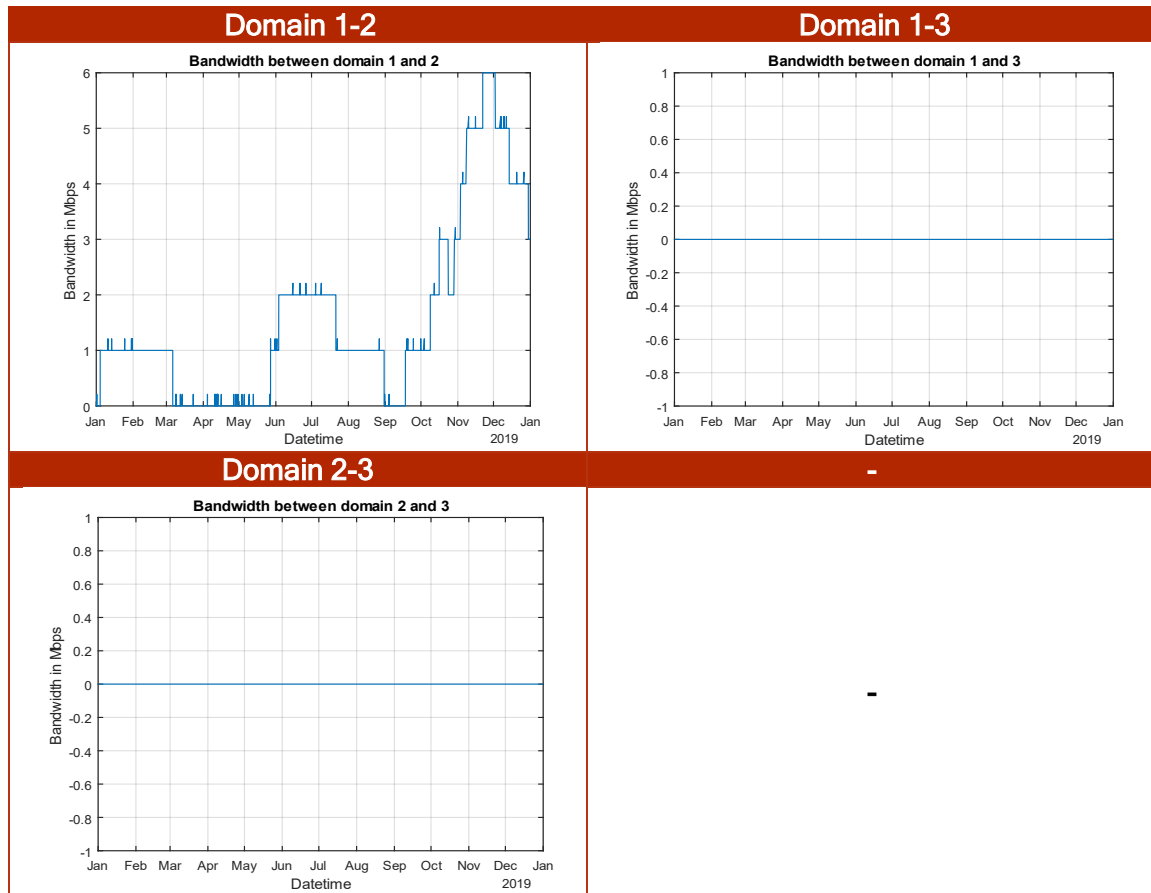


It would have been expected that all the local datacentres were occupied up to 20% and that the overflow domains were not occupied as the demand was designed for that purpose. However, the deployment event flowchart shown in 2.3.4.1.1 manages the demand in a more complex way. The deployment event flowchart places the VNFs of a particular service to the bigger available datacentre that complies with all the latency constrains if possible, in the local domain and otherwise in one of the overflow domains. The occupation that arises in the overflow domains is due to the special requirements of the entertainment service, which requires protection and will deploy the same service in the local domain and an additional domain, and the Emergency eHealth service, which requires domain affinity and will place all its VNFs in the local domain except two VNFs that are place in an overflow domain.

2.3.5.2.3 Interchanged bandwidth among domains

In the pessimistic scenario the interchanged bandwidth among domains is that one coming from the special requirements. As there’s enough vacancy in domain 2 which have a medium datacentre and therefore offers a more convenient price than domain 3 that has only a small datacentre, all the bandwidth is interchanged by domain 1 and 2 as it is shown here:

TABLE 59: PESSIMISTIC INTERCHANGE BANDWIDTH AMONG DOMAINS



2.3.5.2.4 Blocked services

In the pessimistic scenario there are no blocked services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.2.5 Not-Scaled services

In the pessimistic scenario there are no not-scaled services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.2.6 Price analysis

Price is defined by Equation 162 which is dependent on the profit margin and the breakeven price which is defined by Equation 161 and dependant the cost of this infrastructure, the forecast utilization of the local infrastructure, and the real utilization of the local infrastructure. According to what it was stated in 2.2.3.5.1, several combinations of the previous variables are possible and depending on how they are combined the profit will be greater or lower as it was shown in Figure 16. For the experimental analysis, the following values are selected for the MATLAB simulations:

TABLE 60: PESSIMISTIC SCENARIO VALUES

	C_l	u_{rl}	u_{fl}	m
Small	260166.83	0.2	0.3	0.8
Medium	542361.89			
Large	5203616.68			

The maximum local demand for the pessimistic scenario will be:

TABLE 61: MAXIMUM LOCAL DEMAND PESSIMISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small $d_{l_{max}}$	144	48	38	52	33	15
Medium $d_{l_{max}}$	864	288	230	314	203	90
Large $d_{l_{max}}$	21528	7176	5740	7828	5065	2266

The real local demand for the pessimistic scenario will be:

TABLE 62: REAL LOCAL DEMAND PESSIMISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small d_{l_r}	28	9	7	10	6	3
Medium d_{l_r}	172	57	46	62	40	18
Large d_{l_r}	4305	1435	1148	1565	1013	453

With the previous parameters, the prices for each of the services of the 5G-TRANSFORMER ecosystem will be:

TABLE 63: PESSIMISTIC SCENARIO PRICES

	S1	S2	S3	S4	S5	S6
Small	\$10890.71	\$33450.03	\$42572.76	\$31220.02	\$52033.37	\$117075.08
Medium	\$3769.32	\$11351.77	\$14148.58	\$10385.66	\$16270.86	\$36157.46
Large	\$1450.38	\$4352.47	\$5439.33	\$3989.15	\$6166.24	\$13794.57

2.3.5.2.6.1 Cost

In the pessimistic scenario, the only cost is the infrastructure cost as indicated in 2.3.5.1.3. This cost is composed by one datacentre of each type. Therefore, at the end of the year the cost will be:

TABLE 64: PESSIMISTIC COST

	Small	Medium	Large	Total
Cost	\$260.166,83	\$542.361,89	\$5.203.616,68	\$6.006.145,4

2.3.5.2.6.2 Revenue

The revenue for the scenario comes from the price that the tenants pay for a service as those indicated in 2.3.5.2.6. In addition, there's also revenue coming from the licenses, transaction for scaling up services and special requirements as in Table 57. The revenue for the pessimistic scenario is shown in Figure 32:

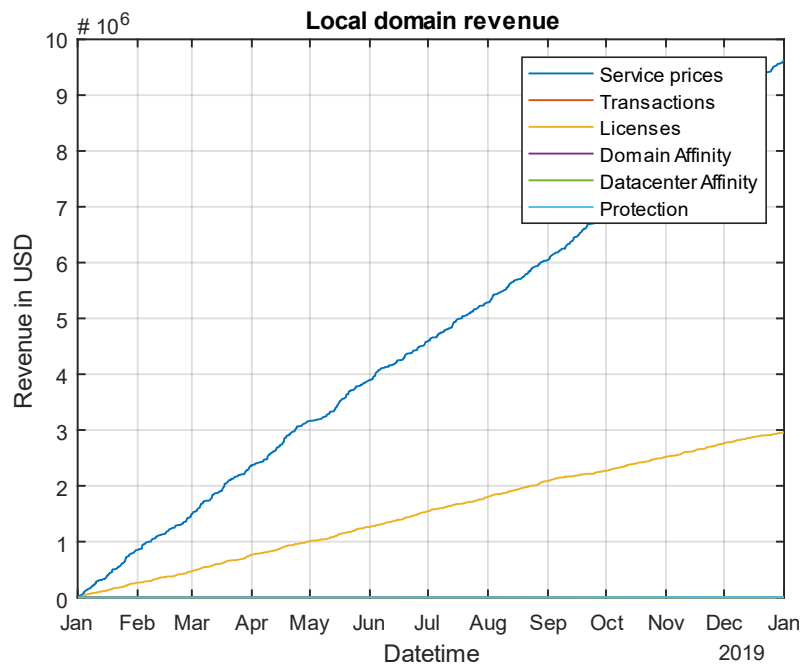


FIGURE 30: PESSIMISTIC REVENUE

2.3.5.2.6.3 Profit

The profit for the realistic scenario is shown in Figure 31:

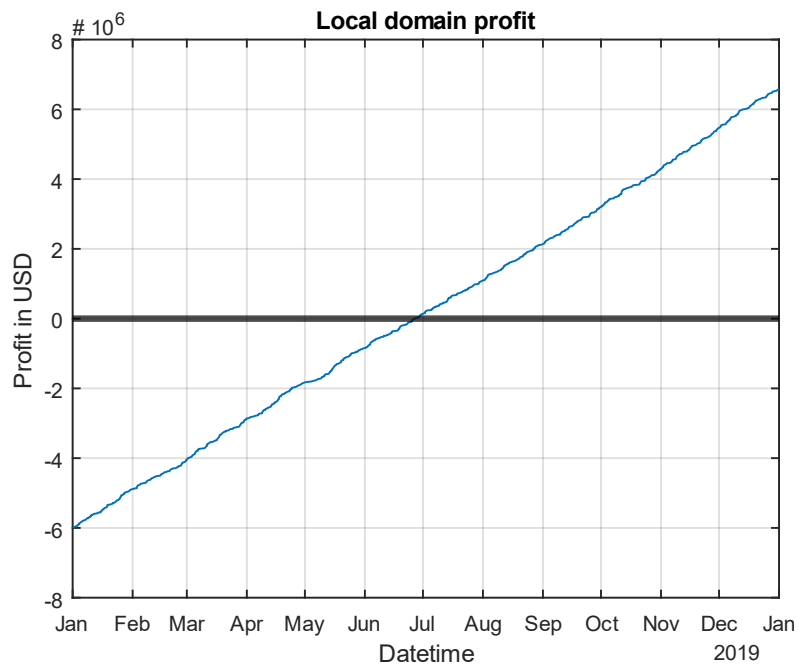


FIGURE 31: PESSIMISTIC PROFIT

2.3.5.3 Realistic scenario

The realistic scenario supposes that the infrastructure provided by the local domain is utilized at 50% and no infrastructure from the overflow domains is used.

Hereafter, the generated events, datacentre states, interchanged bandwidth among domains, blocked-services, not-scaled services and the price analysis for this scenario can be found.

2.3.5.3.1 Generated events

For the pessimistic scenario the number of generated events is 6773 events during a year timeframe. In Figure 32, the Service ID Histogram indicates how many events of each of the different 5G-TRANSFORMER service-like are generated in comparison with the other events can be found:

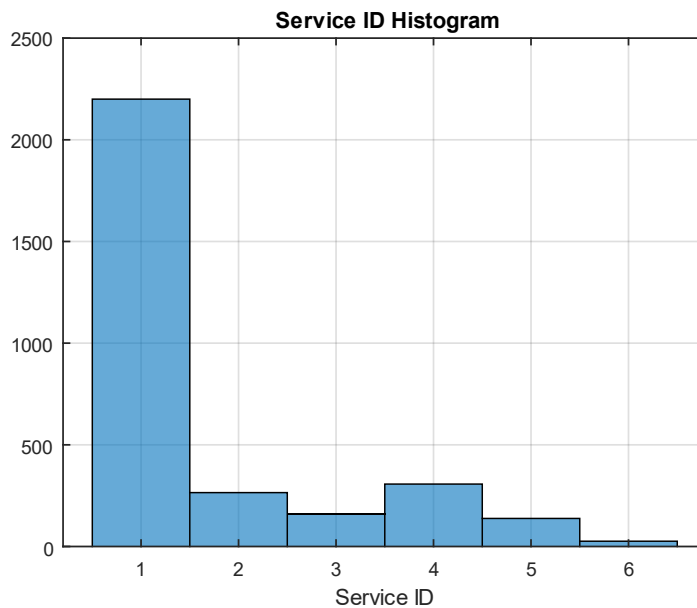


FIGURE 32: REALISTIC SERVICE ID HISTOGRAM

In Figure 33, Life-Time Histogram providing a general view of how much time last the services once they are deployed can be found:

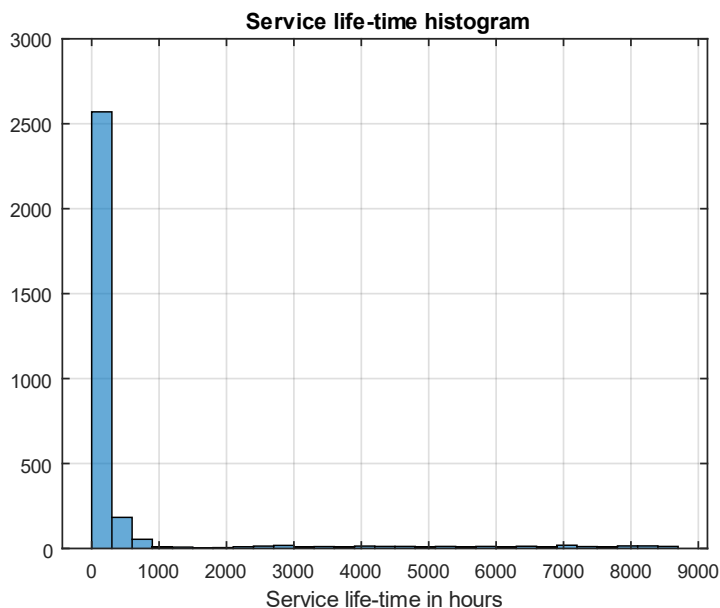
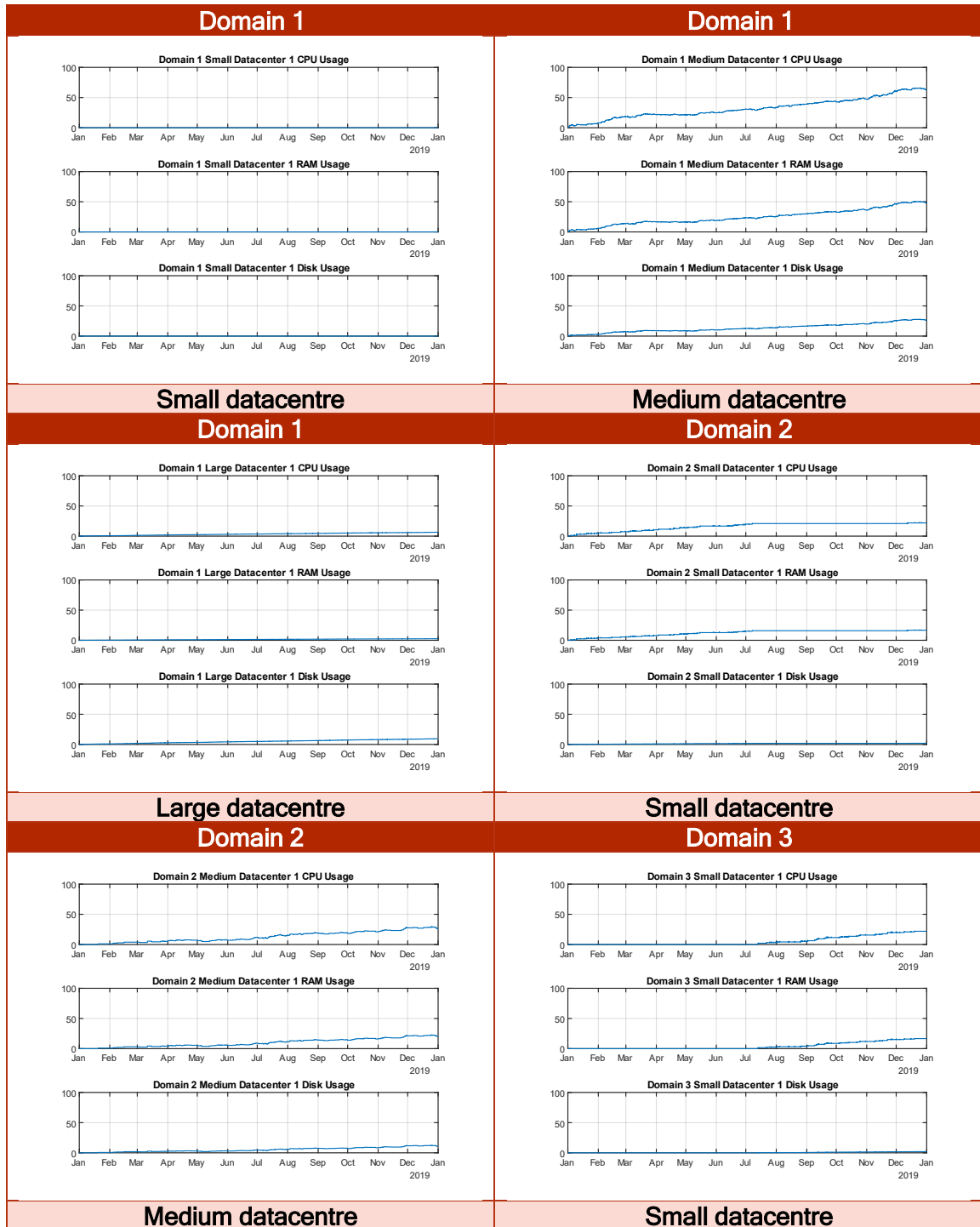


FIGURE 33: REALISTIC SERVICE LIFE-TIME HISTOGRAM

2.3.5.3.2 Datacentres state

The state of the datacentres for the realistic scenario for the local domain can be found in Table 65:

TABLE 65: REALISTIC DATACENTRES STATE



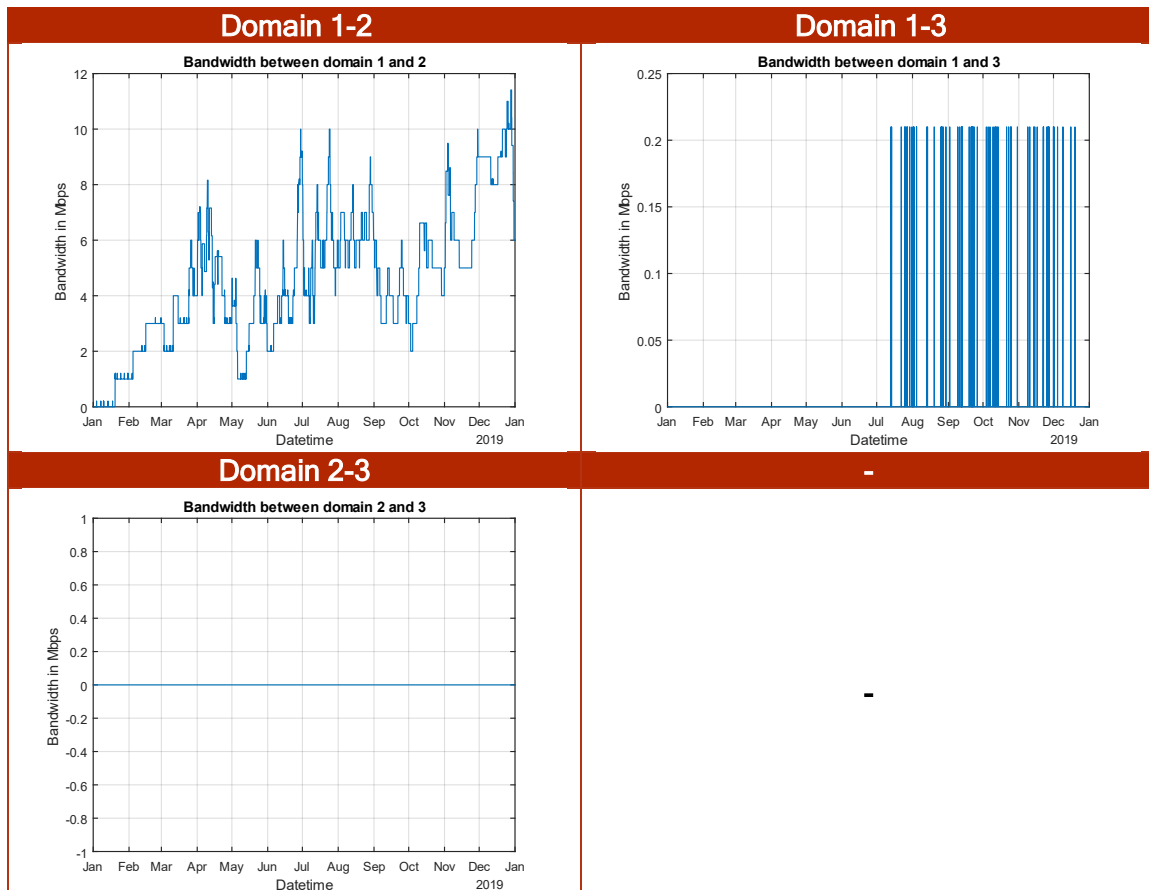
The same as in the pessimistic scenario comments can be applied here. It is important also to note that, as domain 2 rise the use of its the medium datacentre above the 25%

of its capacity, domain 2 and domain 3 small datacentre are also used as it implies a cheaper deployment when below of that 25% even with a smaller infrastructure.

2.3.5.3.3 Interchanged bandwidth among domains

In the realistic scenario the interchanged bandwidth among domains is that one coming from the special requirements. As indicated before, there’s a small bandwidth interchange between domain 1 and 3 that was not seen in the pessimistic scenario, as shown in the following table:

TABLE 66: REALISTIC INTERCHANGE BANDWIDTH AMONG DOMAINS



2.3.5.3.4 Blocked services

In the realistic scenario there are no blocked services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.3.5 Not-Scaled services

In the realistic scenario there are no not-scaled services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.3.6 Price analysis

Price is defined by Equation 162 which is dependent on the profit margin and the breakeven price which is defined by Equation 161 and dependant the cost of this infrastructure, the forecast utilization of the local infrastructure, and the real utilization of the local infrastructure. According to what it was stated in 2.2.3.5.1, several

combinations of the previous variables are possible and depending on how they are combined the profit will be greater or lower as it was shown in Figure 17. For the experimental analysis, the following values are selected for the MATLAB simulations:

TABLE 67: REALISTIC SCENARIO VALUES

	C_l	u_{rl}	u_{fl}	m
Small	260166.83	0.5	0.5	0.3
Medium	542361.89			
Large	5203616.68			

The maximum local demand for the pessimistic scenario will be:

TABLE 68: MAXIMUM LOCAL DEMAND REALISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small $d_{l_{max}}$	144	48	38	52	33	15
Medium $d_{l_{max}}$	864	288	230	314	203	90
Large $d_{l_{max}}$	21528	7176	5740	7828	5065	2266

The real local demand for the pessimistic scenario will be:

TABLE 69: REAL LOCAL DEMAND REALISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small d_{lr}	72	24	19	26	16	7
Medium d_{lr}	432	144	115	157	101	45
Large d_{lr}	10764	3588	2870	3914	2532	1133

With the previous parameters, the prices for each of the services of the 5G-TRANSFORMER ecosystem will be:

TABLE 70: REALISTIC SCENARIO PRICES

	S1	S2	S3	S4	S5	S6
Small	\$4697.46	\$14092.37	\$17800.89	\$13008.35	\$21138.56	\$48316.7
Medium	\$1632.11	\$4896.33	\$6131.05	\$4490.9	\$6980.9	\$15668.24
Large	\$628.46	\$1885.37	\$2357.04	\$1728.34	\$2671.69	\$5970.62

2.3.5.3.6.1 Cost

In the realistic scenario, the only cost is the infrastructure cost as indicated in 2.3.5.1.3. This cost is composed by one datacentre of each type.

Therefore, at the end of the year the cost will be:

TABLE 71: REALISTIC COST

	Small	Medium	Large	Total
Cost	\$260.166,83	\$542.361,89	\$5.203.616,68	\$6.006.145,4

2.3.5.3.6.2 Revenue

The revenue for the scenario comes from the price that the tenants pay for a service as those indicated in 2.3.5.2.6. In addition, there's also revenue coming from the licenses, transaction for scaling up services and special requirements as in Table 57.

The revenue for the realistic scenario is shown in Figure 34:

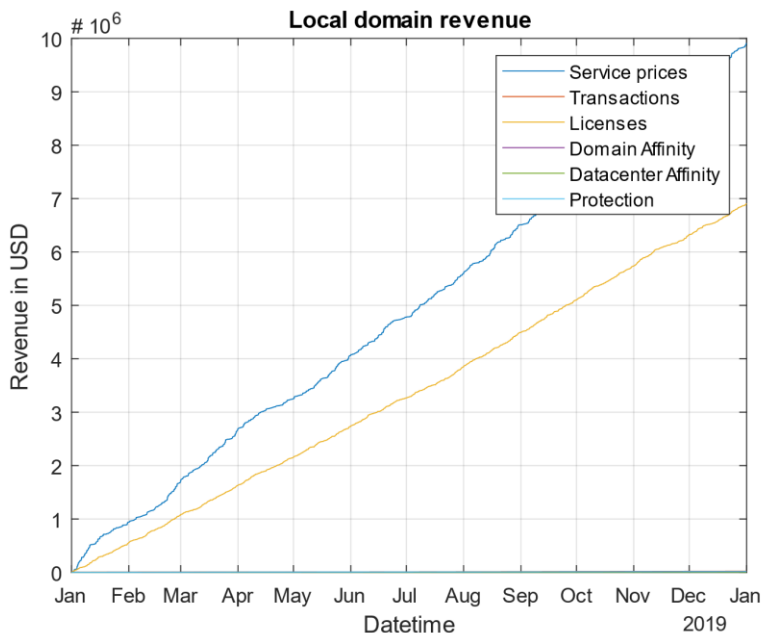


FIGURE 34: REALISTIC REVENUE

2.3.5.3.6.3 Profit

The profit for the realistic scenario is shown in Figure 35:

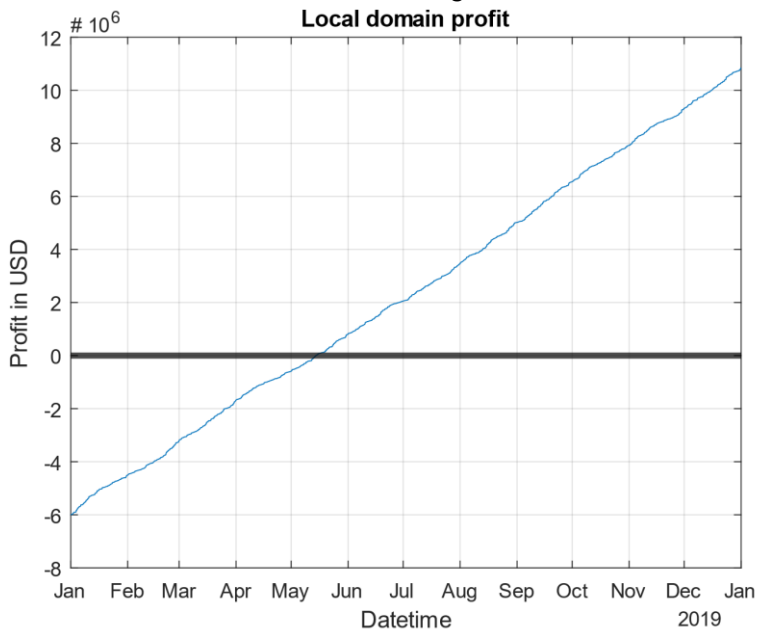


FIGURE 35: REALISTIC PROFIT

The profit is quite similar to that of the pessimistic scenario, even with a smaller margin profit due to the bigger demand and the smaller difference between the real and forecast utilization of the local infrastructure.

2.3.5.4 Optimistic scenario

The optimistic scenario supposes that the infrastructure provided by the local domain is utilized at 80% and no infrastructure from the overflow domains is used.

Hereafter, the generated events, datacentre states, interchanged bandwidth among domains, blocked-services, not-scaled services and the price analysis for this scenario can be found.

2.3.5.4.1 Generated events

For the optimistic scenario the number of generated events is 10840 events during a year timeframe. In Figure 36, the Service ID Histogram indicates how many events of each of the different 5G-TRANSFORMER service-like are generated in comparison with the other events can be found:

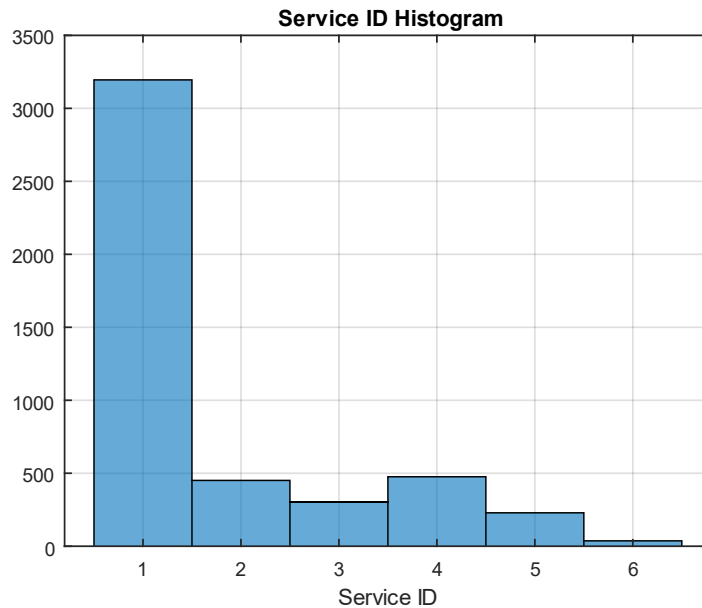


FIGURE 36: OPTIMISTIC SERVICE ID HISTOGRAM

In Figure 37 the Life-Time Histogram providing a general view of how much time last the services once they are deployed can be found:

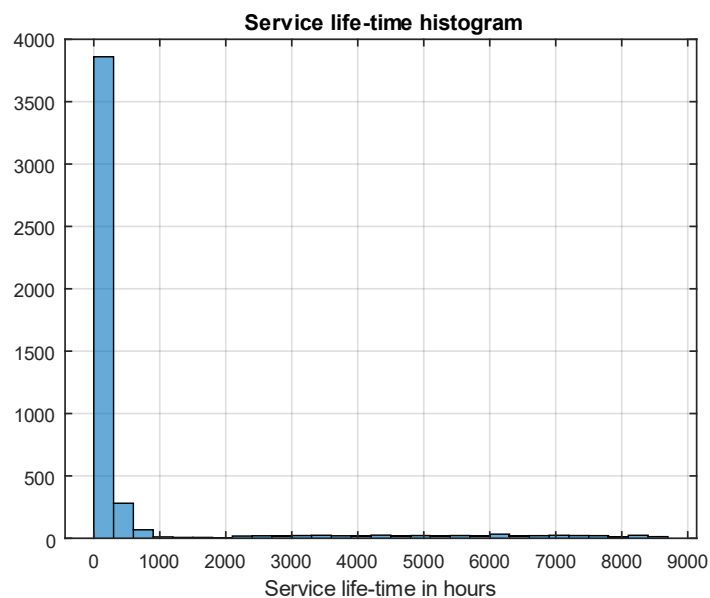
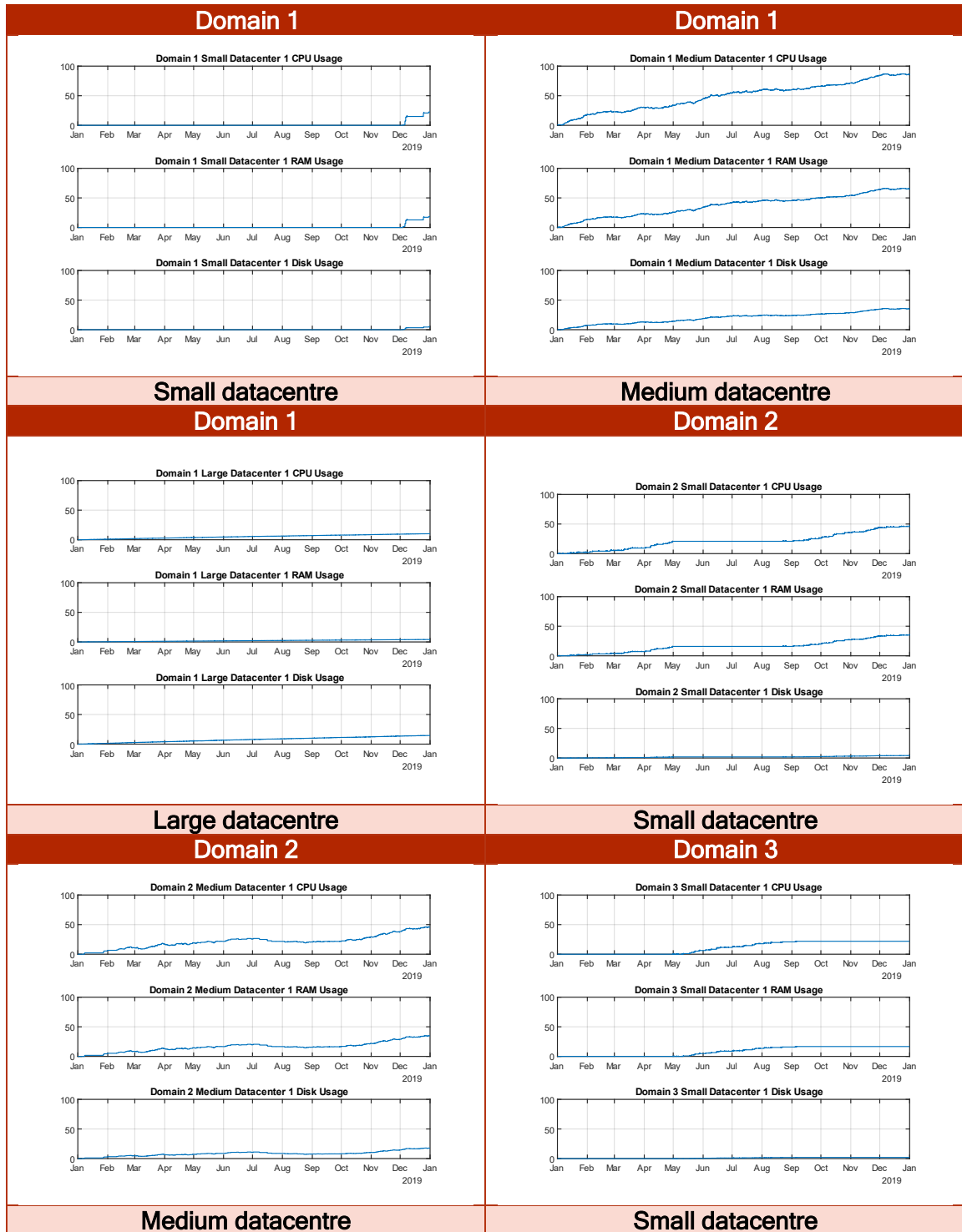


FIGURE 37: OPTIMISTIC SERVICE LIFE-TIME HISTOGRAM

2.3.5.4.2 Datacentres state

The state of the datacentres for the optimistic scenario for the local domain can be found in Table 72:

TABLE 72: OPTIMISTIC DATACENTRES STATE



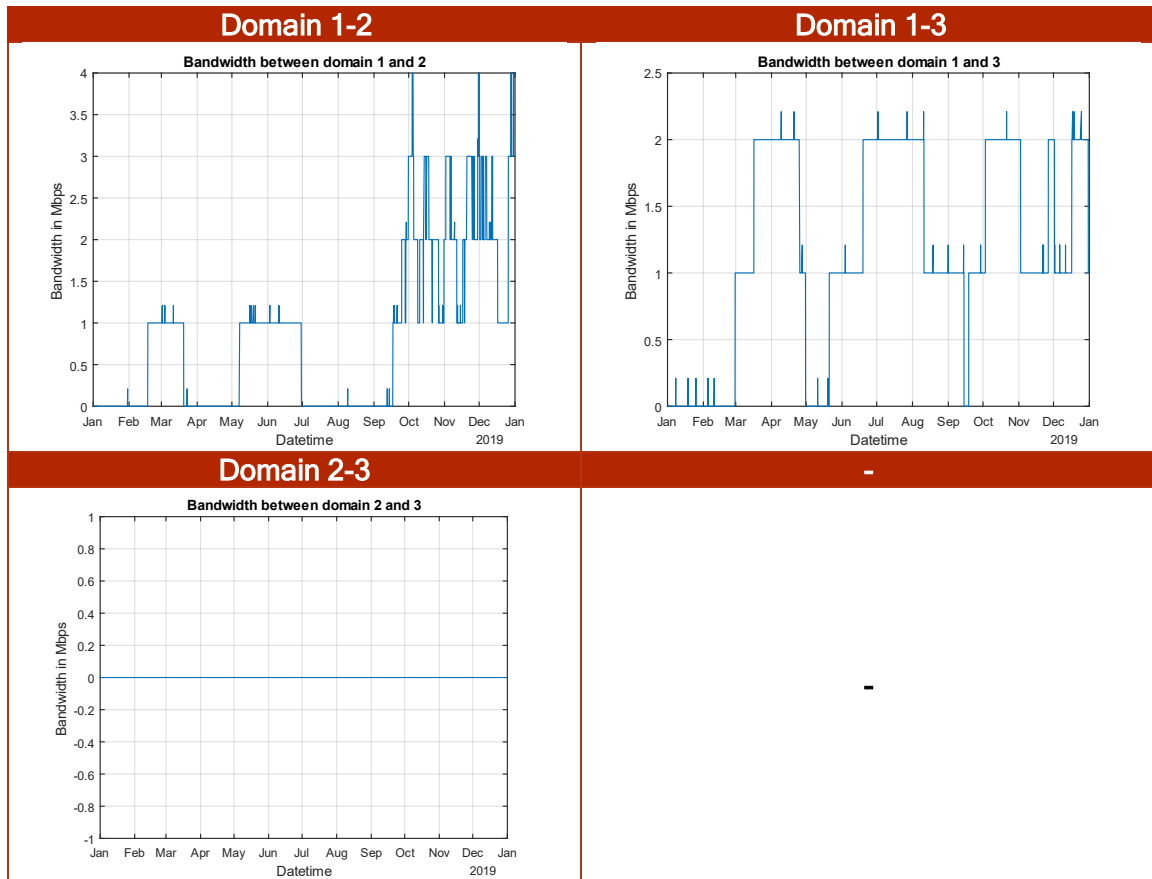
As it can be seen there's a slight difference with the previous scenarios. As the most latency-requiring service demand increments considerably, the use of the overflow

infrastructure is increased due to the impossibility of using bigger datacentres in the federation.

2.3.5.4.3 Interchanged bandwidth among domains

In the optimistic scenario the interchanged bandwidth among domains is that one coming from the special requirements and the federated services. The same comments as in the realistic scenario apply here:

TABLE 73: OPTIMISTIC INTERCHANGE BANDWIDTH AMONG DOMAINS



2.3.5.4.4 Blocked services

In the optimistic scenario there are no blocked services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.4.5 Not-Scaled services

In the optimistic scenario there are no not-scaled services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.4.6 Price analysis

Price is defined by Equation 162 which is dependent on the profit margin and the breakeven price which is defined by Equation 161 and dependant the cost of this infrastructure, the forecast utilization of the local infrastructure, and the real utilization of the local infrastructure.

According to what it was stated in 2.2.3.5.1, several combinations of the previous variables are possible and depending on how they are combined the profit will be greater or lower as it was shown in Figure 18. For the experimental analysis, the following values are selected for the MATLAB simulations:

TABLE 74: OPTIMISTIC SCENARIO VALUES

	C_l	u_{rl}	u_{fl}	m
Small	260166.83	0.8	0.6	0.3
Medium	542361.89			
Large	5203616.68			

The maximum local demand for the pessimistic scenario will be:

TABLE 75: MAXIMUM LOCAL DEMAND OPTIMISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small $d_{l,max}$	144	48	38	52	33	15
Medium $d_{l,max}$	864	288	230	314	203	90
Large $d_{l,max}$	21528	7176	5740	7828	5065	2266

The real local demand for the pessimistic scenario will be:

TABLE 76: REAL LOCAL DEMAND OPTIMISTIC SCENARIO

	S1	S2	S3	S4	S5	S6
Small d_{lr}	115	38	30	41	26	12
Medium d_{lr}	691	230	184	251	162	72
Large d_{lr}	17222	5740	4592	6262	4052	1812

With the previous parameters, the prices for each of the services of the 5G-TRANSFORMER ecosystem will be:

TABLE 77: OPTIMISTIC SCENARIO PRICES

	S1	S2	S3	S4	S5	S6
Small	\$3327.72	\$10220.84	\$13008.35	\$9231.73	\$15062.3	\$31798.17
Medium	\$1151.74	\$3468.6	\$4323.18	\$3173.4	\$4930.57	\$11048.12
Large	\$443.17	\$1329.62	\$1662.02	\$1218.91	\$1883.51	\$4211.91

2.3.5.4.6.1 Cost

In the optimistic scenario, the only cost is the infrastructure cost as indicated in 2.3.5.1.3. This cost is composed by one datacentre of each type.

Therefore, at the end of the year the cost will be:

TABLE 78: OPTIMISTIC COST

	Small	Medium	Large	Total
Cost	\$260.166,83	\$542.361,89	\$5.203.616,68	\$6.006.145,4

2.3.5.4.6.2 Revenue

The revenue for the scenario comes from the price that the tenants pay for a service as those indicated in 2.3.5.2.6. In addition, there's also revenue coming from the licenses, transaction for scaling up services and special requirements as in Table 57.

The revenue for the optimistic scenario is shown in Figure 38:

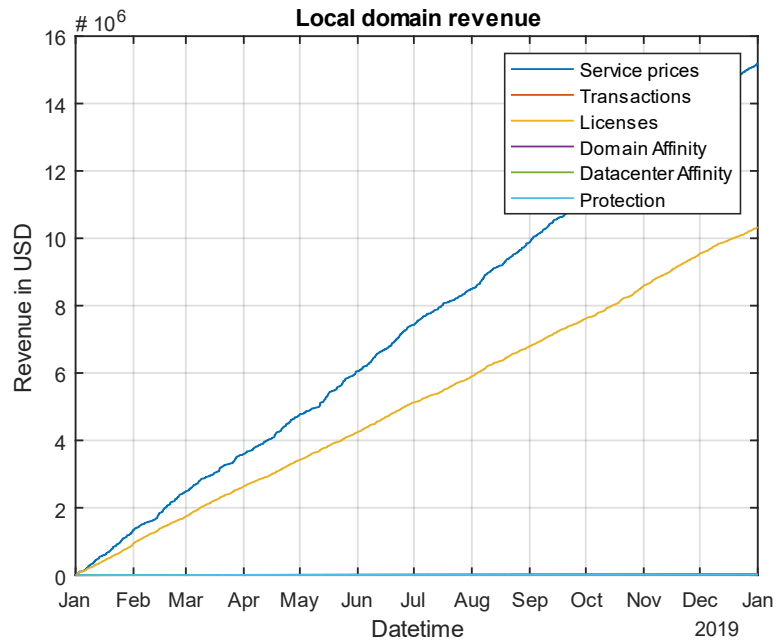


FIGURE 38: OPTIMISTIC REVENUE

2.3.5.4.6.3 Profit

The revenue for the optimistic scenario is shown in Figure 39:

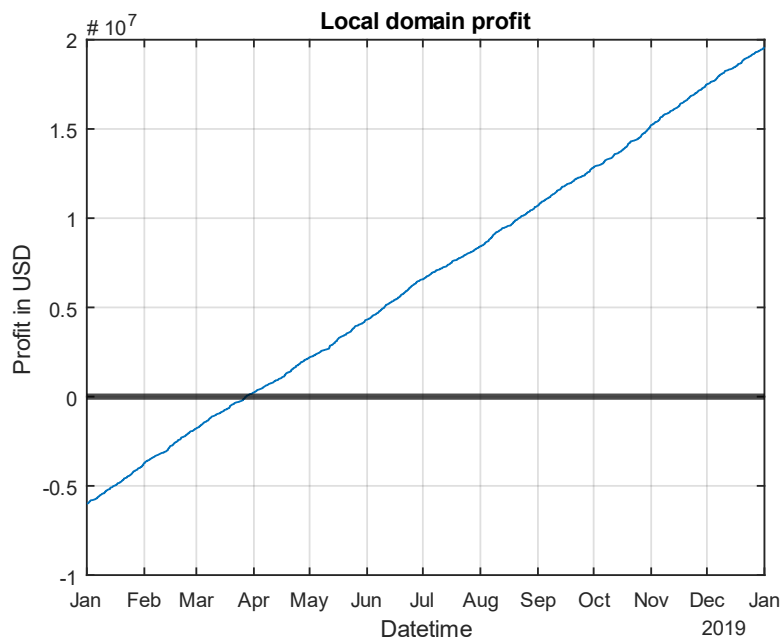


FIGURE 39: OPTIMISTIC PROFIT

2.3.5.5 Ideal scenario

The ideal scenario supposes that the infrastructure provided by the local domain is utilized at 100% and no infrastructure from the overflow domains is used.

Hereafter, the generated events, datacentre states, interchanged bandwidth among domains, blocked-services, not-scaled services and the price analysis for this scenario can be found.

2.3.5.5.1 Generated events

For the ideal scenario the number of generated events is 13547 events during a year timeframe. In Figure 40, the Service ID Histogram indicates how many events of each of the different 5G-TRANSFORMER service-like are generated in comparison with the other events can be found:

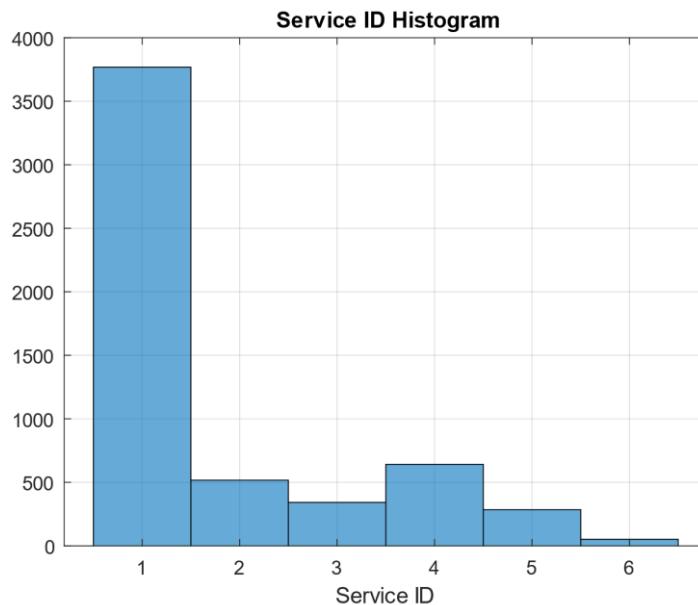


FIGURE 40: IDEAL SERVICE ID HISTOGRAM

In Figure 41, the Life-Time Histogram providing a general view of how much time last the services once they are deployed can be found:

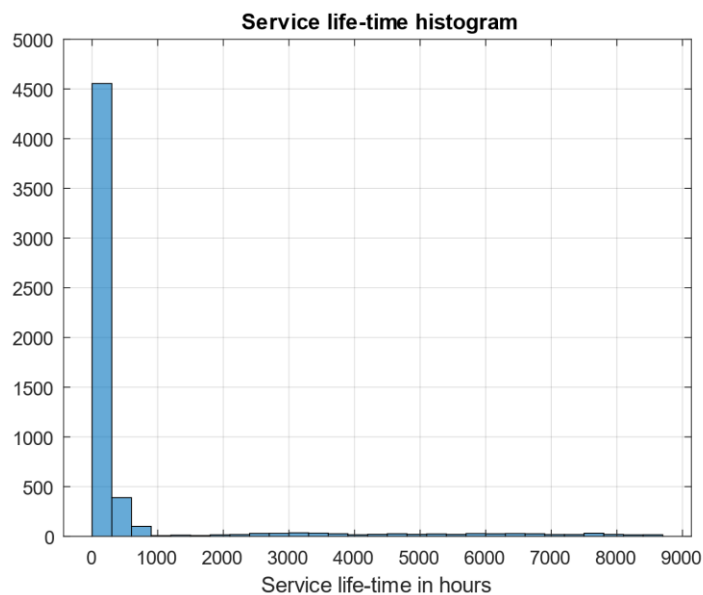
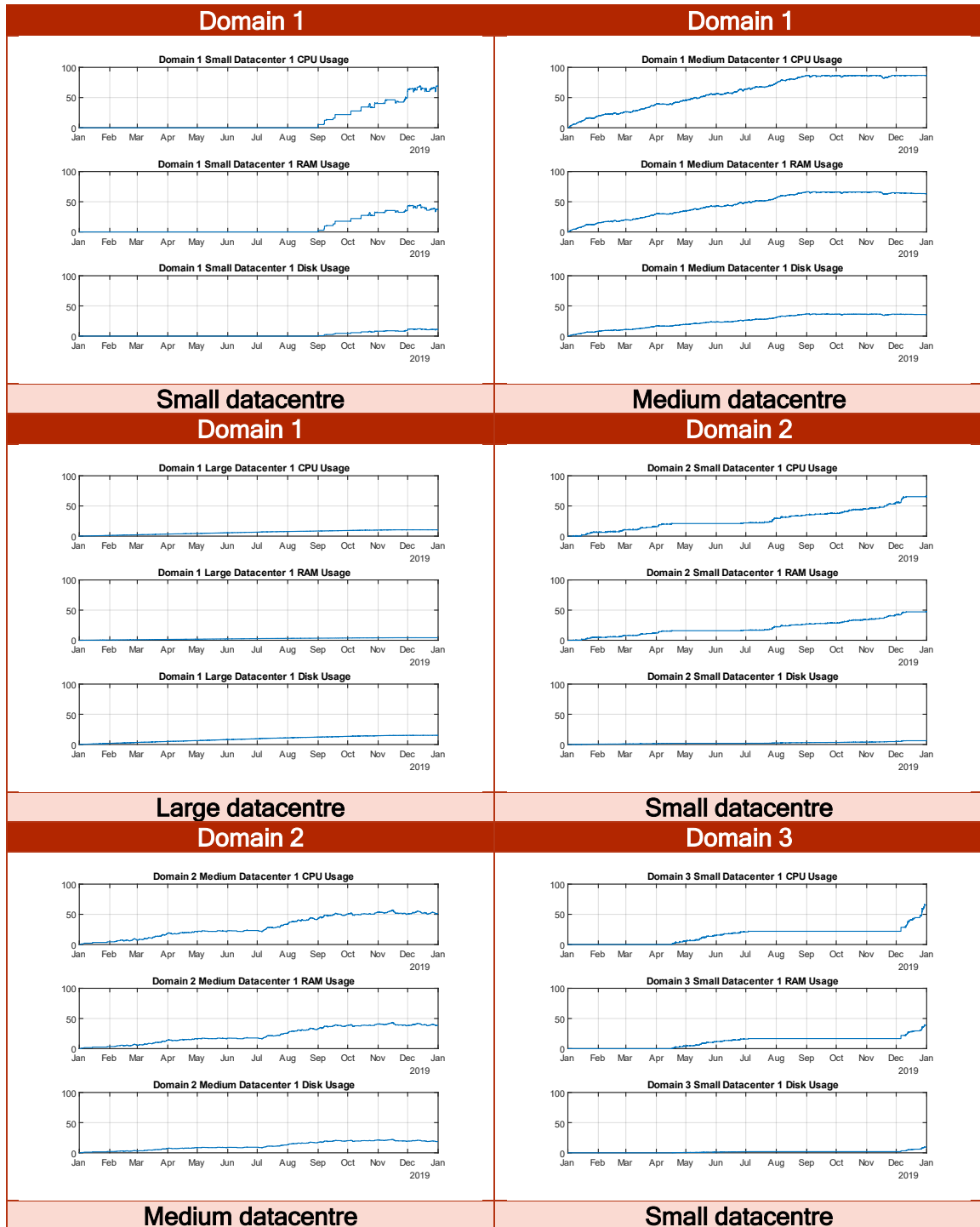


FIGURE 41: IDEAL SERVICE LIFE-TIME HISTOGRAM

2.3.5.5.2 Datacentres state

The state of the datacentres for the ideal scenario for the local domain can be found in Table 79:

TABLE 79: IDEAL DATACENTRES STATE



There's something that may be not well understood in the previous table: Why is the local infrastructure not operating at 100% when the demand is defined for such a purpose? The answer is composed of two parts:

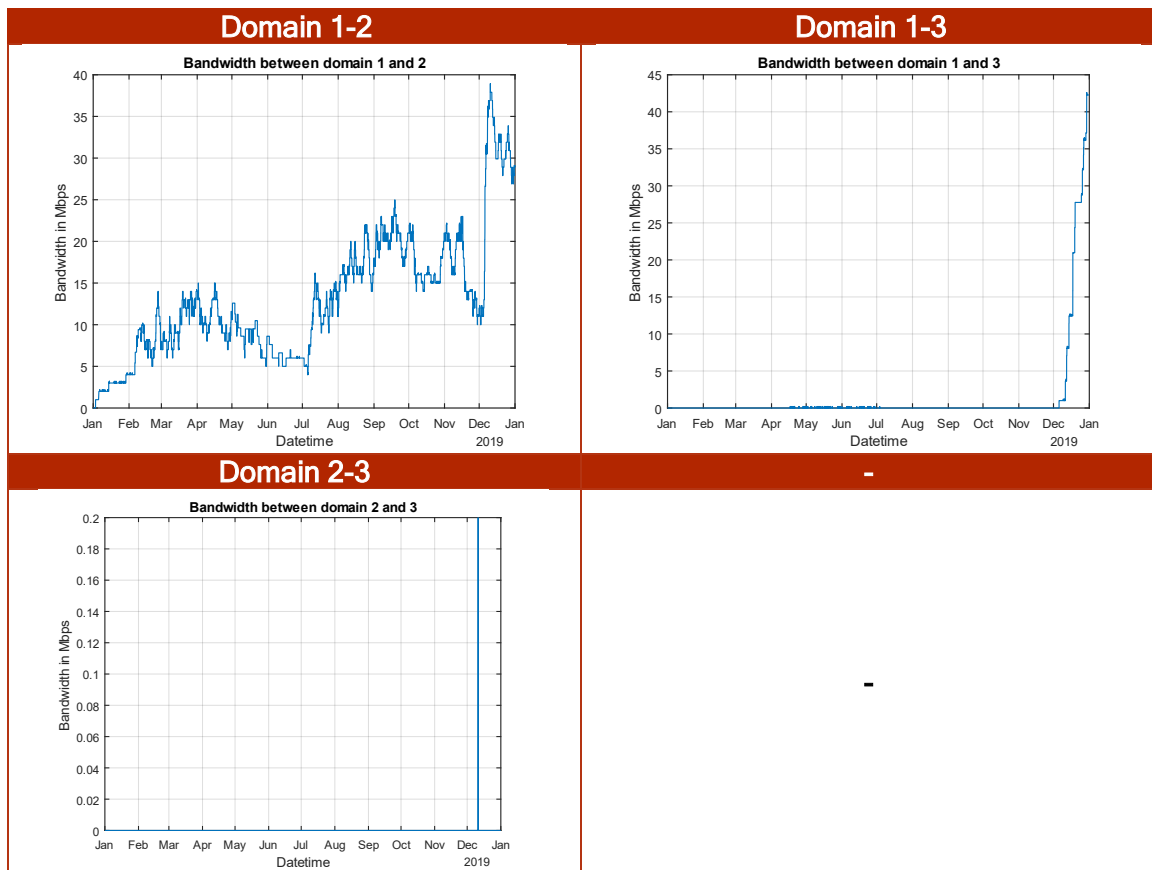
First, not all the services can be run in a large or medium datacentre because of the latency requirements. As a result, the large datacentre is being underused, but an increment in the federated traffic can be shown. This increment is also due to the special requirements: datacentre affinity, domain affinity and protection.

Second, datacentres are never operating at 100%, but they reach their limit as it can be seen in the medium datacentre of domain 1. Exactly at that time-point, the small datacentre starts to take the demand that cannot be process by the medium datacentre. This situation happens also in the overflow domains although it cannot be seen so clearly as both the small datacentre of domain 2 and 3 start taking the not-processed demand at the same time.

2.3.5.5.3 Interchanged bandwidth among domains

In the ideal scenario the interchanged bandwidth among domains is that one coming from the special requirements and all the federated services as shown here:

TABLE 80: IDEAL INTERCHANGE BANDWIDTH AMONG DOMAINS



2.3.5.5.4 Blocked services

In the ideal scenario there are no blocked services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.5.5 Not-Scaled services

In the ideal scenario there are no not-scaled services, as there is still vacancy in the local domain and no service will be federated to any overflow domain.

2.3.5.5.6 Price analysis

Price is defined by Equation 162 which is dependent on the profit margin and the breakeven price which is defined by Equation 161 and dependant the cost of this infrastructure, the forecast utilization of the local infrastructure, and the real utilization of the local infrastructure.

According to what it was stated in 2.2.3.5.1, in this scenario the price is incremented by the profit margin, raising the profit in a linear manner. For the experimental analysis, the following values are selected for the MATLAB simulations:

TABLE 81: IDEAL SCENARIO VALUES

	C_l	u_{r_l}	u_{f_l}	m
Small	260166.83	1	1	0.3
Medium	542361.89			
Large	5203616.68			

The maximum local demand for the pessimistic scenario will be:

TABLE 82: DEMAND IDEAL SCENARIO

	S1	S2	S3	S4	S5	S6
Small $d_{l,max}$	144	48	38	52	33	15
Medium $d_{l,max}$	864	288	230	314	203	90
Large $d_{l,max}$	21528	7176	5740	7828	5065	2266

With the previous parameters, the prices for each of the services of the 5G-TRANSFORMER ecosystem will be:

TABLE 83: IDEAL SCENARIO PRICES

	S1	S2	S3	S4	S5	S6
Small	\$2348.73	\$7046.19	\$8900.45	\$6504.18	\$10249	\$22547.8
Medium	\$816.06	\$2448.17	\$3065.53	\$2245.45	\$3473.26	\$7834.12
Large	\$314.23	\$942.69	\$1178.52	\$864.17	\$1335.58	\$2985.31

2.3.5.5.6.1 Cost

In the ideal scenario, the only cost is the infrastructure cost as indicated in 2.3.5.1.3. This cost is composed by one datacentre of each type.

Therefore, at the end of the year the cost will be:

TABLE 84: IDEAL COST

	Small	Medium	Large	Total
Cost	\$260.166,83	\$542.361,89	\$5.203.616,68	\$6.006.145,4

2.3.5.5.6.2 Revenue

The revenue for the scenario comes from the price that the tenants pay for a service as those indicated in 2.3.5.2.6. In addition, there's also revenue coming from the licenses, transaction for scaling up services and special requirements as in Table 57.

The revenue for the ideal scenario is shown in Figure 42:

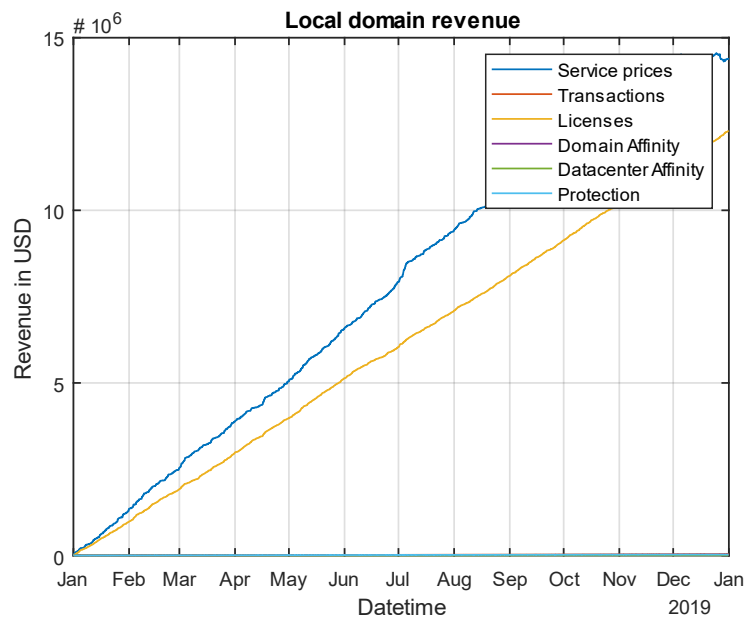


FIGURE 42: IDEAL REVENUE

2.3.5.5.6.3 Profit

The revenue for the ideal scenario is shown in Figure 43:

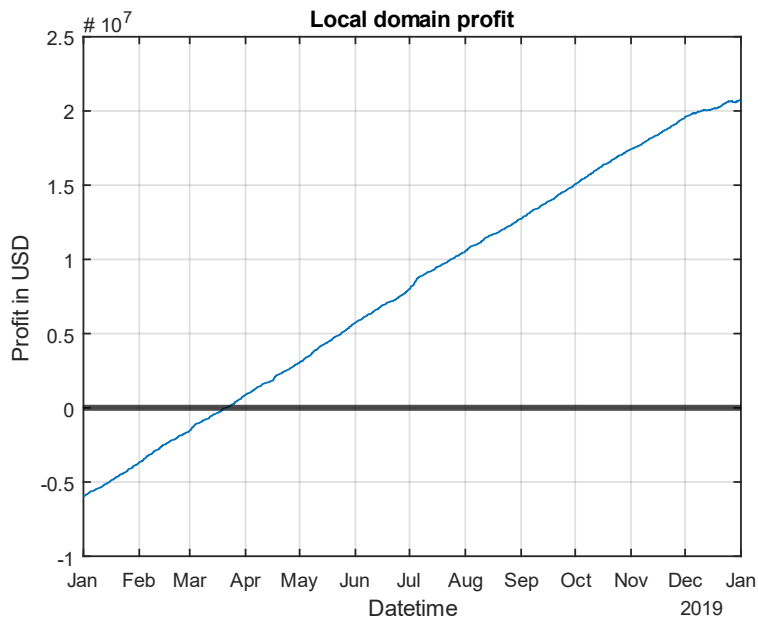


FIGURE 43: IDEAL PROFIT

2.3.5.6 Overflow scenario

This scenario supposes the local infrastructure to be operating as its maximum capacity and some of the services are deployed by the overflow infrastructure managed by the federation. The overflow infrastructure attends an additional 300% of the local demand.

Hereafter, the generated events, datacentre states, interchanged bandwidth among domains, blocked-services, not-scaled services and the price analysis for this scenario can be found.

2.3.5.6.1 Generated events

For the federation scenario the number of generated events is 40641 events during a year timeframe. In Figure 44, the Service ID Histogram indicates how many events of each of the different 5G-T service-like are generated in comparison with the other events can be found:

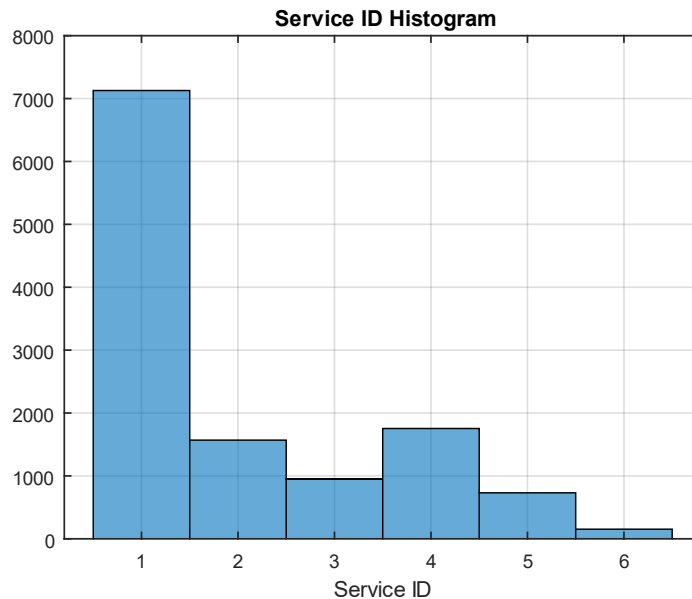


FIGURE 44: OVERFLOW SERVICE ID HISTOGRAM

In Figure 45 the Life-Time Histogram providing a general view of how much time last the services once they are deployed can be found:

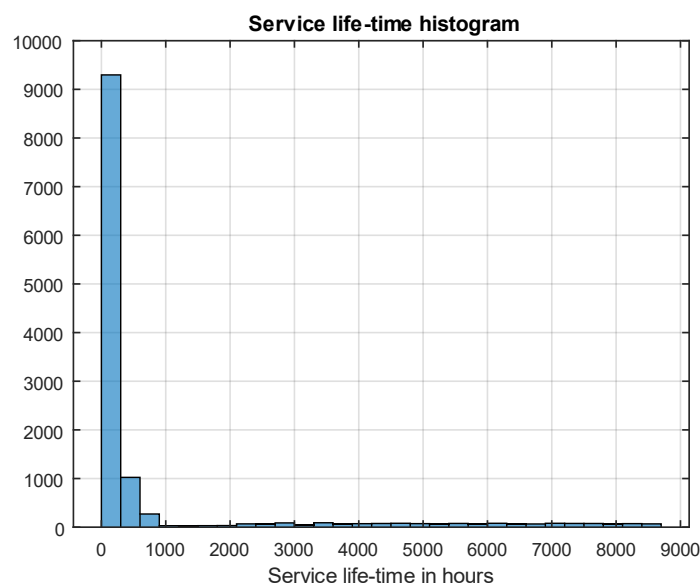
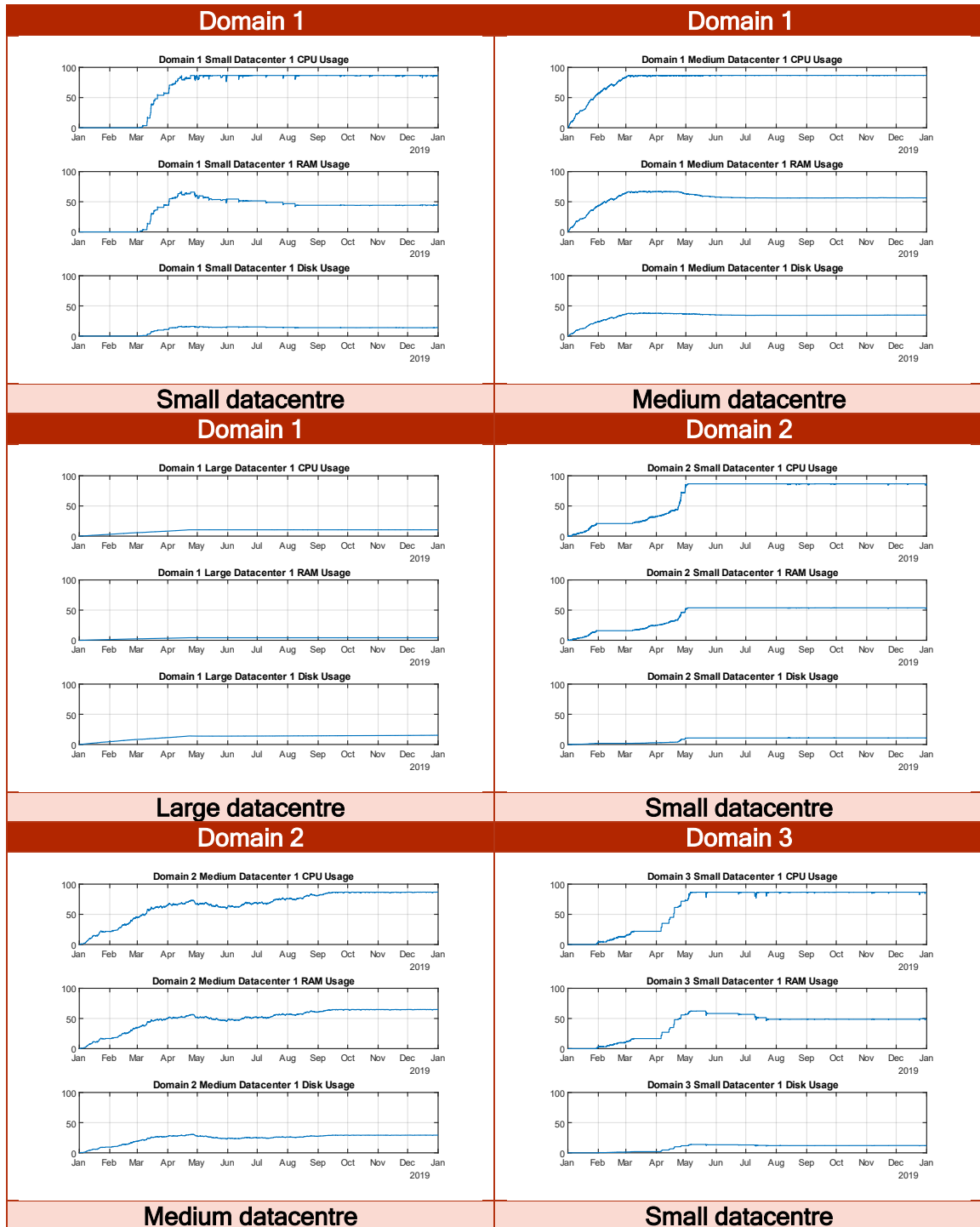


FIGURE 45: OVERFLOW SERVICE LIFE-TIME HISTOGRAM

2.3.5.6.2 Datacentres state

The state of the datacentres for the federation scenario for the local domain can be found in Table 85:

TABLE 85: OVERFLOW DATACENTRES STATE



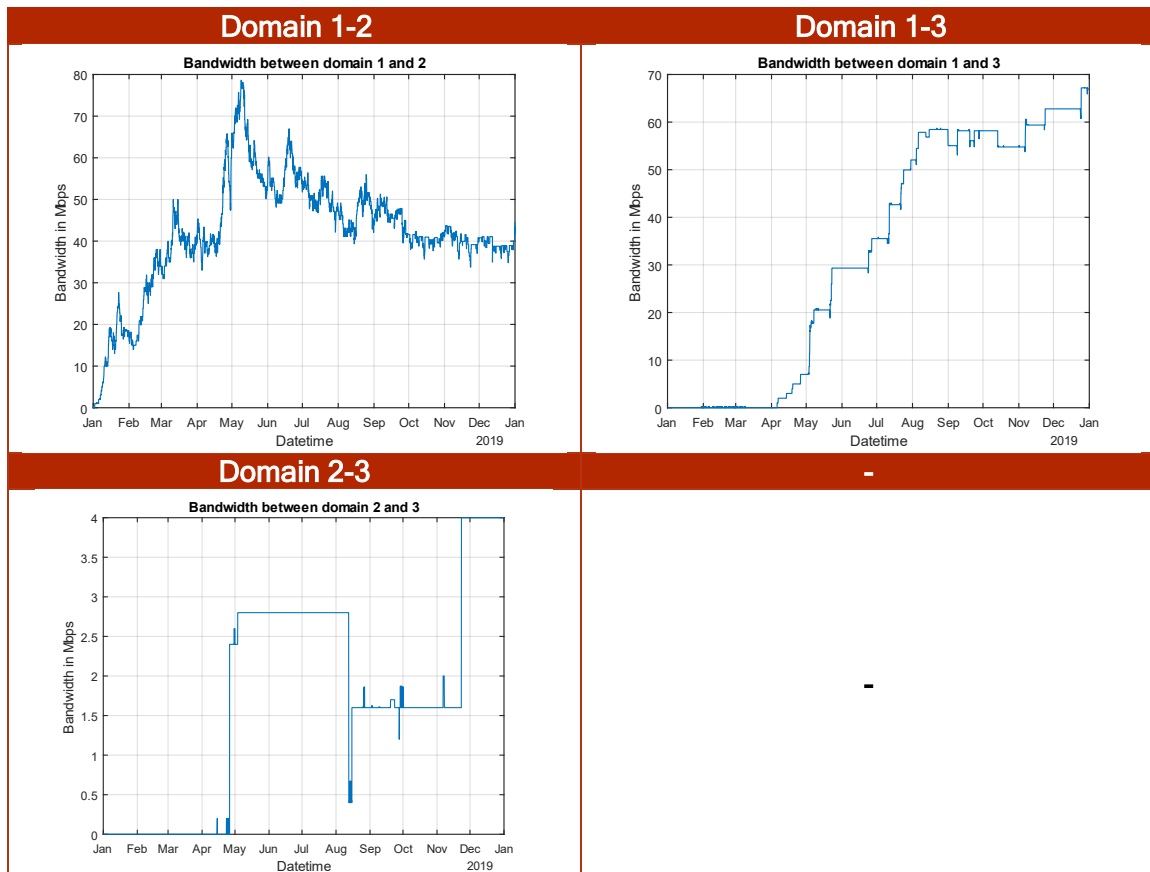
As it can be seen, all the datacentres regardless the domain are operating at its maximum capacity due to the huge demand. Therefore, there's no remaining capacity

and the infrastructure cannot be used until a previous service finishes. This will lead for the first time to blocked and not-scaled services.

2.3.5.6.3 Interchanged bandwidth among domains

In the overflow scenario the interchanged bandwidth is really elevated as it can be seen hereafter:

TABLE 86: OVERFLOW INTERCHANGE BANDWIDTH AMONG DOMAINS



2.3.5.6.4 Blocked services

As indicated before, for the first time all the datacentres regardless the domain are operating at its maximum capacity. This means that there will be blocked services in the federation as the demand excess both the local and the overflow demand.

From all the different causes that can lead to block a service in 2.3.4.2 only the occupation is analysed here. Blocked services must be understood as a under-dimensioning of the datacentre's infrastructure.

Regarding the demand it was stated in 2.3.5.6.1 that 13547 were generated. Not all of them can be deployed and these services that are not deployed must be understood as a lack of revenue, but they won't influence the cost in any way, as the infrastructure remains the same.

According to the previous comments, in Figure 46 the blocked services for each of the services considered in 5G-T can be found:

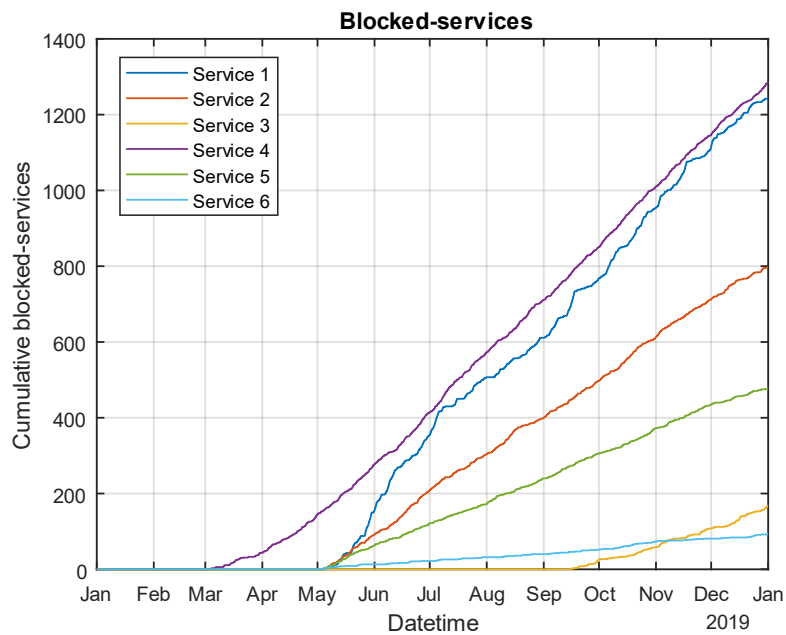


FIGURE 46: OVERFLOW BLOCKED SERVICES

2.3.5.6.5 Not-Scaled services

Not-scaled services in the overflow scenario are due to the same reasons as the blocked services.

In Figure 47, the not-deployed services for each of the services considered in 5G-TRANSFORMER can be found

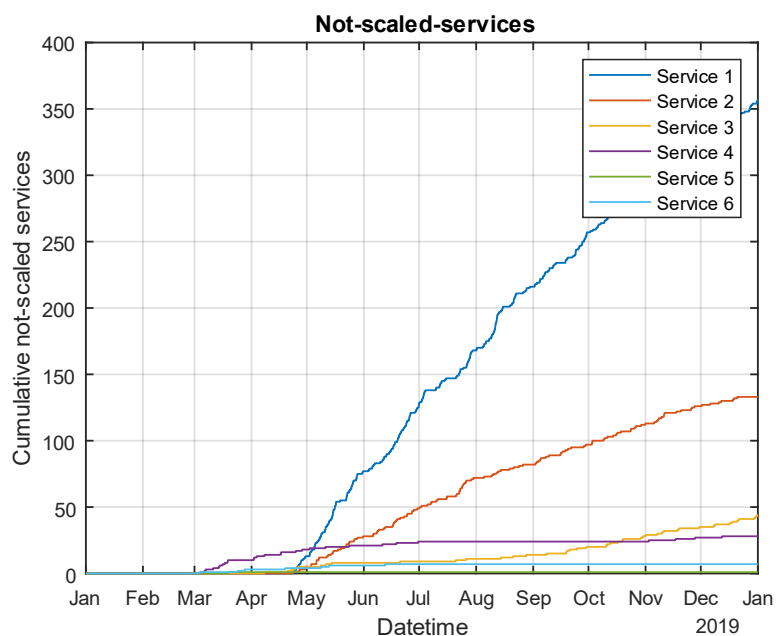


FIGURE 47: OVERFLOW NOT-DEPLOYED SERVICES

2.3.5.6.6 Price analysis

Price is defined by Equation 162 which is dependent on the profit margin and the breakeven price which is defined by Equation 161 and dependant the cost of this

infrastructure, the forecast utilization of the local infrastructure, and the real utilization of the local infrastructure.

According to what it was stated in 2.2.3.5.1, in this scenario the price is incremented by the profit margin, raising the profit in a linear manner. For the experimental analysis, the following values are selected for the MATLAB simulations:

TABLE 87: FEDERATION SCENARIO VALUES

	C_l	u_{rl}	u_{fl}	m
Small	260166.83	1	1	0.1
Medium	542361.89			
Large	5203616.68			

The maximum local demand for the pessimistic scenario will be:

TABLE 88: DEMAND FEDERATION SCENARIO

	S1	S2	S3	S4	S5	S6
Small $d_{l_{max}}$	144	48	38	52	33	15
Medium $d_{l_{max}}$	864	288	230	314	203	90
Large $d_{l_{max}}$	21528	7176	5740	7828	5065	2266

With the previous parameters, the prices for each of the services of the 5G-TRANSFORMER ecosystem will be:

TABLE 89: OVERFLOW SCENARIO PRICES

	S1	S2	S3	S4	S5	S6
Small	\$2710.08	\$8130.22	\$10269.75	\$7504.82	\$11825.77	\$26016.69
Medium	\$941.61	\$2824.81	\$3537.15	\$2590.91	\$4007.61	\$9039.37
Large	\$362.58	\$1087.72	\$1359.84	\$997.12	\$1541.06	\$3444.59

2.3.5.1.1.1 Cost

In the overflow scenario, the only cost is the infrastructure cost, as even without having vacancy in the local domain and federating services to the overflow infrastructure, the cost for the overflow infrastructure is not paid by the local domain as indicated in 2.3.5.1.3.

Therefore, at the end of the year the cost will be:

TABLE 90: OVERFLOW COST

	Small	Medium	Large	Total
Cost	\$260.166,83	\$542.361,89	\$5.203.616,68	\$6.006.145,4

2.3.5.1.1.2 Revenue

The revenue for the scenario comes from the price that the tenants pay for a service as those indicated in 2.3.5.2.6. In addition, there's also revenue coming from the licenses, transaction for scaling up services and special requirements as in Table 57.

The revenue for the overflow scenario is shown in Figure 48:

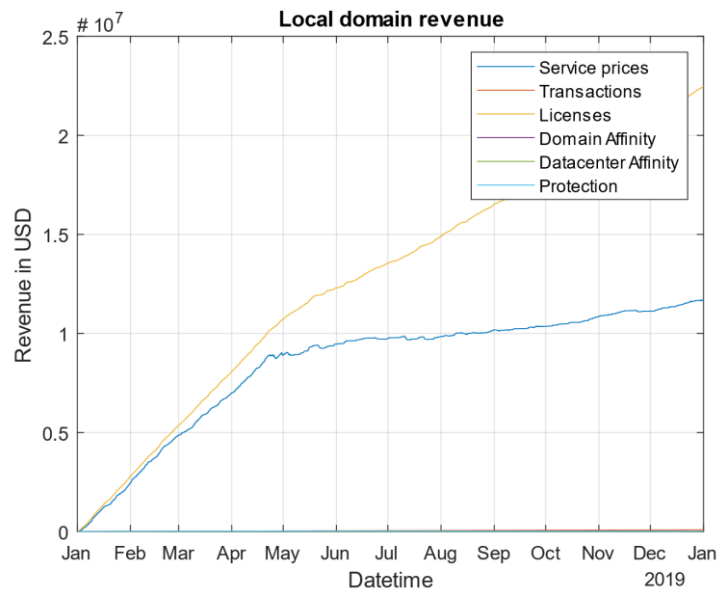


FIGURE 48: OVERFLOW REVENUE

As it can be seen, when the local infrastructure reaches the maximum point of capacity, the price for the services decrease. In addition, the price for the licenses decreases but not as much as the price for the VNFs. This is because not all the VNFs are federated (in this case, the licenses are paid to the overflow domain and they are reduced from the VNFs price), but some of them are still deployed in the local domain.

2.3.5.1.1.3 Profit

The revenue for the overflow scenario is shown in Figure 49:

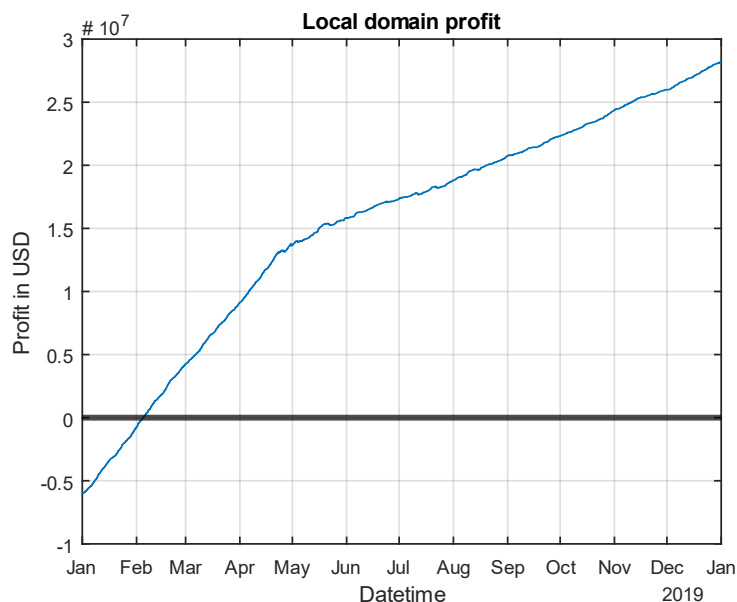


FIGURE 49: OVERFLOW PROFIT

In the previous figure it can be seen that the simulation reports the same profit structure that was predicted in Figure 20. Because of that, we can say that the experimental study validates the analytical study.

2.4 Summary

This section aims at providing some conclusions to the techno-economic analysis in the form of some lessons learnt, the KPIs of the project regarding CAPEX and OPEX expenditures and the business model transformation impacts on the vertical use cases implemented in the 5G-T platform.

2.4.1 Lessons learnt

There are several lessons that can be learned from the techno-economic analysis of the 5G-TRANSFORMER project detailed in this deliverable:

- *Estimating the demand is not an easy task.* As all the 5G-TRANSFORMER are novel services, it is not possible to know a priori how their demand will be. As a result, a sensitivity analysis was proposed, resulting in 5 scenarios: Pessimistic scenario, Realistic scenario, Optimistic scenario, Ideal scenario and Federation scenario which considers a huge demand.
- *Price also feedbacks the demand.* In the techno-economic analysis a price independent from the demand has been proposed. However, as the price increases the demand tends to decrease and vice versa.
- *Weight of the forecast utilization of the local infrastructure in the profit.* As it was indicated in analytical study, making an accurate estimation of the forecast utilization of the local infrastructure will determine which profit margin must be selected to overcome the breakeven and ensure profit.
- *Role of the profit margin in the profit.* Regarding the previous point, there might be some situations in which even when setting a profit margin, the breakeven point is not reached. For example:
 - When the forecast utilization of the local infrastructure is lower than the real utilization of the local infrastructure, then the breakeven can only be reached by increasing significantly the profit margin, which in a real situation is normally not possible.
 - When Equation 163 is not fulfilled and federation takes place, the price that is set results being too low to overcome the high rate of occupation of the overflow infrastructure and its associated increment of price can prevent the local domain from reaching the breakeven point.
- *Weight of the licenses per user in the profit.* As it was shown during the experimental analysis, licenses play an important role in the profit. They provided a revenue that in some cases is similar to that one of the services themselves. Their price however is difficult to be predicted, as the number of users and the associated demand is unknown.
- *The real utilization of the datacentres is more complex than what is shown analytically.* Along all the experimental study it was shown that the utilization of the datacentres is not the same as that one indicated in the analytical study as the placement algorithm is of high complexity. Datacentres are not only selected according to their capacity and price but also by their latency, occupation and special requirements such as datacentre affinity, domain affinity or protection.
- *Need for negotiation of convenient rates for the federation in order to make it profitable.* In order to ensure profitability when federating a service, the price that is offered by an overflow domain to the local domain has to be smaller than

the market price. This is possible when considering a wholesale offer rather than a market offer.

- *Federation is not easy to be given as a really big demand is required, it can imply a change of paradigm:* Federation needs of a really huge demand, as otherwise services are run in the local infrastructure when it is possible. Therefore, in order to ensure a real scenario in which federation takes place as the normal way of operating, a change of paradigm should be considered by reducing the local infrastructure size and relying in the overflow infrastructure.
- *Latency plays the most important role in a federation:* How the services are distributed among the different datacentre types, either in the local or the overflow infrastructure depends on the latency that its VNFs require. Therefore, it can be said that latency plays the most important role in a federation and the administrative domain that reduces it the most will take the majority of the federated demand.
- *The ideal location of the datacentres of each administrative domain can determine the success or failure of an infrastructure provider:* In the experimental study it was considered that the large datacentres are further from the end user than the small datacentres which is the most common approach. It can be interesting study the implications of placing medium datacentres closer to the user, especially in areas where the demand could be considered highly elevated.
- *Not all services are equally profitable when federating:* Services with a high life-time are less profitable than those that have a lower life-time. Therefore, developing some kind of policy that depending on the datacentre infrastructure accepts or rejects deploying a service could be beneficial for the profitability of the 5G-TRANSFORMER ecosystem.

2.4.2 KPIs of the project

Solutions based in 5G-T approach leverage on the multi-tenancy approach allowing the sharing of a given infrastructure simultaneously by multiple verticals. Apart from that, leveraging on a 5G-T provider for the provision and operation of the vertical service can alleviate the vertical customer from the operational tasks, relying directly on the provider. This section analyses how the 5G-T platform can provide beneficial results in terms of CAPEX and OPEX in comparison with a scenario where each vertical customer deploys and operates its own infrastructure for the same service. The Key Performance Indicators KPIs of the 5G-T ecosystem are presented in Table 91:

TABLE 91: KPIs OF THE PROJECT

	Name	Value
KPI 1	Reduction of today's network management OPEX	>20%
KPI 2	Reduction of today's network resource utilization CAPEX	>20%
KPI 3	CAPEX and OPEX savings due to slicing and multi-tenancy (depending on the number of verticals)	>80%

To achieve these objective KPIs, the following assumptions are taken:

- *Infrastructure scalability.* it is commonly accepted that the sharing of infrastructure brings benefits in terms of costs because of the maximization in the usage of resources, assuming a given service demand. In the case of 5G-T, the compared scenarios are, on one hand, a number of vertical customers deploying their own infrastructure, and on the other hand a 5G-T provider deploying an infrastructure dimensioned to host the same number of vertical services. In both cases, the considered infrastructures have to account for both infrastructure acquisition costs of servers, switches and circuits (linked to CAPEX) and infrastructure running costs of hosting and maintenance (linked to OPEX).
- *Service operational tasks.* in case the vertical customer assumes the responsibility of directly deploying the infrastructure and the service on top of it, then it is necessary to have the proper staff (linked to OPEX) and tools (linked to CAPEX) for such purpose. It can be assumed that most verticals will not fully dedicate technical staff only in these tasks, so the costs related to this can be calculated based on a given % of the working time of the employee. Regarding the tools for deploying and operating the service, the vertical customer has to purchase and own the necessary systems facilitating the operation. In contrast to this situation, when moving to the 5G-T approach, the provider is the one requiring to fully allocating expert staff as well as the necessary artefacts for the deployment and operation of the services.

The following sub-sections detail the calculations considered in this analysis.

2.4.2.1 Infrastructure scalability

To study the impacts on costs of the infrastructure scalability, we can directly take the previous calculations performed for characterizing both a small and a large DC. Considering lower or higher needs does not actually vary the conclusions since they are proportional to the final number of vertical customers that could share a single service provider infrastructure. This is why, in our approach, we directly take the small and large DC for their level of details.

As a further simplification, we take as key parameter for dimensioning the number of vCPUs in each setup. Thus, from the dimensioning exercise performed for characterization of the DC infrastructure, a small DC comprises 576 vCPUs, rising a total costs (in a three years window) of \$787.700.

Looking at the large DC dimensioning, such DC provides a capacity of 86.112 vCPUs. This means that in terms of vCPU capacity, the large DC provides a capacity higher than 149 times the small DC.

2.4.2.2 Service operational tasks

As described before, in case each vertical owns the infrastructure it is required to have staff in charge of provisioning and operating the service, even though such staff will not be dedicated full time to the task. On the contrary we can assume that in the side of the service provider, dedicated staff takes care of the provision and operation of the services from the vertical customers full time.

In addition to that, from a cost perspective, it can be also assumed that the skills of the staff devoted to these tasks is higher than the skills required for the maintenance staff considered in the infrastructure TCO analysis. In order to characterize that, we consider

the same amount of persons in each case but with a 30% higher cost because of higher skills are required. Regarding the dedication, we assume that service staff is dedicated 30% of their time in the case in which the vertical operates its own service, while the dedication is 100% time for the service provider.

Finally, an estimation on the tools investment is considered. For the service provision and operation each vertical customer will require having their own set of tools. In the service provider case, only a set of tools will be necessary, however the software tools required have to be properly dimensioned to support the expected number of customers, so a higher cost (e.g., in terms of licences) can be expected, but lower than lineal.

No commercial reference is available for 5G-T components. Then, it is required to elaborate some hypothesis. The most demanding work can be considered to be in the orchestration of services in the compute nodes available. Attending to the number of compute nodes per DC, the small one has 3 while the large one has 300. Software product development is usually scaled in a stepwise manner, with some fixed costs independent of the scale of infrastructure to be managed. For the sake of simplicity, it will be assumed that in the case of the vertical customer owning and operating the infrastructure the minimal number of supported compute nodes is established in units of 10, while for the service provider case it is established in units of 100. However, the cost associated to the products is not usually proportional, thus is, the cost for a unit of 100 it is not 10 times the cost for a unit of 10, but lower. We consider here a value of 7, based on the usual commercial strategies, the power of negotiation of a service provider, etc. For the point of view of pricing in this exercise we assume a cost of \$50.000 for the management and control of up to 10 compute nodes.

2.4.2.3 Results

According to the considerations above, the following calculations are performed:

- Infrastructure related
 - Cost of the DC, including CAPEX and OPEX, taken directly from the TCO analysis for dimensioning the DCs
 - Cost of the circuits necessary for delivering the expected traffic outside the DC. In order to make things comparable we consider that circuits are rented in both scenarios.
- Service provision related
 - Cost of high skilled staff for provisioning and running the services.
 - Cost of the tools required for provisioning and running the services.

The resulting calculations are as followed. In line with the bandwidth calculations, the maximum number of verticals that finally could be handled in a large DC is 103. So, this is the base for comparison.

For the small DC, the TCO for 3 years results in the following cost structure:

TABLE 92: SMALL DATACENTRE ANALYSIS FOR KPI OBJECTIVES

	Infrastructure cost	Staff	Tools	
CAPEX	\$162689	\$0	\$50000	
OPEX	\$617811	\$723320	\$0	TOTAL
TOTAL	\$780500	\$723320	\$50000	\$1553821

In the case of the large DC, the TCO accounts for:

TABLE 93: LARGE DATACENTRE ANALYSIS FOR KPI OBJECTIVES

	Infrastructure cost	Staff	Tools	
CAPEX	\$7322337	\$0	\$1050000	
OPEX	\$8288513	\$8679842	\$0	TOTAL
TOTAL	\$15610850	\$8679842	\$1050000	\$25340693

Then, comparing the cost of the service provider scenario (i.e., 25,3 M\$) versus 149 times the cost for a single vertical (i.e., 231,5 M\$), the 5GT service provider scenario is 89,05% more cost efficient than the case when each vertical owns and operates individual infrastructures.

2.4.3 Business model transformation of 5G-T vertical use cases

The 5G-T system promotes new mechanisms for the monetization of the new network generation and new business models for the vertical industries. Hence, 5GT offers a unique value that can be exploited for revenue generation by operators and application service providers alike, thus creating new value chains and a variety of interaction models with operators and IT actors

As part of the Techno-economic study, and in order to show the added value of the 5G-T platform on the deployment of the 5G-T vertical use cases, we explain hereafter the business model transformation of the 5G-T vertical use cases by answering the two following questions:

- What kind of transformation has been brought by the platform to the classical model of the vertical service?
- What are the economic benefits to the Vertical industry business with the use of the 5G-T platform features?

2.4.3.1 MVNO

The 5G Phase 2 (beyond the 5G NR) is focused on the B2B market (but not the B2C), which presents a significant digitization challenge in Europe. In this frame, MVNO offerings that provide optimized support for B2B services via network slicing are considered as relevant for many vertical domains, mainly, the Industry 4.0, the Public Safety and the Automotive.

In a general perspective, the Network slicing concept that is one of the main features implemented and tested in the 5G-T project has impacts on the business model of the MVNO and the MNO. In fact, it allows the provisioning of network slice as a service (NSaaS) by a mobile network operator (MNO) hosting a virtual mobile network operator (MVNO) and, as a consequence, provides the flexibility to:

- Create a network slice instance with several types of access technologies (currently EPC, WiFi, LoRa, soon NB-IoT, 5G core), features and CUPS resulting from a composition of available VNF.
- Scale the service capacity using deployment flavors
- Deploy network service across multiple sites abstracting the underlying virtualization technologies

With regards to the infrastructure cost analysis performed in section 2.2, and compared to classic infrastructures where the infrastructure is sold and calculated as a whole, the cost for an MVNO/MNO has transformed to include more elements/factors such as:

- The cost of virtualized resources (vcpu, ram, storages) is breakdown into several costs of hardware, space, staff and power that can vary according the size of the datacenter (small, medium, large)
- The cost of the cloud according to its type (private cloud or public cloud), size and its reservation modes.

This transformation of the infrastructure cost calculation allows:

- An infrastructure cost estimation of a network slice while varying its capacity
- A better selection of the NFVI provider according to the use of the resources (federation of VIMs)
- The ability to adjust resource allocation based on service demand for cost saving

2.4.3.2 Entertainment

5G Transformer has opened the possibility of a new paradigm in the massive content distribution in a sport venue, by deploying the content distribution infrastructure in the network (MEC) rather than in the traditional CDNs. In this way, the deployed service allows the attendees to a sport event to access real-time enhanced media content, enriching the sport fan experience.

By using the 5G-T platform to deploy the Entertainment vertical use case, the following economic benefits are noticed from a business model perspective:

- No need to deploy per event expensive ad-hoc high density WiFi infrastructure with a large bandwidth to Internet for CDN support, which allows infrastructure deployments cost saving for the Vertical.
- Cabling cost reduction, especially in the open venues cases, where video and other high bandwidth content needs to be exchanged in open fields

2.4.3.3 Automotive

For the Automotive vertical domain, the main business model transformation brought by the 5G-T platform includes:

- The involvement of new actors and business roles that interacts with the Automotive vertical such as, the Auxiliary Service Providers (3rd Party entity, i.e. CIM), the MEC Providers and the 5G-T provider.
- The automatic deployment of a service where it is needed
- More guaranteed SLAs in the deployment of the service (reliability, latency, density) as they are taken as necessary configuration requirements in the creation of the network slice
- Arbitration, Scalability

The evolutions lead to:

- New business model for connected vehicles involving new actors
- Data monetization (sending data to Neutral Server, e.g CIM)
- More opportunities for the Automotive Industry range from savings in public safety through a reduction in the number and the severity of accidents (both

costs in lives and property damage) to savings in infrastructure planning, deployment and maintenance.

In a longer-term vision of the Automotive industry assuming a complete penetration of V2X applications, more benefits from the 5G-T like systems can be foreseen, including:

- Annual economic damage from accidents reduced by up to 6.5 billion EUR in Europe
- Up to 4.9 billion EUR of economic losses might be avoided due to improved traffic efficiency and reduction of environmental damage (simTD Project, June 2013).

2.4.3.4 eHealth

The 5GT platform provides rapid deployment of emergency services that can potentially save people's lives and help triggering proper reaction to catastrophes. Besides, it presents an integration baseline for multiple emergency services (fire department, police, etc.) and enables the extension and sophistication of the core services. Using the 5GT platform the emergency service can rapidly adapt to new events and utilize the already deployed infrastructure to provide better coordination and support to the emergency teams on site.

Using the 5GT-platform in the eHealth use case contributes to the evolved and rapid service for the emergency events that can potentially reduce the casualties due to slow responses and the time that it takes to transfer patients to the hospital

2.4.3.5 eIndustry

The 5GT-platform brought transformation to the eIndustry with respect to the classical model of their vertical service especially in terms of operational aspects. Actually, the 5GT-platform enables the automation a lot of operations such as request of services, reconfiguration of service. Moreover, the 5GT-platform expose to the vertical interfaces for automatic operations based on SLA parameters that allows hiding the technical details proper of the infrastructure. Such characteristics allow the vertical to improve their business in the following main aspects:

- Reduce time of service configuration and improve the service management;
- Focus on core business activity, delegating to an external actor (usually a network operator) the operation of management. Hence, it avoids having skilled people dedicated to such activity.
- Have a SLA based interfaces to view and monitor the service

The experiment that has been carried out for eIndustry use case, demonstrated the capability to interconnect different geographical sites also in case of latency critical use case. Hence the eIndustry can extend the centralized control in a remote location to more equipment in different sites and enlarge the advantage in terms of monitoring and prediction of possible issues. The project will not demonstrate this aspect but provide a concrete demonstration of capability to go in this direction.

3 5G-TRANSFORMER Final architecture design and refinements

3.1 Summary of final architecture design and refinements

The final design of the 5G-TRANSFORMER architecture, as presented in Figure 50, follows the same concepts and design defined in the refined version of the baseline architecture reported in D1.3 [2], and same function roles for the three main architecture components: Vertical Slicer (5GT-VS), Service Orchestrator (5GT-SO) and Mobile Transport and Computing Platform (5GT-MTP).

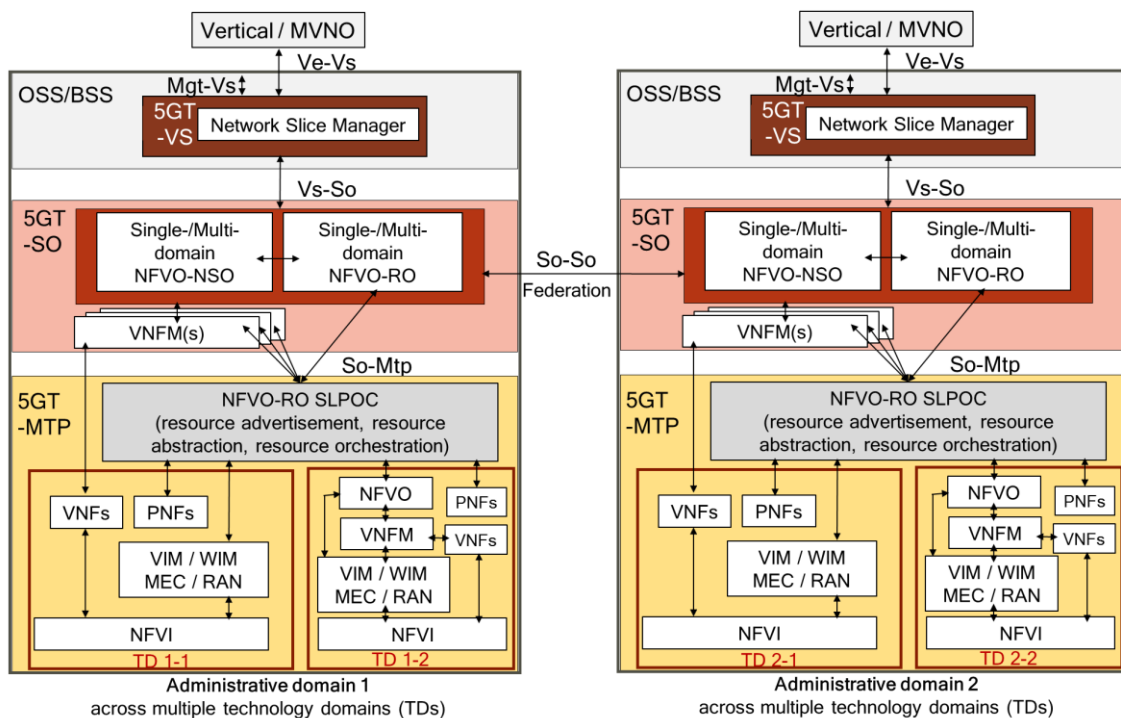


FIGURE 50: 5G-TRANSFORMER SYSTEM ARCHITECTURE

However, there are two added functionalities and features at the 5GT-VS, 5GT-SO and 5GT-MTP in the final architecture design. They are described in detail in the following subsections (Section 3.2 and Section 3.3) and are summarized hereafter:

- RAN (Radio Access Network) abstraction is proposed to enable to connect vertical applications to the 5G radio networks. Also, the 5G-TRANSFORMER architecture supports RAN slicing allowing the separation of the RAN network into several logical networks that support different vertical services. For this purpose, the 5GT-MTP provides an abstraction of the RAN to the 5GT-SO while the verticals describe the coverage area they are interested into. The 5GT-SO requests the RAN resources from the 5GT-MTP and connect them with other VNFs or NFV-NSs. This extension is further described in section 3.2.
- MEC (Multi-Access Edge Computing) solution is proposed to support the deployment of MEC applications and services at the edge of the network. To this aim, the architecture design of 5G-TRANSFORMER is extended to integrate MEC related functions, which includes end-to-end MEC application

and service on-boarding and instantiation workflows, providing traffic offloading and MEC service management. The main approach is to handle the MEC applications the same as the normal VNFs, but with some differences regarding on-boarding and instantiation workflows. How and where to perform the orchestration remains the same. That is, MEC applications do not have any `flavorId` or `instantiationLevel` reference defined in the AppD. Accordingly, the 5G-TRANSFORMER platform proposes extensions of the NSD to handle MEC specific parameters and data models. Further details of the MEC solution are provided in section 3.3.

Table 94 summarizes the main features included in the final 5G-T architecture (described in this deliverable) including the features defined and implemented in the initial design (as described in D1.2 [2]) and those implemented in the refined design (as reported in D1.3 [2]).

TABLE 94: SUMMARY OF FEATURES IN THE 5G-TRANSFORMER DESIGN

Features in initial design	Features added in refined design	Features added in final design
Vertical Service Descriptor to Network Service Descriptor translation	Service scaling (vertical-driven, 5GT-VS driven and 5GT-SO auto-scaling)	Enhanced RAN abstraction
Initial interface towards verticals	Enhancement to vertical support with enhanced NBI towards the verticals based on the REST-based API and web-based GUI for the vertical service definition (VSB, VSD) and service configuration, creation of new instances of vertical services at runtime directly triggered by vertical applications	Full MEC support
Basic arbitration	Enhanced arbitration	
Basic authentication and per-tenant authorization	Policy management	
Basic arbitration mechanisms	Vertical service composition	
Lifecycle management of simple Vertical Services and Network Slices (no service composition)	Advanced Lifecycle management of Network Slices, Network Slice Subnets with service compositions	
Initial Vs-So and So-Mtp interfaces (single PoP)	Updated Vs-So and So-Mtp interfaces (also support inter-PoP and federation)	
NV-NS creation, instantiation, termination and query operational status operations	NFV-NS on-boarding, NFV-NS scaling, NFV-NS service composition and federation	
Resource orchestration	Initial RAN abstraction	

functions	
Basic monitoring platform	Enhanced monitoring platform fully integrated to the system to enable: <ul style="list-style-type: none"> • monitor the vertical service (5GT-VS) • service assurance operations and automated SLA management (5GT-SO) • monitor the domain resources (5GT-MTP)
Cloud abstraction	Support for inter-NFVI-PoP orchestration
Resource allocation and termination for VIM and WIM domains	Enhanced placement algorithms (5GT-MTP), including functional splits
	Enhanced placement algorithms (5GT-SO) to handle the location and latency constraint
	Dynamic VNF configuration on instantiation
	Initial MEC support

3.2 Radio Abstraction

An integral part of any vertical service within the 5G-TRANSFORMER system is connecting the mobile devices via the 5G air interface to the vertical applications. The 5G RAN has been defined with support for network slicing, allowing the separation of the RAN into several logical networks, see [5] for a description of slicing in the RAN. Hereafter, we describe how the RAN is included in the various descriptors and how this information is processed in the 5G-TRANSFORMER system to automatically create slices in the RAN for different vertical services. We also describe how this has been implemented and used in one of the PoCs of 5G-TRANSFORMER.

We describe the general concepts in Section 3.2.1 and thereafter describe how RAN abstractions has been or could be introduced to each of the 5G-TRANSFORMER components 5GT-VS, 5GT-SO, and 5GT-MTP.

3.2.1 General Concepts

A vertical service includes the vertical applications as VNFs. When mobile devices or user equipment (UE) are connected to the vertical applications via a 5G mobile network, then this 5G network is part of the vertical service as well. For an exemplary vertical service instance, the diagram in Figure 51 shows one UE, two gNBs in the RAN, core network functions for U-plane (UPF), session handling (AMF, SMF) and more generic C-plane functions (AUSF, LCM, etc.), and the vertical application (VA). For sake of simplicity just one VA is shown here, also in practice several different VNFs constitute the VA.

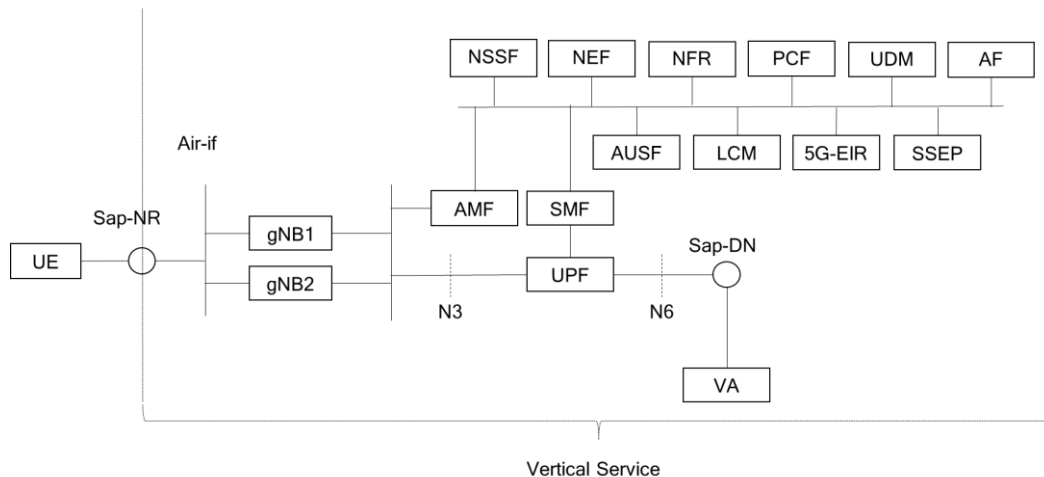


FIGURE 51: RAN AND CN ENTITIES INVOLVED IN THE DEPLOYMENT OF A VERTICAL SERVICE

The UE is considered as not belonging to the vertical service and not under control of the 5G-TRANSFORMER system. All the other network functions including the the gNBs are under control of the 5G-TRANSFORMER system. Note that this diagram shows a simplified view, e.g. multiplicities of network functions are not indicated. Only two gNBs are shown, but all network functions in the diagram may have more than one instance. Also, one may expect separate management connections to all these network functions; these have been omitted for the sake of simplicity as well.

Whereas core network functions and vertical applications can be handled as virtual network functions, this is different for the RAN. The RAN is deployed to a much larger extend with PNFs, which implies it has to be handled differently. E.g., in case of capacity shortage, an operator cannot simply instantiate another instance of a VNF, service personal would actually go to the field and physically deploy an additional gNB.

In the following, we describe how the RAN can be abstracted and how this abstracted description could be handled in each of the main 5G-TRANSFORMER system components. One design goal for this abstraction is to keep the 5GT-SO free of radio knowledge, allowing it to remain useful for vertical services with different radio technologies.

The support of RAN abstraction in vertical services requires extending the abstraction concept to consider the following two points:

- The radio and transport should be exposed as unique infrastructure to provide “connectivity” among the end points of the service, such as vertical user devices (e.g. robots, etc.) and its applications running remotely in some server.
- The vertical service usually is defined in a certain geographical area that doesn’t to coincide with the radio-coverage areas

Considering the previous two points RAN abstraction was defined in a similar way as for MEC where the mobility of a service is considered and defined by the concept of tracking area [5]. Actually, an abstract view of a RAN should be defined in a certain geographical area suitably defined representing the area where the service should be provided and the corresponding radio coverage. Therefore, it can happen that the radio

coverage areas do not match with the service areas, so an interaction among 5GT-SO and 5GT-MTP is necessary to deal with this situation.

In line with the abstraction method defined in [4] and [5] for the transport and data-centres, also for the RAN we adopted the concept to hide from the 5GT-SO all the technical parameters of the resources and instead exposing service parameters. Knowledge on the radio aspects is very specific to the used technologies and is better kept in the respective 5GT-MTPs. The 5GT-SO cannot - and should not - decide which specific gNBs to use for a specific vertical service instance. Instead, the 5GT-MTP will provide an initial abstraction of the RAN to the 5GT-SO. The 5GT-SO will then request a specific RAN instance for vertical service instance from the 5GT-MTP, which will provide such RAN instance including the corresponding SAPs. Eventually, the 5GT-SO can connect the CN and the VAs to these SAPs to establish the overall service. This is a slightly more complex procedure as for storage and network resources, where the 5GT-SO could make more decisions using appropriate placement functions.

In Figure 51, we did not indicate which of the network functions actually constitute the RAN. Actually, this can be defined differently. Two approaches are shown in Figure 52:

- The first, indicated in blue, contains the gNBs and the connections among them. This abstraction is more aligned with the usual separate into RAN and CN as defined by 3GPP.
- The second, indicated in red, includes as well the core network functions most relevant for connecting UEs to a data network as well. These functions are the user plane function (UPF), session management function for handling individual PDU sessions (SMF), and the access management function for handling the connectivity of UEs at all (AMF). This abstraction terminates at the Sap-DN, where VAs, especially at the mobile edge can be connected to the RAN.

The second abstraction has been used in use cases with very low latency requirements such as the eIndustry PoC, see [10]. In this use case, for latency and security reasons, the vertical application runs in a server located into the vertical premises, so the UPF can be considered the natural border node for the abstraction view between the extended RAN dedicated for such vertical and the rest of core network that can be shared with other customers/clients.

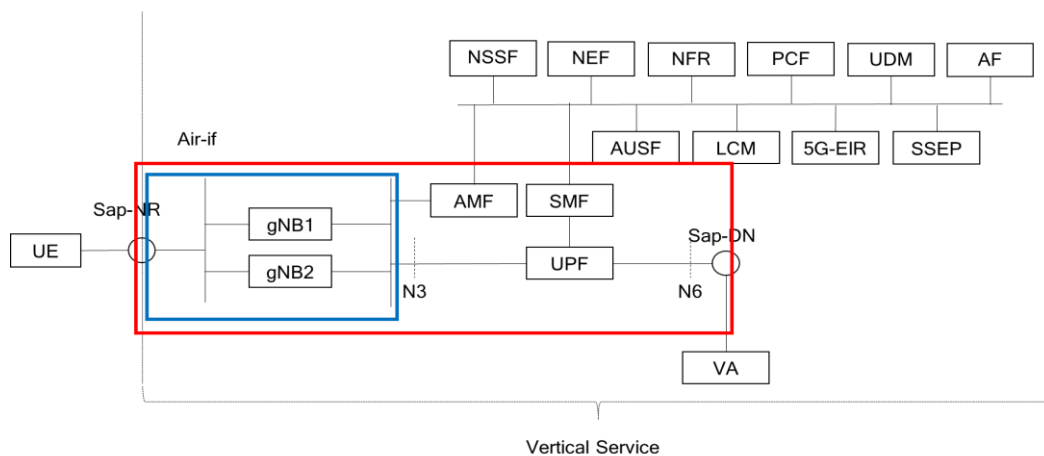


FIGURE 52: RAN ABSTRACTION APPROACHES

In the following, we do not distinguish between these two abstraction approaches; all the issues and the described solutions apply equally to both of them. Most importantly, for both approaches, the gNBs are included and there is a SAP representing the air interfaced (Sap-NR) as a demarcation of the vertical service. We refer to those CN functions, which are not included in the RAN abstraction as the core network. Note, although we have used NR in the diagrams above, similar abstractions can be made for other radio technologies such as LTE or WiFi.

Hereafter, we describe for each of the 5G-TRANSFORMER components - 5GT-VS, 5GT-SO and 5GT-MTP:

- Which information on the RAN is included in the descriptors used by the 5GT-VS;
- How this information is passed within the 5G-TRANSFORMER system to other components;
- How the 5GT-SO invokes the 5GT-MTP, which eventually controls the radio resources.

3.2.2 Radio Abstraction in the 5GT-VS

For a vertical defining a vertical service allowing devices to connect wirelessly, the most relevant property of the SAP for this connection is that it has a spatial extension. I.e. this SAP is associated with a coverage area. The vertical is not interested in how many BTSs or other wireless access points are needed to provide connectivity in such area, it just wants to define the area within which it wants its devices to be able to connect wirelessly to the vertical applications. In [7], we have described already a locationInfo attribute that can be used to define a coverage area associated with an SAP. For the sake of simplicity, we have limited ourselves to circular coverage areas, but by appropriate extension of the corresponding data type any shape of coverage areas could be described.

SLA requirements such as throughput per user, throughput per area, amount of users, and latency can be expressed already for vertical services as such. These SLA requirements apply implicitly to the RAN as well.

Although the vertical should not have to define detailed properties of the RAN, some abstract information might be required from the vertical. As an example, for some values of the SLA requirements the radio access technology (RAT) such as LTE, 5G, or Wi-Fi can be inferred from them. E.g. when the required RTT among UE and VA is 5ms, then mobility has to be provided by 5G. Neither LTE nor Wi-Fi could satisfy such an SLA requirement. In this case there is no need for the vertical to request a specific RAT. For values such as a RTT of 20ms both LTE and 5G would be applicable. In case the mobile devices used by the vertical support just one of these RATs, the vertical might actually want to require a specific RAT. This ensures that the requested coverage area would be provided with a RAT suitable for the devices of the vertical and on the other hand would prevent providing coverage with a RAT the vertical cannot use at all.

In the network service descriptors (NSD) passed from the 5GT-VS to the 5GT-SO, a required RAT can be encoded as a specific value of the layerProtocol attribute of the SAP and the Single Network Slice Selector Assistance Information (S-NSSAI) as an address of this SAP. This corresponds to defining a SAP with IP protocol and a specific IP address for a fixed line SAP.

After defining the main abstract information on the RAN such as the coverage area, optionally the RAT, and the more generic SLA requirements, these have to be described within the different descriptors used by the 5GT-VS and the other components of the 5G-TRANSFORMER system. Vertical Service Blueprints and descriptors can easily be extended with such information. For Network Service Descriptors according to ETSI NFV [22] we have described above or in [7] how the information can be encoded in the NSDs.

Figure 53 shows a graphical representation of a generic, nested, NSD for a vertical service including the RAN and CN functions. This representation is more abstract than the vertical services shown in Figure 51 and Figure 52 as the gNBs and the CN functions are represented by PNFS nested network services, resp. Note, for the CN we distinguish among functions expected to be instantiated per vertical service instance, such as the network exposure function (NEF) or the policy control function (PCF) and functions we expect to be commonly used for all vertical service instances using the 5G network of one operator. Examples of such functions are the network slice selection function (NSSF) and the 5G Equipment Identity Register (5G-EIR). The RAN itself is represented by a PNF.

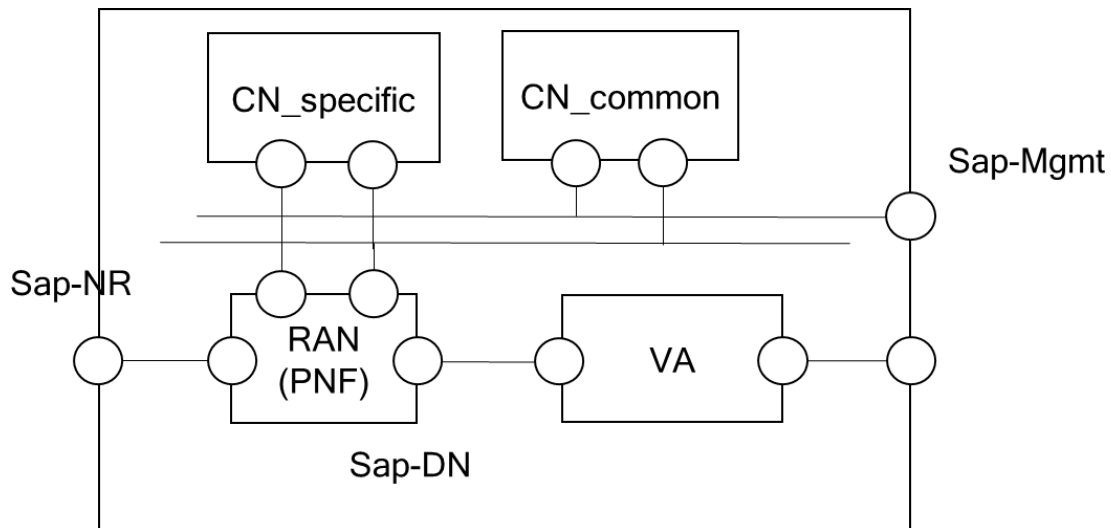


FIGURE 53: EXEMPLARY NSD FOR A VERTICAL SERVICE INCLUDING RAN

The 5GT-SO will orchestrate also the RAN part of the vertical services based on the information contained in NSD as shown above.

3.2.3 Radio Abstraction in the 5GT-SO

The 5GT-SO has to perform four tasks to instantiate a RAN for a vertical service. The 5GT-SO has to:

1. Match the coverage area as reported by the 5GT-MTP with the coverage area requested for a specific service and to instantiate a RAN in that specific area.
2. Request an additional RAN from a federated domain in case the coverage area reported by the 5GT-MTP is smaller than the requested one.
3. Handle the identifier of the RAN; this might require interaction of the 5GT-SO with some component in the RAN or CN.

4. Connect the RAN with other VNFs or NFV-NSs such as the CN or vertical applications.

Matching the provided coverage area by a 5GT-MTP and the requested one for a specific NSI and eventually instantiate the RAN uses a different approach than for other resources. E.g. for compute resources, a 5GT-MTP indicates to the 5GT-SO the available compute resources and the 5GT-SO decides on the compute resources in which datacentre to use for this NSI. Deciding on specific BTS and their configuration would be too complex for the 5GT-SO to handle, therefore this logic is kept within the 5GT-MTP. Instead the abstract view of a RAN as a coverage area is used. In a first step, the 5GT-MTP indicates the coverage area it can provide. Then, for a new NSI to be deployed, the 5GT-SO compares this with the requested coverage area, as expressed as locationInfo of the SAP corresponding to the air interface as well as with information on requested RATs. If the 5GT-MTP actually provides the correct RAT(s) and coverage area, the 5GT-SO requests the instantiation of the corresponding RAN from the 5GT-SO. Eventually, and if this instantiation was successful, the 5GT-MTP returns a list of connection points of the RAT to which the CN and the vertical applications could be connected.

Matching provided and requested coverage areas actually requires that the 5GT-SO is able to handle this notion of areas. This is a new concept for the 5GT-SO, but it is needed to decide whether the required coverage area is within the provided one. If it is not, the 5GT-SO has to determine which required area could not be covered and the 5GT-SO could try to get a RAN for the uncovered area instantiated from a federated domain. Similarly, if the 5GT-MTP provides several RANs, the 5GT-SO has to divide the required area accordingly and request the instantiation of several RANs from the 5GT-MTP. Although the coverage area can be seen as an attribute of a RAN resource similar to the bandwidth of a logical link, the operations on this attribute are different to adding bandwidth, storage capacities, etc.

The RAN identifier, e.g. the S-NSSAI, is passed from the 5GT-SO to the 5GT-MTP as address of the SAP. The 5GT-MTP can then configure this address to the actual gNBs. In the core network, this identifier needs to be configured to some network elements, such as the NSSF, as well. This can be provided either by the 5GT-SO to a controlling entity of the core network - which is not foreseen yet or by explicit configuration through external element management functionality.

Eventually, the 5GT-SO has to connect the RAN of a specific VSI with the other functions of the VSI. Here, the already available functionality of the 5GT-SO to connect child services of a composed service can be used.

3.2.4 Radio Abstraction in the 5GT-MTP

The main operations performed by 5GT-MTP are:

- Define the abstract view to expose to the 5GT-SO on the basis of several inputs that are out of the scope of this description (e.g. measurement, historical data, contract, etc.)
- Receive the request for a service on such coverage area in terms of service parameters as described in the previous section
- Translate the 5GT-SO request in commands for configuring the resources.

As previously described, two types of abstract view can be provided: one without core network functions and terminating at the 3GPP N3 interface, and the other with some core network functions and terminating at the N6 interface. In both cases, as reported in Figure 52, several interfaces to the other nodes managed at 5GT-SO level will be provided.

For sake of simplicity but without loss of generality let's consider the case of abstract view at N6 interface, as schematized in Figure 52 and which includes information about the coverage area and UPFs. The abstracted coverage area at 5GT-SO is the union of the several cells drawing the perimeter of the area as in Figure 54. To allow any shape, the perimeter is described by a set of ordered geographical points (P_1 - P_N). By connecting these ordered geographical points, the perimeter is identified. Then, RAN abstraction also includes information of the UPFs gateways (which expose N6 interfaces) based on the vertical service requirements (e.g., MVNO vs. other user services). The RAN abstraction is elaborated by the 5GT-MTP, which holds a detailed knowledge of the cells, DUs, CUs, and EPCs/UPFs, as described in the following. Note that the mapping of Radio into NFVIPoP is described in 5GT D2.3 [5].

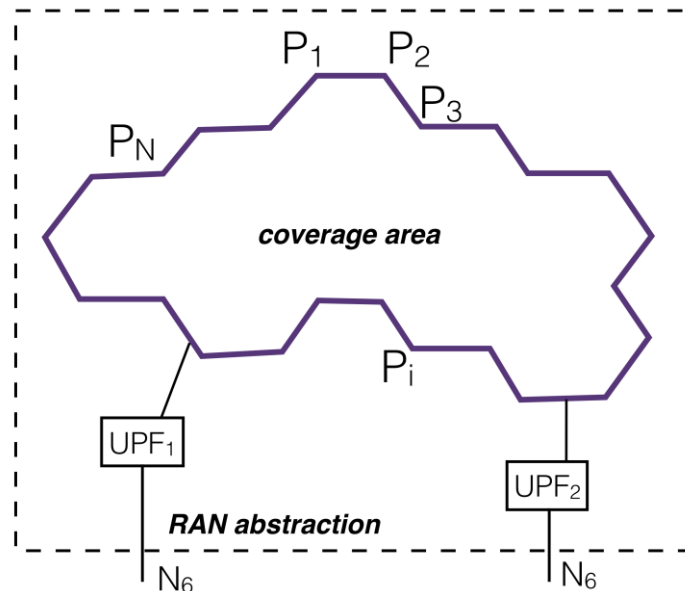


FIGURE 54: 5GT-SO VIEW OF THE RAN AND THE COVERAGE AREA

Table 95 shows the view at the 5GT-MTP of the radio coverage area (formerly introduced in [5] as RAN abstraction), including cells, DUs, CUs, and optionally of EPCs including User Plane Function (UPF¹), as shown in Figure 55. In the following we assume that the EPC is included in the radio abstraction. Thus, 5GT-MTP also has the view of the connectivity to the EPCs and such connectivity is in charge of the radio controller. The abstracted coverage area sent to the 5GT-SO is the union of the several cells drawing the perimeter of the area. Again, the perimeter of the area is described by a set of ordered geographical points, as illustrated in Figure 55 (P_1 - P_N) and reflected in Table 95 (line 2). The information model in Table 95 also includes bandwidth values,

¹ UPF is the External PDU Session point of interconnect to Data Network and it includes, among other functions, the support of packet routing & forwarding, packet inspection, QoS handling; UPF has part of the P-GW functionality from EPC in 4G.

which depend on the signal-to-interference-noise-ratio that is a function of several parameters such as the distance between user equipment and antennas and fading. Examples of propagation models can be found in the ITU Radio Communication Sector document ITU-R M.2412-0 [17] for several scenarios such as rural, urban, and indoor environments. Then, the signal-to-interference-plus-noise ratio imposes a proper modulation and code for transmission, which finally determines the bandwidth. B-max is defined as the maximum among the maximum bandwidth values involving all the cells. B-min is the minimum bandwidth at the perimeter.

TABLE 95: ABSTRACTION OF THE RADIO COVERAGE AREA

Parameter	Description
Id	Identifier of the area
geographicalAreaInfo	It provides information about the covered area expressed as a list ordered points identifying a perimeter. Such points are identified with geographical coordinates (latitude and longitude)
B-min	Minimum bandwidth (Mb/s) that can be allocated to a service, e.g. along the perimeter
B-max	Maximum bandwidth (Mb/s) that can be allocated to a service in the area, e.g. in proximity of an antenna
Gateways	List of IP addresses associated to UPFs gateways
RAT	RAT type, encoded as layerProtocol as defined for Cpd, [20]
RanIdentifier	'address' of the RAN, corresponding to address parameter in SapData [21]

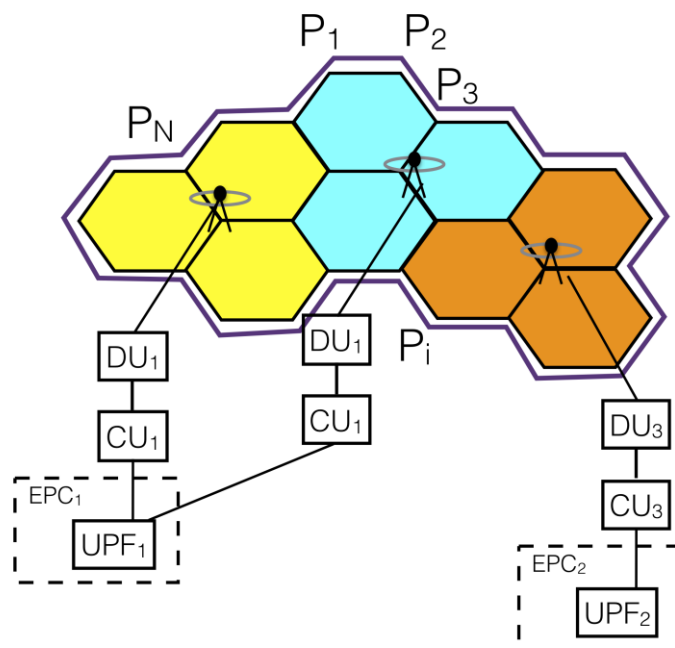


FIGURE 55: 5GT-MTP VIEW

3.3 MEC integration

In this section, we describe how MEC support is built into the 5G-TRANSFORMER architecture. In particular, we describe the role of each entity (5GT-VS, 5GT-SO, 5GT-MTP) in the end-to-end onboarding and instantiation workflows, discuss specific aspects, such as implementing traffic offloading and managing MEC service requirements, and present the necessary descriptor-level extensions in order to integrate MEC applications in 5G-TRANSFORMER service definitions.

3.3.1 MEC capabilities in 5G-TRANSFORMER descriptors, and VS-level operations

3.3.1.1 VNFD vs. AppD

To ease the integration, VNFs and MEC applications are treated in 5G-TRANSFORMER in a similar manner; but with some distinct differences. These differences are reflected in the onboarding and instantiation workflows, as well as in the descriptors of VNFs and MEC applications, vertical and network services.

Based on the related ETSI standard specifications, a critical difference between a MEC application and a VNF is the fact that the former is standalone, and can be considered to be equivalent to a VNF's VDU. Therefore, the notion of a VNF-internal virtual link (VL) is not relevant for MEC applications. Furthermore, a MEC application does not support flavors and instantiation levels. This requires some extensions at the NSD level to support VNFs and MEC applications in a unified manner.

A NSD includes one or more VNF profiles (VnfProfile information element) per VNF, which are linked to the respective VNFD via the vnfId field. A VNF profile also includes a flavourId and an instantiationLevel, elements which are defined inside the VNFD. For MEC applications, however, there are no flavourId and instantiationLevel references defined in the AppD. To deal with this issue, 5G-TRANSFORMER proposes a specific NSD extension. In particular, it introduces a MecAppProfile information element in the NSD definition, with the same semantics as those of Vnfprofile, but without the flavourId and instantiationLevel fields. This can be seen as an adaptation of the PnfProfile information element that is defined for physical network functions included in a network service [22], since this is closer conceptually to a MEC application profile (no support for flavours and instantiation levels, "static" resources).

Furthermore, an appToLevelMapping field has been added to the nsInstantiationLevel element of the NS deployment flavor. In a similar fashion as with vnfToLevelMapping, it can be used to specify the number of MEC application instances to launch.

3.3.1.2 Handling location constraints

At the VS level, a MEC application package including an AppD is on-boarded and included in an NSD in a similar manner as with classical VNFDs. To apply location constraints on the placement of MEC applications, the approach described below is followed. Note that this is similar to the approach for implementing location constraints for regular VNFs. The focus here is on the actual information elements that have to be used in order to implement this support for MEC applications:

- The NSD includes a list of VNFDs, AppDs, and SAPDs.
- An AppD includes one or more external connection points (*appExtCpd* field in the AppD). For presentation simplicity, we consider a scenario where there is a single

external connection point in the AppD. This CP can be considered to denote the *edge*. In practice, other external connection points may be present in an AppD, to implement the other virtual links of a MEC application.

- The connection point of the MEC application is associated with one of the SAPDs included in the NSD. This takes place by the VS setting appropriately the *associatedCpdId* field of the SAPD to the value of the *appExtCpd* of the MEC application descriptor.

The above take place at onboarding time. Note that no location information has yet been encoded; the latter happens at instantiation time. In particular, in 5G-TRANSFORMER, the *SapData* information element that is included in a VSI/NSI has been extended with a *LocationInfo* field that encodes geographical coordinates and a radius, thus indicating a region of coverage. At MEC application instantiation, the VS adds in the request a *SapData* element which includes (i) the identifier of the SAPD to which the MEC app connection point has been associated, and (ii) a *LocationInfo* element, which indicates the geographical location where to deploy the MEC application. The relationships between AppD, SAPD, and NSI are shown in Figure 56.

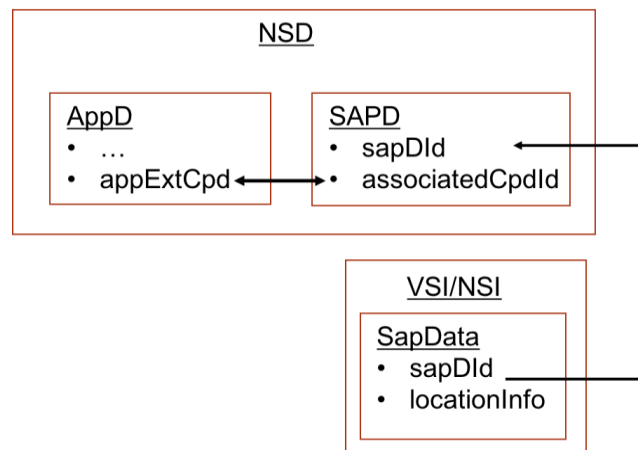


FIGURE 56: RELATIONSHIP BETWEEN APPD, SAPD, AND NSI WITH A FOCUS ON LOCATION CONSTRAINTS

The *LocationInfo* element in 5GT is indicated by the Vertical when instantiating the MEC application. *LocationInfo* is composed of three coordinates X,Y,Z and a radius R; X,Y,Z correspond to the GPS coordinates of a central point, where R is the radius of the area.

3.3.1.3 Dealing with virtual links

It should be noted that the VnfProfile and PnfProfile elements include an nsVirtualLinkConnectivity field, which can be used to map connection points of a VNF or PNF to virtual link profiles defined in the NSD, and therefore associate the connection points of VNFs or PNFs to VLs. A MEC application can also have virtual links to VNFs, PNFs, or other MEC applications. In the case of MEC applications, however, these connection points can only be external, and are defined in the AppD's appExtCpd list. Our approach to associate MEC application CPs with VLs is the following:

- A mecAppProfile element has been introduced to the NS deployment flavor element of the NSD (nsDf), corresponding to a MEC application. A mecAppProfile has a

mapping of cpdlids to virtualLinkProfileids. The latter are parts of the nsVirtualLinkConnectivity information element.

- Each CP included in the AppD is linked in the NSD to a virtualLinkProfileid.

Note that there is a special appExtCpd in the AppD, which is used for implementing location constraints; this is not included in the above mapping procedure.

Figure 57 presents how information elements of the NSD, VNFD, and AppD are linked, as well as the necessary extensions to include MEC application profiles and handle VLs that involve MEC applications.

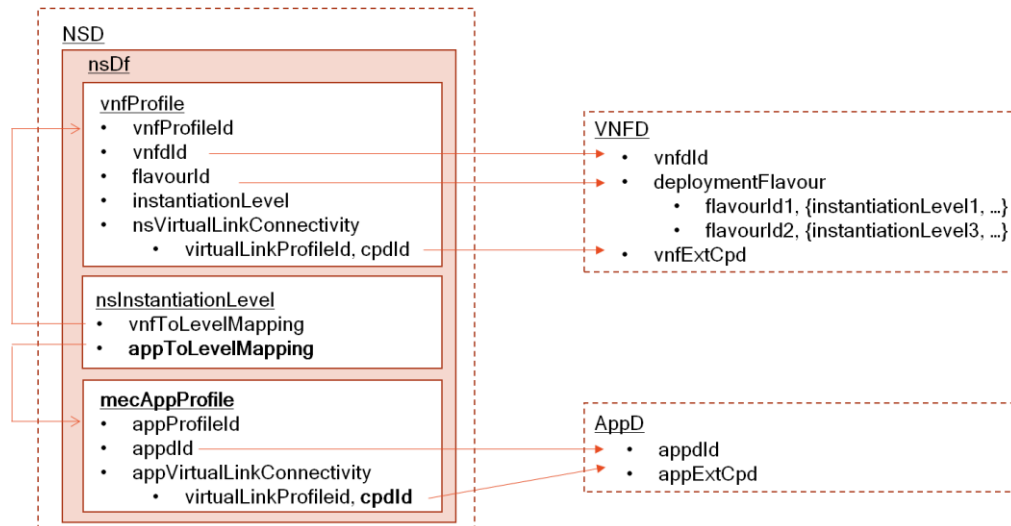


FIGURE 57: DESCRIPTOR EXTENSIONS FOR MEC SUPPORT

3.3.2 SO-level operations

At the SO level, when receiving the request to deploy a service including a MEC application, the SO parses the NSI to find and extract AppD references and *SapData* elements, in order to correlate each one of the latter with one of the SAPDs included in the NSD, and in turn “link” it with the corresponding AppD external connection point using the *associatedCpdId* of the SAPD. Note that this procedure *has to be implemented anyway for regular VNFs with location constraints*.

The second step consists in selecting the NFVI-PoP where to deploy the MEC application. To do so, the Placement Algorithm (PA) running at the SO will use the *LocationInfo* included in the NSI, some AppD fields (compute requirements, as with regular VNFs), and the information provided by the MTP on the different NFVI-PoPs. For each NFVI-PoP, among others, the MTP information includes the following elements: RegionID, resources (CPU, storage, etc), MEC capabilities. The RegionID encodes the GPS coordinates of the NFVI-PoP, Resources indicate the available resources at the NFVI-PoP, and a binary value “MEC” indicates if the NFVI-PoP is an edge cloud or not (e.g., if it has MEC capabilities).

The SO PA will select the appropriate NFVI-PoP where to instantiate the MEC application. The SO PA needs to take into consideration the location, the “edge” flag and the available resources of the NFVI-PoP, as also in the case of regular VNFs. If no appropriate NFVI-PoP is found to instantiate the MEC application, the SO PA rejects the MEC application and the VS is informed about the decision. Otherwise, the SO

requests the deployment of the MEC application to the MTP. This will follow the same process as for classical VNFs.

3.3.3 MTP-level operations

3.3.3.1 MEC application placement

At the MTP level, when receiving the request to instantiate a MEC application in a specific NFVI-PoP, the MTP uses a specific intra-pop Placement Algorithm (PA) algorithm to select the appropriate edge host to run the MEC application. Our design aims to keep the SO with minimal MEC-awareness (i.e., only knowledge of whether an NFVI-PoP supports MEC), and instead push MEC specific operations (i.e., functionality that depends on information found in the AppD, such as latency-aware placement and implementation of traffic redirection) to the MTP, that handles them transparently.

The MTP edge PA algorithm may consider resources but also host/DC latencies to ensure that the *appLatency* of the application expressed in the AppD is not exceeded. In other words, the *appLatency* field of the AppD is ignored by the SO and is handled internally by the MTP to select the appropriate host (intra-PoP) which can guarantee the required UE-to-MEC application latency. This latency is not straightforward to measure, but a reasonable estimate can be assumed to be available at the MTP level. This could correspond to the latency between the eNodeB(s) that are served by a MEC host and the host where the MEC application is deployed.

The MTP edge PA may be similar to the one used to select the NFVI-PoP or an adaptation that adds latency as constraint. The MTP will then start by instantiating the MEC application by using the VIM plug-in, and then access the MEC plug-in to install the traffic redirection rule as indicated in the AppD *appTrafficRule* elements, and DNS rules as specified in the *appDNSRule* elements.

3.3.3.2 Implementing traffic offloading

To implement traffic offloading towards a MEC application instance, information found in the AppD (description of traffic flows that should be offloaded, DNS names that need to be resolved to the MEC application instance) and information about the actual MEC application instance (IP address of the instance) need to be combined. The 5GT-SO needs to carry out a two-step procedure:

- Request the 5GT-MTP to launch a MEC application instance and retrieve the identifier, IP address and other information about the instance.
- Request the 5GT-MTP to set up the appropriate MEP configuration for the instance. In this case, it sends to the 5GT-MTP the identifier of the instance, its IP address, and a reference to the corresponding AppD, which is onboarded to the 5GT-MTP.

The 5GT-MTP then retrieves the AppD from its internal database, checks if *appTrafficRule* or *appDNSRule* elements are included there, and *complements them* with the information about the MEC application instance (IP address) that is necessary to pass on to the MEC plugin to apply the appropriate MEP configuration over the 5GT-MTP's SBI towards the MEC plugin. For instance, if in the AppD there is only an *appDNSRule* to register a specific DNS name for the corresponding MEC application, one *appTrafficRule* that matches a specific traffic flow, and no *appServiceRequired* or *appServiceProduced* elements, the request towards the MEC plugin to apply the MEP configuration should look like the following. *RegionId* includes the identifier of the

region where the MEC application is instantiated, in order for the MEC plugin to retrieve the API endpoint of the MEP responsible for this region. The highlighted addresses are added by the 5GT-MTP after successful instantiation and belong to the MEC application instance.

```
{
  "RegionId": "2",
  "appTrafficRule": [
    {
      "trafficRuleId": "tr1",
      "filterType": "FLOW",
      "priority": 0,
      "trafficFilter": [
        {
          "srcAddress": [
            "208930100001114"
          ],
          "dstAddress": [
            "195.251.255.142"
          ],
          "dstPort": [
            "80"
          ],
          "Protocol": [
            "tcp"
          ]
        }
      ],
      "action": "FORWARD",
      "dstInterface": [
        {
          "interfaceType": "IP",
          "dstMACAddress": "00:11:22:33:44:55",
          "dstIPAddress": "172.24.248.5"
        }
      ]
    }
  ],
  "appDNSRule": [
    {
      "dnsRuleId": "111",
      "domainName": "pfrag.gr",
      "ipAddressType": "IP_V4",
      "ipAddress": "172.24.248.5",
      "ttl": 0
    }
  ]
}
```

3.3.3.3 MEP discovery

A MEC application is assumed to be able to discover the API endpoint of the MEC platform in order to consume services such as the RNIS. This is not specified in the relevant standard document (ETSI MEC 011 [23]), and we assume that it is implemented by DNS: the MEC application requests the resolution of a well-known

domain name (e.g. “mep.mec”), to which the local DNS replies with the IP address of the MEC platform.

3.3.4 Service onboarding

With respect to service onboarding, the main differences compared with the case of regular VNFs are the following:

- A MEC application package is onboarded only as part of a Virtual Application (VA), a process initiated by the vertical. As in the case of VNFs, the onboarding of MEC application precedes the onboarding of an NFV-NS, which is initiated by the vertical slicer.
- A MEC application is also onboarded at the 5GT-MTP. This takes place via an additional exchange between the 5GT-SO and the 5GT-MTP, after the MEC application has been onboarded at the 5GT-SO. The AppD of the MEC application is maintained in the AppD catalogue and AppD database of both the 5GT-SO and the 5GT-MTP, respectively.

This design option was followed to simplify the implementation of the 5GT-SO and the SO-MTP interface. There are specific configuration options, such as the implementation of traffic offloading and the provision/consumption of MEC services that are kept outside the scope of the 5GT-VS and 5GT-SO and are considered 5GT-MTP-internal.

The onboarding workflow is detailed in D1.3 (Section A7.1) [5].

3.3.5 End-to-end instantiation workflow

The implementation and configuration of MEC-specific functionality is handled internally in the MTP, via interactions with the MEC plugin. The MEC plug-in is in charge of requesting the MEC platform to apply traffic redirection, DNS rule management, and other MEC platform configurations for the MEC applications. We assume that the MEP and the vEPC functions (particularly the SPGW-U) are already deployed at the edge NFVI. We further assume a multi-tenancy environment where the MEP and vEPC are shared among the 5G-TRANSFORMER network slices. Furthermore, as mentioned before, we have opted for a solution where MEC application descriptors are also onboarded and maintained by the MTP.

Therefore, when the 5GT-SO requests the deployment of a MEC application at a specific edge NFVI, it includes in the request a reference to the corresponding AppD. The 5GT-MTP may then extract the relevant information elements from the AppD to apply MEC-specific configurations. These elements are *appTrafficRule*, *appDNSRule*, *appServiceRequired*, *appServiceProduced*, and *appLatency*. The first indicates the type of traffic that needs to be offloaded to the MEC application, the second includes a set of DNS rules to be added to the MEC DNS service,² while the third and fourth indicate the MEC service(s) that the MEC application requires and provides, respectively. The *appLatency* element represents the maximum latency tolerated by the MEC application and can be used by the 5GT-MTP for intra-PoP application placement or other optimizations. The exact semantics of *appLatency* are not specified in the related ETSI MEC 010-2 standard [24]. In 5G-TRANSFORMER, we consider this field to refer to the maximum latency from the UE to the MEC application.

² This can be used for DNS-based traffic offloading, by resolving DNS requests from UEs to the newly created MEC application instance.

Figure 58 presents a high-level view of the process of instantiating a NS that includes MEC applications. The 5GT-VS initiates the process of network service instantiation³ by communicating to the SO the NSI identifier that corresponds to an onboarded NSD which includes AppD references (1). Location constraints for the placement of MEC applications are included in the instantiation request, and in particular in the *SapData* information element.

Upon the reception of this message, the 5GT-SO executes the necessary steps required to derive a resource allocation and placement decision for the VNFs and MEC applications included in the NS. These steps include the execution of a placement algorithm that also satisfies location constraints and considers in its decision the MEC capabilities of the underlying NFVI-PoPs, as these are exposed by the 5GT-MTP. For a MEC application, the 5GT-SO sends to the 5GT-MTP a deployment request (2) which includes, among others, the MEC application's AppD identifier and the desired location.

The 5GT-MTP then requests (3) the instantiation of the MEC application to the VIM plug-in, by indicating the location (i.e., the selected edge NFVI, also including the precise location information specified in the NS instantiation request) and the image of the application. It is assumed that the image is already on-boarded at the edge NFVI. Note that, at this level, the 5GT-MTP uses an intra-PoP internal placement algorithm that selects the appropriate NFVI to host the MEC application instance, which supports the maximum tolerated latency and provides the MEC service(s) requested by the MEC application. This information is obtained from the fields found in the AppD stored at the 5GT-MTP, specifically *appLatency* and *appServiceRequired*. The VIM plug-in requests (4) the creation of the application instance at the selected edge NFVI and receives (5) the identifier and the IP address of the instantiated application. Then, it acknowledges (6) the creation of the instance to the 5GT-MTP, indicating this IP address and application instance identifier. This confirmation is eventually returned to the 5GT-SO (7). As soon as the MEC application instance compute resources have been allocated and the instance is up and running, the 5GT-SO requests the 5GT-MTP to install the VFs between the MEC application and other MEC applications, VNF instances, or PNFs (8).

With an additional message (10), the 5GT-SO requests the 5GT-MTP to set up any necessary MEC-level configuration, as this is specified in the application's AppD. In this message, the AppD identifier, the application instance identifier, and the IP address information of the MEC application are included. The 5GT-MTP then retrieves the AppD that corresponds to the given identifier, extracts the *appTrafficRule*, *appDNSRule*, *appServiceProduced* and *appServiceRequired* fields, complements them with information about the IP address of the instance, which is not available in the AppD at onboarding time, and requests the configuration of the MEP for the new MEC application to the MEC plug-in (11).

The MEC plug-in, which acts as a MEC platform manager (MEPM), as per the ETSI MEC architecture specification [25], applies the requested configuration to the MEP over the Mm5 reference point. This involves a number of message exchanges between the MEC plug-in and the MEP. First, the MEC plug-in requests the setup of traffic rules (12). In response, the MEP uses the Mp2 interface to apply traffic redirection rules to the SPGW-U deployed at the edge (13). In a similar fashion, the MEC plug-in

³ For simplicity, the message to create a NSI identifier that precedes *InstantiateNs* is not shown.

configures the MEC DNS service appropriately (16) based on the presence of the *appDNSRule* field in the request, and notifies the MEP of potential services provided or required by the MEC application (18). After the successful execution of the above steps, the MEC plug-in confirms the successful MEP configuration to the 5GT-MTP (20). The latter then acknowledges the deployment and configuration of the MEC application to the 5GT-SO (21).

The above procedure is repeated for each AppD referenced in the NSD that corresponds to the NSI that is instantiated. The successful creation of the NSI is eventually signaled to the 5GT-VS (22).

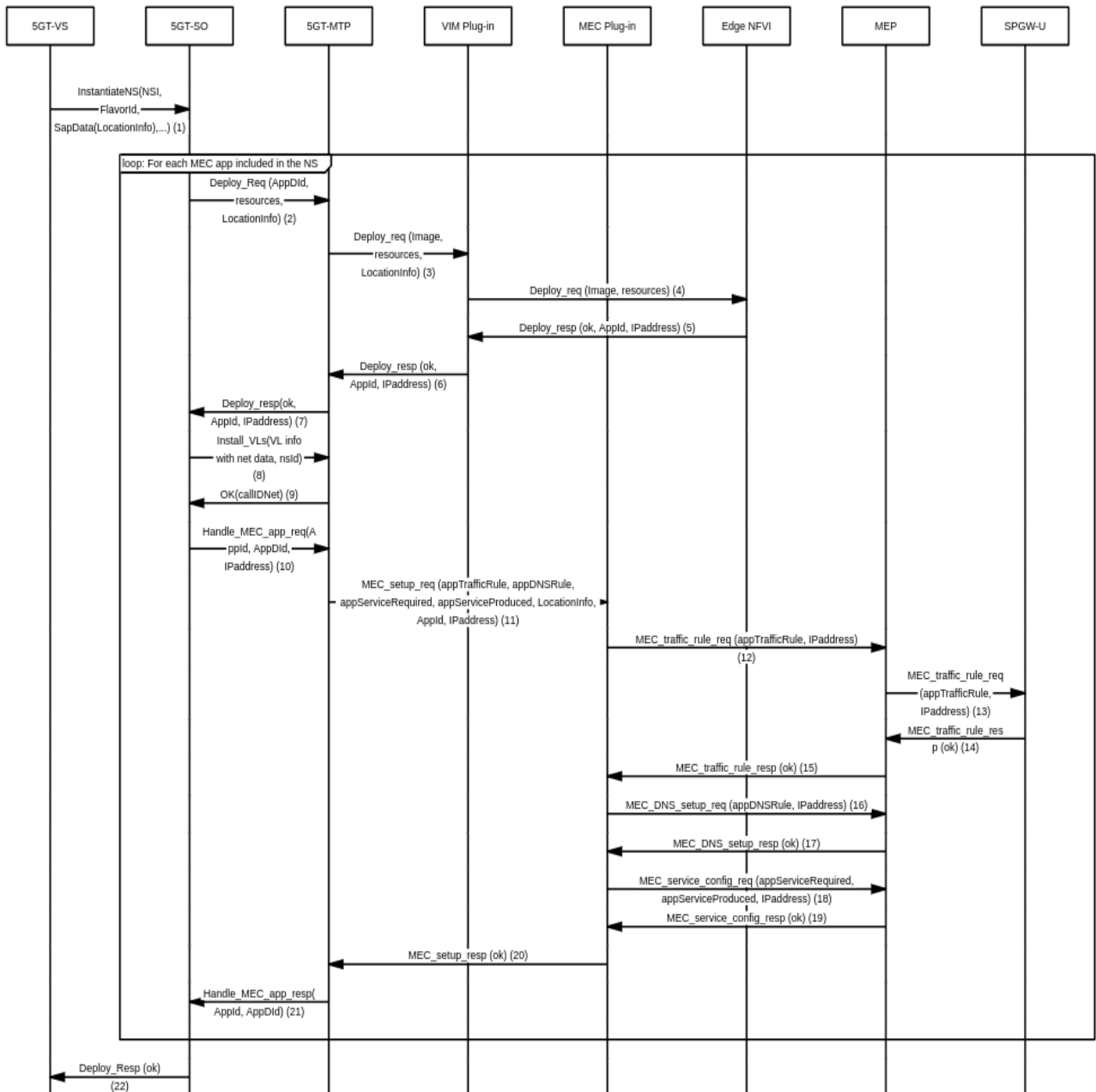


FIGURE 58: WORKFLOW OF DEPLOYING A NETWORK SERVICE INSTANCE THAT INCLUDES MEC APPLICATIONS

4 Conclusions

In the first deliverable of the project D1.1 [1], we have proposed an initial study of the stakeholders and target (vertical) services of the 5GT system as well as a bench of business and functional requirements that have ultimately driven the architectural design of the 5GT architecture. We proposed in this deliverable to complete this study with a detailed techno-economic analysis that investigates deeper the cost models of the 5GT vertical use case services selected for implementation in the final proof of concepts, in order to evaluate and/or promote some new business models for the vertical actors involved in the project. The study presented in section 2, contained; an analytical part that allowed modelling the monetary flow for each vertical use case and the definition of the different kind of costs that constitutes the 5GT service price; and an experimental part in which we performed simulations of the vertical services deployment in order to estimate over different scenarios the deployment costs. Simulations were performed on a MATLAB based tool developed for the purpose of the study and described as part of section 2.

The study is finally concluded with some recommendations that analyses the extent of the economic benefits of business models that includes federation between different stakeholders, and the economic consequences of the use of cloudified and virtualized infrastructure resources on the profitability of services.

The second part of this deliverable presented in section 3 is the final 5G-TRANSFORMER architecture design in which we summarize the latest enhancements of the platform with regards to the refined design proposed in D1.3 [3]. The two main proposed extensions are RAN abstraction and MEC support. For both features, we provided a detailed description of the extensions to the design and operations of the 5GT platform main building blocks; the 5GT-VS, the 5GT-SO and the 5GT-MTP.

5 References

- [1] 5G-TRANSFORMER, D1.1, Report on vertical requirements and use cases, November 2017.
- [2] 5G-TRANSFORMER, D1.2, 5G-TRANSFORMER initial system design, May 2018.
- [3] 5G-TRANSFORMER, D1.3, 5G-TRANSFORMER refined system design, May 2019
- [4] 5G-TRANSFORMER, D2.1, Definition of the Mobile Transport and Computing Platform, March 2018.
- [5] 5G-TRANSFORMER, D2.3, Final design and implementation report on the MTP, May 2019.
- [6] 5G-TRANSFORMER, D3.1, Definition of vertical service descriptors and SO NBI, March 2018.
- [7] 5G-TRANSFORMER, D3.3, Final design and implementation report on the Vertical Slicer, May 2019.
- [8] 5G-TRANSFORMER, D4.1, Definition of service orchestration and federation algorithms, service monitoring algorithms, March 2018.
- [9] 5G-TRANSFORMER, D4.3, Final design and implementation report on service orchestration, federation and monitoring platform, May 2019.
- [10] 5G-TRANSFORMER, D5.4, 5G-TRANSFORMER Reports on trials results, November 2019.
- [11] "Gartner Hype Cycle for Cloud Computing, 2018," Amazon Web Services, Inc. [Online]. Available: https://pages.awscloud.com/Analyst_Reports_Gartner-Hype-Cycle-for-Cloud-Computing.html?trk=ar_card. [Accessed: 21-Oct-2019].
- [12] C. Wu, R. Buyya, and K. Ramamohanarao, "Cloud Pricing Models," ACM Computing Surveys, vol. 52, no. 6, pp. 1-36, 2019.
- [13] P. Iovanna, F. Cavaliere, F. Testa, S. Stracca, G. Bottari, F. Ponzini, A. Bianchi, and R. Sabella, "Future Proof Optical Network Infrastructure for 5G Transport," IEEE/OSA Journal of Optical Communications and Networking, vol. 8, no. 12, pp. B80-B92, 2016.
- [14] 3GPP TS 22.261 Version 15.5.0 Release 15, Jul. 2018, [online] Available: https://www.etsi.org/deliver/etsi_ts/122200_122299/122261/15.05.00_60/ts_122261v150500p.pdf.
- [15] Kenton, W. (2019). Sensitivity Analysis Definition. [online] Investopedia. Available at: <https://www.investopedia.com/terms/s/sensitivityanalysis.asp> [Accessed 29 Sep. 2019].
- [16] <https://www.cisco.com/c/en/us/products/collateral/switches/nexus-7000-series-switches/white-paper-c11-737022.html>
- [17] "Guidelines for evaluation of radio interface technologies for IMT-2020," Report ITU-R M.2412-0.
- [18] <https://www.penguincomputing.com>
- [19] ETSI GS NFV-EVE 003 v1.1.1 (2016-01)
- [20] ETSI GS NFV-IFA 011, V2.3.1, Management and Orchestration; VNF Packaging Specification, 2017.
- [21] ETSI GS NFV-IFA-013, V3.1.1, Management and Orchestration; Os-Ma-Nfvo reference point - Interface and Information Model Specification, 2018.
- [22] ETSI GS NFV-IFA 014, "Network Functions Virtualization (NFV) Release 2; Management and Orchestration; Network Service Templates Specification," v2.5.1, August 2018.
- [23] ETSI GS MEC 011, "Mobile Edge Computing (MEC); Mobile Edge Platform Application Enablement," v1.1.1, July 2017.

-
- [24] ETSI GS MEC 010-2, "Mobile Edge Computing (MEC); Mobile Edge Management; Part 2: Application lifecycle, rules and requirements management," v1.1.1, July 2017.
 - [25] ETSI GS MEC 003, "Multi-access Edge Computing (MEC); Framework and Reference Architecture," v2.1.1, January 2019.

A.1 Infrastructure cost modelling examples

In section 2.2.2, we presented guidelines for assessing the cost of an infrastructure taking into account various parameters. These parameters include the composition of a datacentre in network equipments, computing servers and storage elements as well as the rental of space for racks, wages of maintenance staff and the electricity consumption. For the sake of simplicity, this annexe contains further numerical applications from the equations derived in section 2.2.2 to calculate the total cost of a datacentre according to capacity (small, medium and large) and based on the actual prices of hardware, rent, salaries and energy power. The results of these calculations allow us to know the cost of a virtualized resource (vcpu, ram, and disk) as well as the ratio of each expense in the overall cost according to the size of the infrastructure.

A.1.1 Management switch pricing

A.1.1.1 Commercial Spine Switch Example

The Arctica 3200xlp is a network switch for the aggregation spine layer. It provides 32 ports of 40 Gigabite Ethernet connectivity. Commercialized by the end of 2014, this spine switch costs \$11 093, it has the following characteristics: Provides 32 ports wherein 2 of them are used for MLAG; It does not provide aggregation ports; This switch consumes around 300 Watts per hour; and contains one rack unit. If we consider 2 spine switches in the infrastructure, we can estimate the costs for these two switches using the Equations in Table 15. The results are shown in Table 96

TABLE 96: COSTS FOR TWO ARCTICA 3200XLP COMMERCIAL SPINE SWITCH

Entry	Value
Number of spine ports	60
Cost for the spine switches per month	\$674.94
Cost for the rack units per month	\$57.14
Cost for the power per month	\$109.58
Total cost per month	\$841.66
Cost per core port per month	\$14.03

A.1.1.2 Commercial Leaf Switch Example

Arctica 4806xp is another open network switch using the Broadcom StrataXGS Trident II chipset. This switch provides 10/40 Gigabit Ethernet Top-of-Rock (TOR) open switch, a combination of 6 ports with 40 Gbps and 48 ports of 10 Gbps. With this configuration, the switch supports cost-efficient intra-rack connectivity using industry standard Twinax cables, while offering multiple 40 Gigabit copper and optical links options for distribution layer connectivity. With the features supporting VxLAN, and native data plane, this switch is suited for modern scale-out architectures and virtualized multi-tenant environments found in service provider data centres. As additional configuration, this switch has 6 aggregation ports from the 48, 2 of the 48 ports are used for the MLAG, and 2 others are used to connect to spine, this switch is contained in one rack with a power consumption of 305 Watts per hour. For such configuration, the switch costs \$9 752. By using 2 Arctica 4806xp switches, we will estimate the costs per month using equations from Table 16. The results are presented in Table 97.

TABLE 97: COSTS FOR TWO ARCTICA 4806XP COMMERCIAL LEAF SWITCH

Entry	Value
Max number of leaf switches that can be connected to the spine	30
Max number of leaf ports in this configuration	1440
Number of leaf ports	96
Cost for the leaf switches per month	\$593.35
Cost for the rack units per month	\$57.14
Cost for the power per month	\$111.40
Cost for the spine ports per month	\$28.06
Total cost per month	\$789.95
Cost per core port per month	\$8.23

A.1.1.3 Commercial Management Switch Example

Arctica 4804ip Network Switch is a cost-effective 48 port 1GbE Layer 2+ network switch in a compact 1U form factor, suited for data centre Out-Of-Band (OOB) networks. The switch implements an x86-based control plane to ease the integration of linux and Openstack automation tools, it comes ready to use with a pre-installed Cumulus RMP (Rack Management Platform) software. With a cost of \$1 327, this management switch is configured with 48 ports providing 1Gbps Ethernet, and 4 aggregation ports with 10Gbps Ethernet. The switch consumes 52Watts/hours, and it is holding on one rack unit. We may use 2 ports to connect to the spine and 0 ports for MLAG.

Similarly, as for Spine and Leaf switches estimations, by considering 2 management switches and the equations in Table 19, we can estimate the costs for management switches per month. The results are shown in Table 98:

TABLE 98: COSTS FOR TWO ARCTICA 4804IP MANAGEMENT SWITCHES

Entry	Value
Number of management ports	96
Cost for the management switches per month	\$80.74
Cost for the rack units per month	\$57.14
Cost for the power per month	\$18.99
Cost for the leaf ports per month	\$16.46
Total cost per month	\$173.33
Cost per management port per month	\$1.81

A.1.2 Network Node Pricing

We present herein an example of a network node. Its cost is composed of the cost of several components as presented in 2.2.2. Such a network node is holding on one rack unit, with a consumption of 180Watts/hour, it provides 4 leaf network ports and 1 for the management. This network node is about \$3 412.00. The estimated costs represented by Equation 23 to Equation 28 are shown in Table 100.

TABLE 99: BILL OF MATERIALS FOR NETWORK NODE

Part name	Part number	Cost per unit	Quantity	Subtotal
1U Dual Xeon Intel Server Platform, 8 x 2.5" HDDs, Dual PSU	R1208WTTGSR	\$1473.00	1	\$1473.00

2nd 750W PSU for R1208WTTGS server	FXX750PCRPS	\$220.00	1	\$220.00
Premium Rail Kit	AXXPRAIL	\$103.00	1	\$103.00
Intel Xeon E5-2623V4 2.60GHz 10MB LGA2011-3 4C/8T	CM8066002402400	\$489.00	2	\$978.00
Intel Certified 16GBDDR4 2133MHz ECC Reg DIMM	KVR21R15D4/16	\$157.00	2	\$314.00
Intel Ethernet X540- T2 Dual 10-Gb PCIe x8 Card	X540-T2	\$194.00	1	\$194.00
300GB Seagate Savvio 10k rpm SAS Drives	ST300MM0026	\$65.00	2	\$130.00
Subtotal per node				$d_1 = \$3412.00$

TABLE 100: COSTS ESTIMATED FOR THE CHOSEN NETWORK NODE DEFINED IN TABLE 13

Costs	Values
Cost per month for the hardware	\$103.80
Cost per month for the rack space	\$28.57
Cost per month for the power	\$32.87
Cost per month for the leaf network ports	\$32.91
Cost per month for the management network ports	\$1.81
Cost per month per controller node	\$199.96

A.1.3 Controller node pricing

An example of controller node is shown in Table 101.

TABLE 101: CONTROLLER NODE'S MATERIAL BILL

Part name	Part number	Cost per unit	Quantity	Subtotal
2U Dual Xeon Intel Server platform, 12 x 3.5" HDDs, Dual PSU	R2312WTTYSR	\$1855.00	1	\$1855.00
2nd 1000W PSU for R2312WTTYSR server	AXX1100PCRPS	\$233.00	1	\$233.00
Premium Rail Kit	AXXPRAIL	\$103.00	1	\$103.00
Intel Xeon E5-2697V4 2.30GHz 45MB LGA2011-3 18C/36T	BX80660E52697V4	\$2470.00	2	\$4940.00
Intel Certified 32GBDDR4 2133MHz ECC Reg DIMM	KVR21L15Q4/32	\$335.00	4	\$1340.00
Intel Ethernet X540-T2 Dual 10-Gb PCIe x8 Card	X540-T2	\$194.00	1	\$194.00
Intel Dual 10Gbps	AXX10GBTWLIOM	\$413.00	1	\$413.00

X540-BT2 I/O Module				
1TB Seagate Constellation SATA 7200rpm HDD 2.5"	ST91000640NS	\$70.00	2	\$140.00
480GB Intel SSD S3610 series SATA HDD 2.5"	SSDSC2BX480G401	\$414.85	2	\$829.70
Subtotal per node				\$10047.70

A.1.4 Compute node pricing

Considering the example of compute node composed of the components presented in 2.2.2 the total cost for this compute node is about \$15 523.00. Its configuration is composed of 1 rack unit, 2 CPU sockets, 18 cores for each CPU with 2 threads, the overcommit ratio is equal to 4. The compute node provides 768GB of RAM, with a consumption of 350Watts/hour. In addition, 12GB of the Ram is reserved to the hypervisor, the RAM overcommit ratio is equal to 1. The CPU and RAM pricing weight are equal to 50%. This compute node provides 4 leaf network nodes, and 1 for the management. In total, we may request for 288 vCPU, and 756GB RAM.

TABLE 102: COMPUTE NODE'S BILL OF MATERIALS

Part name	Part number	Cost per unit	Quantity	Subtotal
1U Dual Xeon Intel Server Platform, 8 x 2.5" HDDs, Dual PSU	R1208WTTGSR	\$1473.00	1	\$1473.00
2nd 750W PSU for R1208WTTGS server	FXX750PCRPS	\$220.00	1	\$220.00
Premium Rail Kit	AXXPRAIL	\$103.00	1	\$103.00
Intel Xeon E5-2697V4 2.30GHz 45MB LGA2011-3 18C/36T	BX80660E52697V4	\$2470.00	2	\$4940.00
Intel Certified 32GBDDR4 2133MHz ECC Reg DIMM	KVR21L15Q4/32	\$335.00	24	\$8040.00
Intel Ethernet X540-T2 Dual 10-Gb PCIe x8 Card	X540-T2	\$194.00	1	\$194.00
Intel Dual 10Gbps X540-BT2 I/O Module	AXX10GBTWLIOM	\$413.00	1	\$413.00
1TB Seagate Constellation SATA 7200rpm HDD 2.5"	ST91000640NS	\$70.00	2	\$140.00
Subtotal per node				\$15523.00

According to this configuration, we compute the different monthly costs for this compute node using the equations in Table 25. The results are highlighted in Table 103:

TABLE 103: COSTS ESTIMATED FOR THE COMPUTE NODE

Cost	Value
Cost per month for the hardware	\$472.24
Cost per month for the rack space	\$28.57
Cost per month for the leaf network ports	\$32.91

Cost per month for the management network ports	\$1.81
Cost per month for the power	\$63.92
Cost per month per compute node	\$535.53
Cost per vCPU per month	\$0.93
Cost per vCPU per hour	\$0.00127
Cost per GB of RAM per month	\$0.35
Cost per GB of RAM per hour	\$0.00048

A.1.5 Block storage HDD pricing

TABLE 104: BLOCK STORAGE HDD'S BILL OF MATERIALS

Part name	Part number	Cost per unit	Quantity	Subtotal
2U Dual Xeon Intel Server platform, 12 x 3.5" HDDs, Dual PSU	R2312WTTYSR	\$1855.00	1	\$1855.00
2nd 1000W PSU for R2312WTTYSR server	AXX1100PCRPS	\$233.00	1	\$233.00
Premium Rail Kit	AXXPRAIL	\$103.00	1	\$103.00
Intel Xeon E5-2620V4 2.10GHz 20MB LGA2011-3 8C/16T	BX80660E52620V4	\$419.00	2	\$838.00
Intel Certified 16GBDDR4 2133MHz ECC Reg DIMM	KVR21R15D4/16	\$157.00	4	\$628.00
LSI MegaRAID SAS 9361-8i SGL RAID Controller	9361-8i	\$752.00	1	\$752.00
LSI CacheVault Flash Module LSICVM02	LSICVM/02	\$159.00	1	\$159.00
Intel 12Gb/s RES3FV288 RAID Expander LP PCI Card	RES3FV288	\$346.00	1	\$346.00
Intel Ethernet X540-T2 Dual 10-Gb PCIe x8 Card	X540-T2	\$194.00	1	\$194.00
400GB Intel SSD DCP3700 series PCIe SATA HDD 2.5"	SSDPE2MD400G40 1	\$758.00	1	\$758.00
1TB Seagate Constellation SATA 7200rpm HDD 2.5"	ST91000640NS	\$70.00	2	\$140.00
4TB Seagate Constellation SATA 7200rpm HDD 3.5"	ST4000NM0053	\$215.00	12	\$2580.00
Subtotal per node				\$8586.00

A.1.6 Block storage SSD pricing

TABLE 105: BLOCK STORAGE SSD'S BILL OF MATERIALS

Part name	Part number	Cost per unit	Quantity	Subtotal
Supermicro 2028U-TN24R4T+ Barebone	SYS-2028U-TN24R4T+	\$3595.00	1	\$3595.00

Intel Xeon E5-2620V4 2.10GHz 20MB LGA2011-3 8C/16T	BX80660E52620V4	\$419.00	2	\$838.00
Micron 32GB DDR4 2400 MT/s	MEMCRU671834	\$367.00	4	\$1468.00
1TB Seagate Constellation SATA 7200rpm HDD 2.5"	ST91000640NS	\$70.00	2	\$140.00
Micron 9100 Pro U.2 Ssd 3.2Tb	MTFDHAL3T2MCE- 1AN1ZABYY	\$1499.00	12	\$17988.00
Subtotal per node				\$24029.00

A.1.7 Bill of materials for object storage low density

TABLE 106: BILL OF MATERIALS FOR OBJECT STORAGE LOW DENSITY

Part name	Part number	Cost per unit	Quantity	Subtotal
2U Dual Xeon Intel Server platform, 12 x 3.5" HDDs, Dual PSU	R2312WTTYSR	\$1855.00	1	\$1855.00
2nd 1000W PSU for R2312WTTYSR server	AXX1100PCRPS	\$233.00	1	\$233.00
Premium Rail Kit	AXXPRAIL	\$103.00	1	\$103.00
Intel Xeon E5-2620V4 2.10GHz 20MB LGA2011-3 8C/16T	BX80660E52620V4	\$419.00	2	\$838.00
Intel Certified 16GBDDR4 2133MHz ECC Reg DIMM	KVR21R15D4/16	\$157.00	4	\$628.00
LSI MegaRAID SAS 9361-8i SGL RAID Controller	9361-8i	\$752.00	1	\$752.00
LSI CacheVault Flash Module LSICVM02	LSICVM/02	\$159.00	1	\$159.00
Intel 12Gb/s RES3FV288 RAID Expander LP PCI Card	RES3FV288	\$346.00	1	\$346.00
Intel Ethernet X540-T2 Dual 10-Gb PCIe x8 Card	X540-T2	\$194.00	1	\$194.00
1TB Seagate Constellation SATA 7200rpm HDD 2.5"	ST91000640NS	\$70.00	2	\$140.00
8TB Seagate Enterprise Series Hard Drive SATA 7200rpm	ST8000NM0055	\$264.00	12	\$3168.00
Subtotal per node				\$8416.00

A.1.8 Bill of materials for object storage high density

TABLE 107: BILL OF MATERIALS FOR OBJECT STORAGE WITH HIGH DENSITY

Part name	Part number	Cost per unit	Quantity	Subtotal
Supermicro 6048R	6048R-E1CR60L	\$4755.00	1	\$4755.00
Intel Xeon E5-2620V4 2.10GHz 20MB LGA2011-3 8C/16T	BX80660E52620V 4	\$419.00	2	\$838.00
Intel Certified 16GBDDR4 2133MHz ECC Reg DIMM	KVR21R15D4/16	\$157.00	4	\$628.00
Supermicro AOC-STG- 12T Server NIC	AOC-STG-12T	\$156.00	1	\$156.00
300GB Seagate Savvio 10k rpm SAS Drives	ST300MM0026	\$65.00	2	\$130.00
8TB Seagate Enterprise Series Hard Drive SATA 7200rpm	ST8000NM0055	\$264.00	60	\$15840.00
Subtotal per node				\$22347.00

A.1.9 Bill of materials for object storage archiving

TABLE 108: BILL OF MATERIALS FOR OBJECT STORAGE ARCHIVING

Part name	Part number	Cost per unit	Quantity	Subtotal
Supermicro 6048R	6048R-E1CR60L	\$4 755.00	1	\$4755.00
Intel Xeon E5-2620V4 2.10GHz 20MB LGA2011-3 8C/16T	BX80660E52620V4	\$419.00	2	\$838.00
Intel Certified 16GBDDR4 2133MHz ECC Reg DIMM	KVR21R15D4/16	\$157.00	4	\$628.00
Supermicro AOC-STG- 12T Server NIC	AOC-STG-12T	\$156.00	1	\$156.00
300GB Seagate Savvio 10k rpm SAS Drives	ST300MM0026	\$65.00	2	\$130.00
8TB Seagate Enterprise Series Hard Drive SATA 7200rpm	ST8000NM0055	\$264.00	60	\$15840.00
Subtotal per node				\$22347.00

A.1.10 TCO pricing

In the following sub-sections, we present some examples of the private small, medium, and large Cloud and we estimate the TCO for three Clouds based on the previous examples of hardware nodes.

A.1.10.1 Small private cloud cost

TABLE 109: SMALL PRIVATE CLOUD (8 X 5) COST

	Quantity	HW cost per month	Rack units	RU cost	Power (in Watts)	Power cost	10Gbps ports	1Gbps ports	Cost per month	Cost in 3 years
Spine switches	0	\$0	0	\$0	0	\$0			\$0	\$0
Leaf switches	2	\$593	2	\$57	610	\$111			\$762	\$27428
Management switches	1	\$40	1	\$29	52	\$9	2		\$78	\$2824
Controller nodes	3	\$917	6	\$171	780	\$142	18	3	\$1231	\$44312
Network nodes	0	\$0	0	\$0	0	\$0	0	0	\$0	\$0
Compute nodes	3	\$1417	3	\$86	1050	\$192	12	3	\$1694	\$60991
Block storage nodes (HDD)	3	\$784	6	\$171	1596	\$291	6	3	\$1247	\$44874
Object storage nodes (Low density)	3	\$768	6	\$171	1245	\$227	6	3	\$1167	\$42008
Staff	1								\$15456	\$556400
VMWare Licences	6	\$0.00	0	\$0.00	0	\$0.00	0	0	\$200.00	\$7200.00
TCO			24		5333		44	12	\$21834	\$786037

Table 110 summarizes the capacity provided by the small private Cloud, when we use the precedent examples of Bill of materials. The cost breakdown per month and per three years is given in Table 111:

TABLE 110: THE CAPACITY PROVIDED BY SMALL PRIVATE CLOUD

Capacity provided	Values
10Gbps ports	96
1Gbps ports	48
vCPUs	576
RAM	1512
Block storage (in GB)	43200
Object storage (in GB)	64000

TABLE 111: COSTS PROPORTION BY REGARD TO THE TCO

Cost Breakdown (per 3 years)	In \$	In %	Formula (for \$)
Hardware	\$162689	20.70%	$\sum_{per\ HW} HW_{Cost} \times 36$
Rack units	\$24686	3.14%	$\sum_{per\ HW} RU_{Cost} \times 36$

Power	\$35062	4.46%	$\sum_{per\ HW} Power_{Cost} \times 36$
Staff	\$556400	70.79%	$Staff_{Cost} \times 36$

Figure 59 shows the portion of cost in percentage for each of the hardware, rack units, power, and staff. Almost three quarters of the TCO is consumed by the staff composed of one person FTE. The rest of the TCO is distributed among hardware with up to 21%, power consumption with 4.46%, and the 3% remaining is allocated for the rack units. We may observe that the best way to save or reduce the TCO is by reducing the teams size (FTEs). This can be done by using tools that orchestrate and manage automatically the infrastructure.

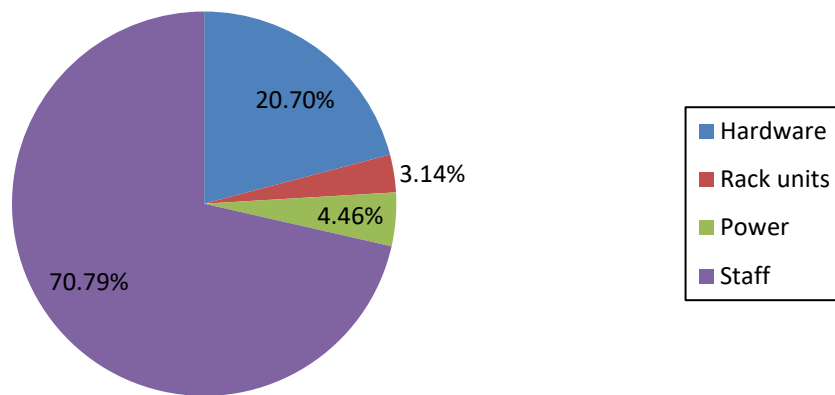


FIGURE 59: COST BREAKDOWN FOR A SMALL PRIVATE CLOUD (IN %)

In term of dollars, Figure 60 shows the costs breakdown for the TCO. The Staff spent about \$550000, while for the hardware is only about \$160000, the power consumption and the rack units consume together less than \$36000 which represent almost an tenth by report to the cost for staff. Therefore, it will be interesting for the infrastructure provider, when trying to save cost to reduce the staff by orchestrating and monitoring by tools the infrastructure.

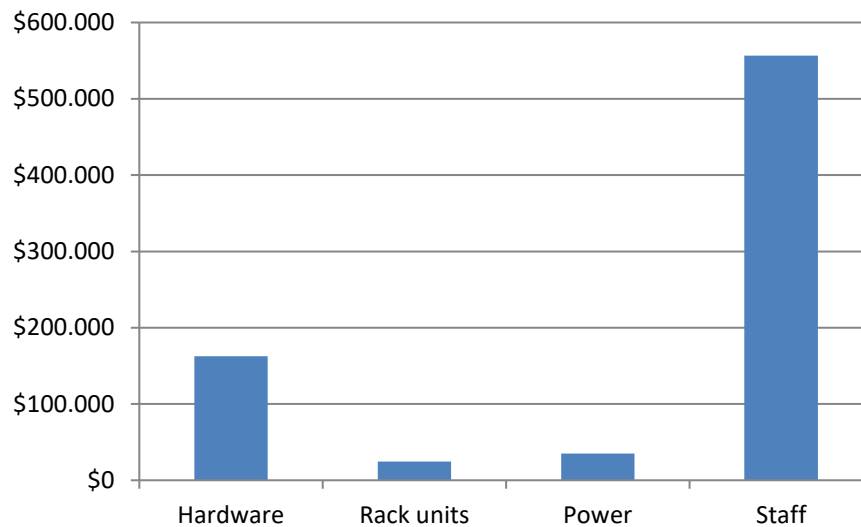


FIGURE 60: COST BREAKDOWN FOR SMALL PRIVATE CLOUD (IN \$)

A.1.10.2 Medium Private Cloud

We did the same exercise for the private medium Cloud. Results are obtained in Table 112 for the TCO and for the capacity provided by this medium Cloud.

TABLE 112: MEDIUM PRIVATE CLOUD (8 X 5) COST

	Quantity	HW cost per month	Rack units	RU cost	Power (in Watts)	Power cost	10Gbps ports	1Gbps ports	Cost per month	Cost in 3 years
Spine switches	0	\$0	0	\$0	0	\$0			\$0	\$0
Leaf switches	2	\$593	2	\$57	610	\$111			\$762	\$27428
Management switches	1	\$40	1	\$29	305	\$56	2		\$125	\$4487
Controller nodes	3	\$917	6	\$171	780	\$142	18	3	\$1231	\$44312
Network nodes	0	\$0	0	\$0	0	\$0	0	0	\$0	\$0
Compute nodes	13	\$6139	13	\$371	4550	\$831	52	13	\$7341	\$264294
Block storage nodes (HDD)	6	\$1567	12	\$343	3192	\$583	12	6	\$2493	\$89748
Object storage nodes (Low density)	6	\$1536	12	\$343	2490	\$455	12	6	\$2334	\$84016
VMWare Licences	26	\$0.00	0	\$0.00	0	\$0.00	0	0	\$866.67	\$31200.00
Staff	2								\$30911	\$1112800
TCO			46		11 927		96	28	\$46063	\$1658286

TABLE 113: THE CAPACITY PROVIDED BY MEDIUM PRIVATE CLOUD

Capacity provided	Values
40Gbps ports	0
10Gbps ports	96
1Gbps ports	48
vCPUs	3456
RAM	9072
Block storage (in GB)	108000
Object storage (in GB)	160000

Figure 61 depicts the breakdown for the costs in percentage for the medium private Cloud. As for Figure 61, the staff cost represents the most important expenditure. However, the hardware cost increases, which is obvious as the infrastructure is bigger.

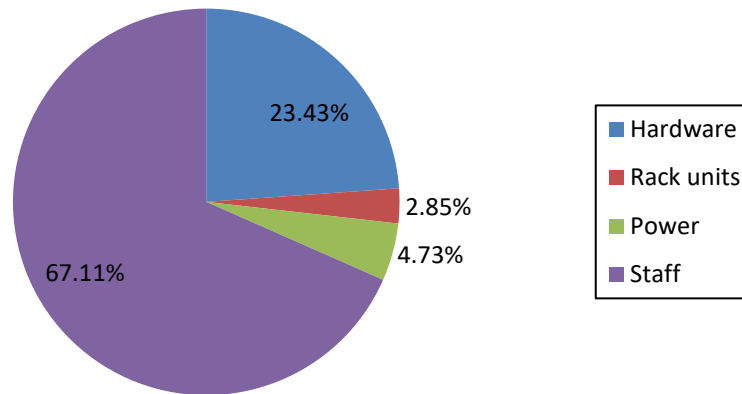


FIGURE 61: COST BREAKDOWN FOR MEDIUM PRIVATE CLOUD (IN %)

Figure 62 shows the costs breakdown for the TCO in dollars. The Staff spent about 1.6 billion \$, while for the hardware is about \$400 000, the power consumption and the rack units consume together less than \$200 000 which represent almost an eighth by report to the cost for staff. Therefore, it will be interesting for the infrastructure provider, when trying to save cost to reduce the staff by orchestrating and monitoring by tools the infrastructure.

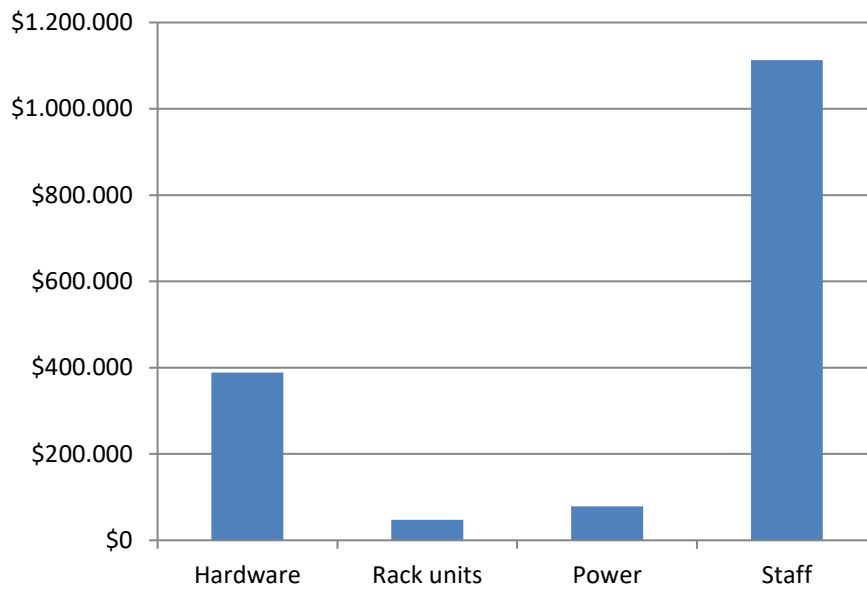


FIGURE 62: COST BREAKDOWN FOR MEDIUM PRIVATE CLOUD (IN \$)

A.1.10.3 Large Private Cloud (24 x 7)

We repeated the previous analysis for large private Cloud. The results are shown in Table 114 for the TCO and for the capacity provided.

TABLE 114: LARGE PRIVATE CLOUD (24 x 7) COST

	Quantity	HW cost per month	Rack units	RU cost	Power (in Watts)	Power cost	10Gbps ports	1Gbps ports	Cost per month	Cost in 3 years
Spine switches	2	\$675	2	\$57	600	\$110			\$842	\$30300
Leaf switches	31	\$997	31	\$886	9455	\$1727			\$11809	\$425137
Management switches	9	\$363	9	\$257	2745	\$501	2		\$1122	\$40384
Controller nodes	6	\$1834	12	\$343	1560	\$285	36	6	\$2462	\$88624
Network nodes	5	\$519	5	\$143	900	\$164	20	5	\$826	\$29744
Compute nodes	300	\$141672	300	\$8571	105000	\$19176	1200	300	\$169419	\$6099083
Block storage nodes (HDD)	10	\$7310	20	\$571	3850	\$703	40	10	\$8585	\$309046
Object storage	30	\$7836	60	\$1714	15960	\$2915	60	30	\$12465	\$448742

nodes (Low density)										
VMWare Licences	50	\$33992	100	\$2857	20750	\$3789	100	50	\$40639	\$1462988
Staff	12								\$18467	\$6676802
TCO			539		160820		1458	401	\$43635	\$1610850

TABLE 115: THE CAPACITY PROVIDED BY LARGE PRIVATE CLOUD

Capacity provided	Values
40Gbps ports	64
10Gbps ports	1488
1Gbps ports	432
vCPUs	86112
RAM	226044
Block storage (in GB)	233280
Object storage (in GB)	626400
Object storage (in GB)	7840000

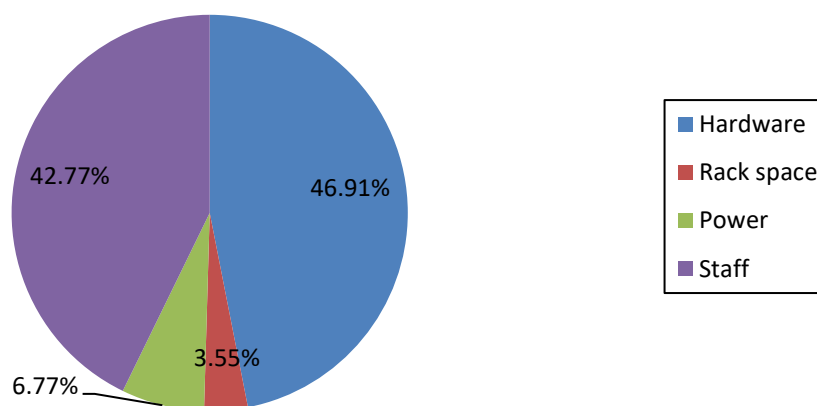


FIGURE 63: COST BREAKDOWN FOR LARGE PRIVATE CLOUD (IN %)

In Figure 63, we present the breakdown costs for the large private Cloud in percentage during a period of three years. In this figure we may notice that the hardware cost is the biggest one. This is quite intuitive as the resource for such large Cloud are supposed to be infinite (i.e., very large), hence the cost for such investment is important. Just after, we have the staff cost with almost 43% of the TCO. We believe that maintaining such big infrastructure needs considerable amount of FTE working people. For the power consumption, for our surprise, is still less than 7%, which is very small in comparison to the staff cost. Figure 64 presents the breakdown costs in \$, where more than \$7 300 000 are spent for the infrastructure.

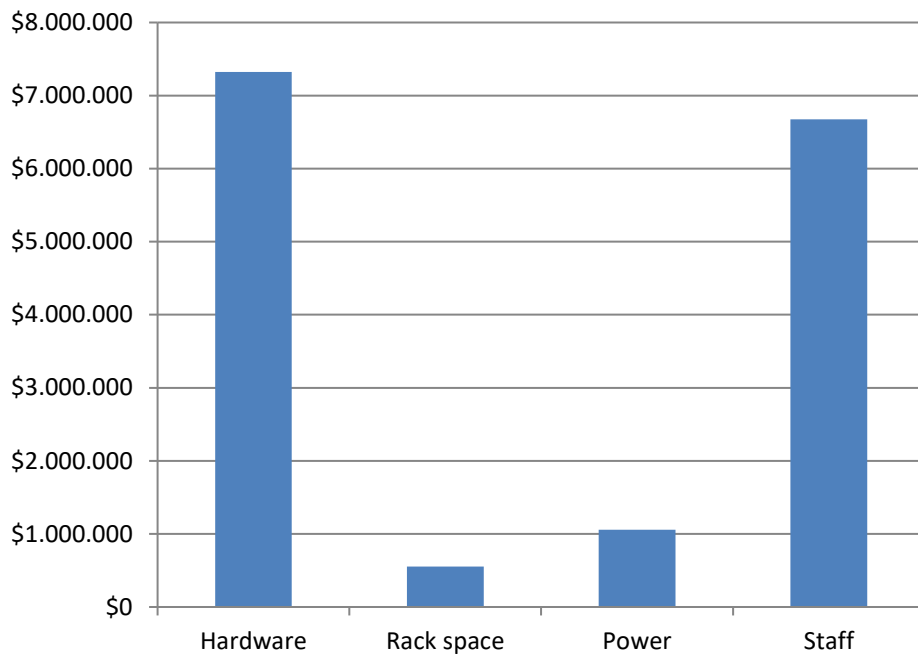


FIGURE 64: COST BREAKDOWN FOR LARGE PRIVATE CLOUD (IN \$)

A.1.10.4 Bottom Line

By this breakdown of costs by sector of expenditure, we are able to see the one that is the most important compared to the size of the data centre. Staff salary expenditures are the largest for small and medium data centres, while the cost of purchasing hardware is predominant for large DCs that provide more virtualized resources. The Figure 65 that makes a comparison of the cost weights shows staff and hardware expenditures have the most impact when varying the size of clouds.

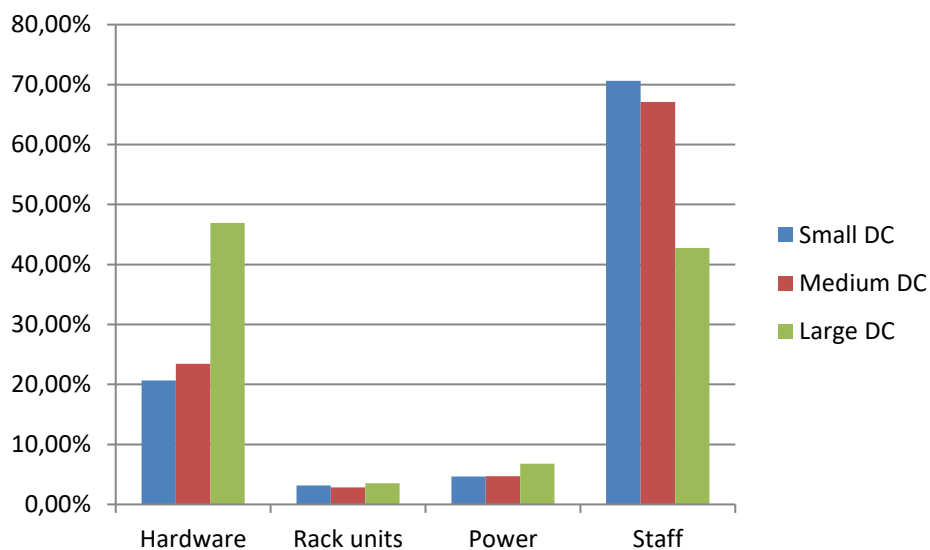


FIGURE 65 : COMPARISON OF THE COST BREAKDOWNS OF CLOUD SIZES

Another interesting analysis is to determine the variation of these costs over time. For this, we project spending to a date past (year 2013 for example) then we use the available statistical data to describe the evolution of the various costs until today (year 2019). Figure 66 and Figure 67 show the fluctuation for the period 2013-2019 of:

1. Hardware comprising Intel Xeon processor (cpubenchmark, s.d.), DIMM RAM, HDD and SSD disks (jcmitt, s.d.),
2. The Reference rent index (ANIL, s.d.)
3. The Electricity price (developpement-durable.gouv, s.d.)
4. and the minimum wage (INSEE, s.d.)

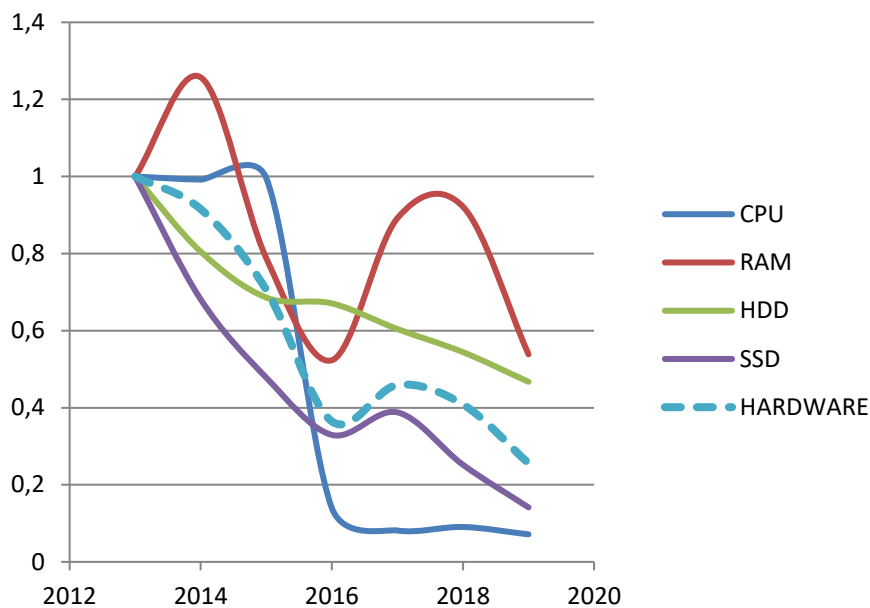


FIGURE 66: EVOLUTION OF THE HARDWARE COSTS FROM 2013 TO 2019

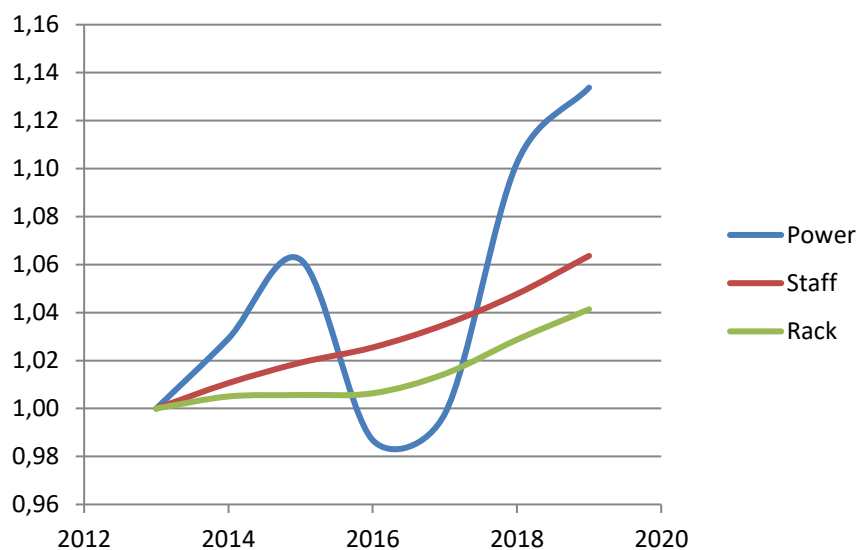


FIGURE 67: EVOLUTION OF THE HARDWARE COSTS FROM 2013 TO 2019

Thanks to the results obtained previously which make it possible to know the relative weight of each of these elements in the composition of a cloud, we easily project the evolution of its total cost by size of cloud as shown Figure 68, Figure 69 and Figure 70.

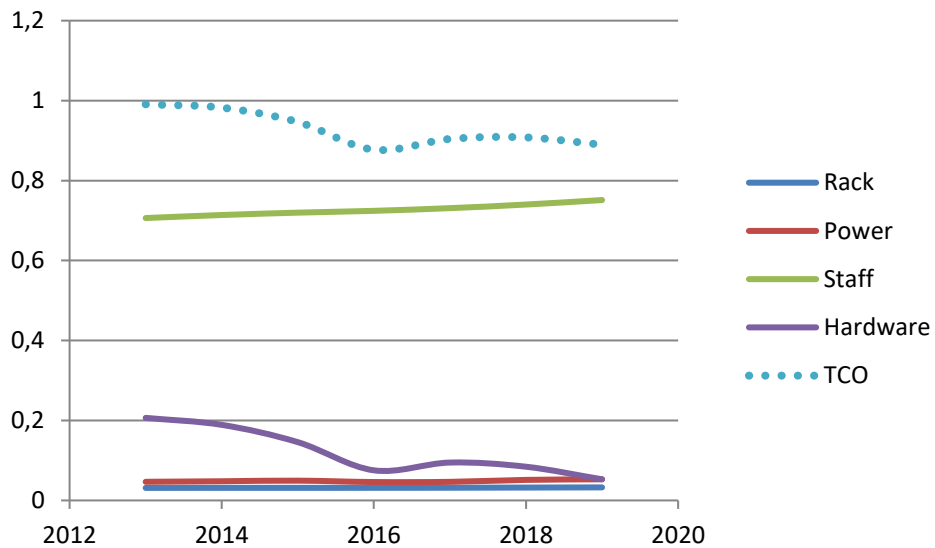


FIGURE 68: EVOLUTION THE TCO OF A SMALL DC FROM 2013 TO 2019

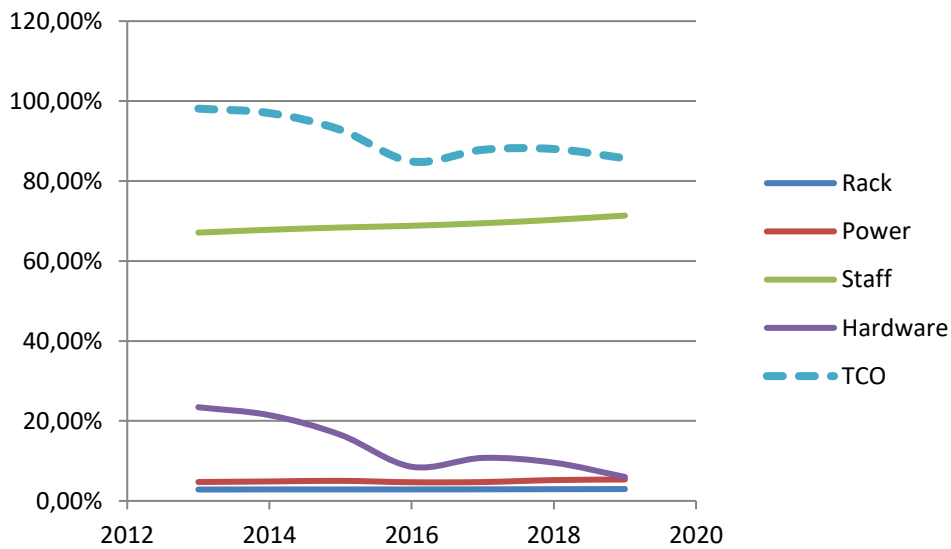


FIGURE 69: EVOLUTION THE TCO OF A MEDIUM DC FROM 2013 TO 2019

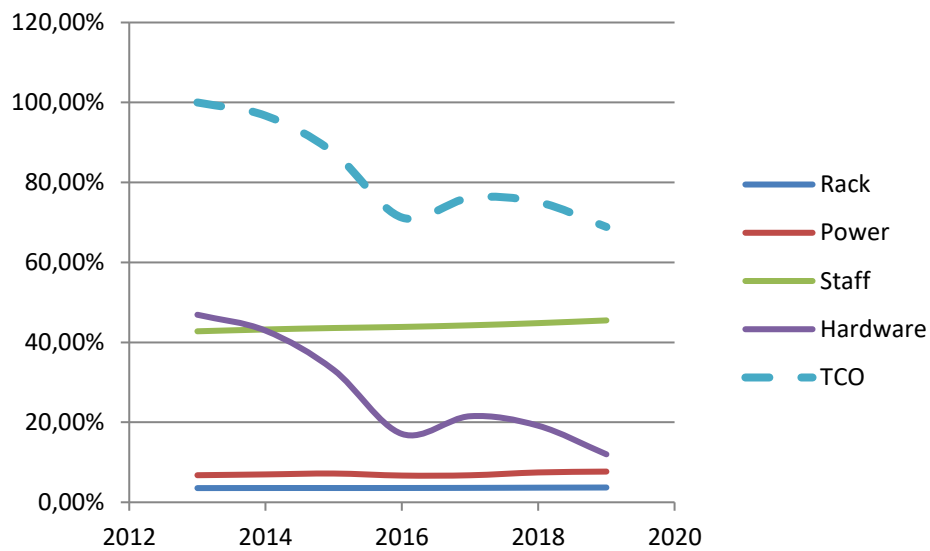


FIGURE 70: EVOLUTION THE TCO OF A LARGE DC FROM 2013 TO 2019

The larger the cloud, the faster the TCO decreases. This reduction is obtained thanks to the volume effect of the hardware whose price decreases more drastically than the other expenses increase.