



NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS
SCHOOL OF SCIENCE
DEPARTMENT OF PHYSICS

SECTION OF ELECTRONIC PHYSICS AND SYSTEMS

Distributed Resource Management in Converged Telecommunication Infrastructures

DOCTORAL DISSERTATION

by

Viktoria-Maria Alevizaki

Athens, June 2023.

This research is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme «Human Resources Development, Education and Lifelong Learning» in the context of the project “Strengthening Human Resources Research Potential via Doctorate Research – 2nd Cycle” (MIS-5000432), implemented by the State Scholarships Foundation (IKY).



Operational Programme
Human Resources Development,
Education and Lifelong Learning
Co-financed by Greece and the European Union





NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS
SCHOOL OF SCIENCE
DEPARTMENT OF PHYSICS
SECTION OF ELECTRONIC PHYSICS AND SYSTEMS

Distributed Resource Management in Converged Telecommunication Infrastructures

DOCTORAL DISSERTATION

by

Viktoria-Maria Alevizaki

Supervisory Committee: Anna Tzanakaki, Chair
Georgios Tombras, Committee Member
Dimitra Simeonidou, Committee Member

Doctoral Committee:

.....
Anna Tzanakaki
Associate Professor NKUA

.....
Georgios Tombras
Professor NKUA

.....
Dimitra Simeonidou
Professor University of Bristol

.....
Markos Anastasopoulos
Associate Professor NKUA

.....
Dionysios Reisis
Professor NKUA

.....
Spyros Denazis
Professor University of Patras

.....
Athanasios Korakis
Professor University of Thessaly

Athens, June 2023.

DECLARATION OF AUTHORSHIP

I, Viktoria-Maria ALEVIZAKI, declare that this thesis titled, “Distributed Resource Management in Converged Telecommunication Infrastructures” and the work presented in it are my own. I confirm that:

- This research was conducted during my candidature for a research degree at this University, either wholly or mainly.
- In cases where any section of this thesis has previously been submitted for a degree or qualification at this University or any other institution, it has been explicitly declared.
- Where I have consulted the published work of others, this is always clearly attributed.
- The source of any quoted material from the work of others is always provided. Aside from such quotes, this thesis represents solely my own original work.
- I have acknowledged all significant sources of assistance.
- If this thesis is based on collaborative work between myself and others, I have clearly delineated the contributions made by others and myself.

Signed:

Date:

ΠΕΡΙΛΗΨΗ

Η πέμπτη γενιά (5G) των ασύρματων και κινητών επικοινωνιών αναμένεται να έχει εκτεταμένο αντίκτυπο σε τομείς πέρα από αυτόν της τεχνολογίας πληροφοριών και επικοινωνιών (Information and Communications Technology - ICT). Το 5G ευθυγραμμίζεται με την 4η βιομηχανική εξέλιξη (4th industrial evolution), θολώνοντας τα όρια μεταξύ της φυσικής, της ψηφιακής και της βιολογικής σφαίρας. Σχεδιάστηκε για να προσφέρει δυνατότητες πολλαπλών υπηρεσιών και χρηστών, εκπληρώνοντας ταυτόχρονα πολλαπλές απαιτήσεις και επιχειρηματικά οικοσυστήματα. Ωστόσο, ορισμένες υπηρεσίες, όπως η επαυξημένη πραγματικότητα (Augmented Reality -AR), το εργοστάσιο του μέλλοντος (Factory of the Future) κ.λπ. θέτουν προκλήσεις για την ανάπτυξη μιας ενιαίας 5G υποδομής με βάση την ενεργειακή και οικονομική αποδοτικότητα. Σε αυτή τη κατεύθυνση, η παρούσα διδακτορική διατριβή υιοθετεί την ιδέα μιας καθολικής πλατφόρμας 5G που ενσωματώνει μια πληθώρα τεχνολογιών δικτύωσης (ασύρματες και ενσύρματες), και στοχεύει στην ανάπτυξη μαθηματικών εργαλείων, αλγορίθμων και πρωτοκόλλων για την ενεργειακή και λειτουργική βελτιστοποίηση αυτής της υποδομής και των υπηρεσιών που παρέχει. Αυτή η υποδομή διασυνδέει υπολογιστικούς, αποθηκευτικούς και δικτυακούς πόρους μέσω του προγραμματιζόμενου υλισμικού (hardware-HW) και της λογισμικοποίησης του δικτύου (network softwarisation). Με αυτό τον τρόπο, επιτρέπει την παροχή οποιασδήποτε υπηρεσίας με την ευέλικτη και αποτελεσματική μίξη και αντιστοίχιση πόρων δικτύου, υπολογισμού και αποθήκευσης.

Αρχικά, η μελέτη επικεντρώνεται στις προκλήσεις των δικτύων ραδιοπρόσβασης επόμενης γενιάς (NG-RAN), τα οποία αποτελούνται από πολλαπλές τεχνολογίες δικτύου για τη διασύνδεση ενός ευρέος φάσματος συσκευών με υπολογιστικούς και αποθηκευτικούς πόρους. Η ανάπτυξη μικρών κυψελών (small cells) είναι ζωτικής σημασίας για τη βελτίωση της φασματικής απόδοσης και της ρυθμαπόδοσης και μπορεί να επιτευχθεί είτε μέσω παραδοσιακών κατανεμημένων δικτύων ραδιοπρόσβασης (D-RAN) είτε μέσω δικτύων ραδιοπρόσβασης νέφους (C-RAN). Ενώ το C-RAN προσφέρει μεγάλα οφέλη όσο αφορά την επεξεργασία σήματος και τον συντονισμό σε σχέση με τα D-RAN, απαιτεί υψηλό εύρος ζώνης μετάδοσης και επιβάλλει σοβαρούς περιορισμούς καθυστέρησης στο δίκτυο μεταφοράς. Για την αντιμετώπιση αυτών των ζητημάτων, προτείνεται μια νέα αρχιτεκτονική «αποσύνθεσης των πόρων». Σύμφωνα με αυτήν, οι λειτουργίες βασικής επεξεργασίας σήματος (BBU functions) μπορούν να διαχωριστούν και να εκτελεστούν είτε στην ίδια θέση με τη κεραία (RU), είτε απομακρυσμένα σε κάποια μονάδα επεξεργασίας που βρίσκεται κοντά (DU) ή μακριά (CU) από την κεραία. Αυτή η έννοια της «αποσύνθεσης των πόρων» επιτρέπει την πρόσβαση σε κοινόχρηστους πόρους που παρέχονται από κέντρα δεδομένων μικρής ή μεγάλης κλίμακας, χωρίς να απαιτείται ιδιοκτησία των πόρων. Ωστόσο, η προσέγγιση αυτή απαιτεί την ανάπτυξη νέων πλαισίων βελτιστοποίησης για τη βελτίωση της αποδοτικότητας και της ευελιξίας των υποδομών 5G, ώστε να διαχειρίζονται αποτελεσματικά τους διαχωρισμένους πόρους. Καθοριστικό ρόλο σε αυτό αποτελεί η αρχιτεκτονική της Δικτύωσης Καθορισμένης από Λογισμικό (SDN), η οποία στοχεύει να επιτρέψει την προγραμματιζόμενη και δυναμική διαχείριση των πόρων του δικτύου μέσω κεντρικού ελέγχου. Έχοντας υπόψη τα παραπάνω, στο πρώτο μέρος της διατριβής αναπτύσσεται ένα πλαίσιο βελτιστοποίησης που προσδιορίζει το βέλτιστο λειτουργικό διαχωρισμό μεταξύ των λειτουργιών βασικής επεξεργασίας σήματος, σε συνδυασμό με τη βέλτιστη τοποθέτηση του SDN ελεγκτή, λαμβάνοντας επίσης υπόψη τη σταθερότητα του συνολικού συστήματος και τη μείωση των συνολικών λειτουργικών

δαπανών. Η ανάλυση επεκτείνεται περαιτέρω με προηγμένα σχήματα βελτιστοποίησης, με σκοπό την προσέγγιση ενός πιο ρεαλιστικού περιβάλλοντος 5G, όπου η ραγδαία αύξηση της κίνησης συνεπάγεται την ανάγκη για μεγαλύτερες δυνατότητες κλιμάκωσης για τη διαχείριση των χωρικών και χρονικών μεταβολών της, καθώς και τερματικών με διαφορετικές απαιτήσεις ποιότητας.

Στη συνέχεια μελετούνται τα δίκτυα πυρήνα του 5G. Στα δίκτυα πυρήνα 5G κάθε λειτουργία είναι λογισμικοποιημένη (softwarized) και απομονωμένη, επιτρέποντας την ανάπτυξη της σε υλικό γενικής χρήσης. Επίσης εισάγεται ένας νέος διαχωρισμός μεταξύ των λειτουργιών του επιπέδου ελέγχου και του επιπέδου δεδομένων (Control and User Plane Separation – CUPS) με βάση την SDN αρχιτεκτονική. Με τον τρόπο αυτό διαχωρίζεται η δικτυακή κίνηση μεταξύ των διαφορετικών 5G οντοτήτων (επίπεδο ελέγχου) και η δικτυακή κίνηση των χρηστών (επίπεδο χρήστη). Κρίσιμο ρόλο στο χειρισμό σημαντικού μέρους του επιπέδου χρήστη στα συστήματα 5G διαδραματίζει η οντότητα «λειτουργία επιπέδου χρήστη» (User Plane Function – UPF). Το UPF είναι υπεύθυνο για την προώθηση της πραγματικής κίνησης χρηστών με πολύ αυστηρές απαιτήσεις απόδοσης. Ανάλογα με τον τύπο της απαιτούμενης υπηρεσίας και την αρχιτεκτονική του δικτύου ραδιοπρόσβασης, οι κόμβοι UPF μπορούν να βρίσκονται είτε πιο κοντά είτε πιο μακριά από αυτό, ανακατευθύνοντας την κυκλοφορία σε διακομιστές κοντά στην άκρη του δικτύου για μείωση του χρόνου καθυστέρησης ή σε κεντρικές εγκαταστάσεις. Ως εκ τούτου, προκύπτει το ερώτημα της επιλογής των βέλτιστων στοιχείων UPF, καθώς η επιλογή ενός μη διαθέσιμου υπολογιστικού πόρου UPF μπορεί να οδηγήσει σε μπλοκάρισμα και καθυστερήσεις της υπηρεσίας. Για την αντιμετώπιση αυτού του ζητήματος, προτείνουμε ένα μοντέλο ειδικά σχεδιασμένο για δυναμική επιλογή βέλτιστων στοιχείων UPF με στόχο την ελαχιστοποίηση της συνολικής καθυστέρησης της υπηρεσίας. Αναπτύσσουμε συναρτήσεις κόστους για το μοντέλο χρησιμοποιώντας εργαστηριακές μετρήσεις που ελήφθησαν από μια πλατφόρμα 5G ανοιχτού κώδικα που φιλοξενείται σε περιβάλλον νέφους οπτικού κέντρου δεδομένων. Με το προτεινόμενο μοντέλο, μπορούμε να επιλέξουμε δυναμικά το καταλληλότερο στοιχείο UPF για τη χρήση υπολογιστικών πόρων, μειώνοντας τη καθυστέρηση εξυπηρέτησης.

Επεκτείνοντας την έννοια αποσύνθεσης των δικτυακών πόρων, η ανάλυση εστιάζει στα συστήματα 6G, τα οποία αναμένεται να υποστηρίξουν ένα ευρύ φάσμα υπηρεσιών μέσω μιας κοινής υποδομής που διευκολύνεται από τον τεμαχισμό δικτύου (network slicing). Τα συστήματα 6G προβλέπεται να λειτουργούν με αποκεντρωμένο τρόπο, που επιτρέπει στις εφαρμογές να παρεμβαίνουν άμεσα στις διαδικασίες ελέγχου για την πιο αποτελεσματική διασφάλιση της ποιότητας εμπειρίας (Quality of Experience – QoE) των τελικών χρηστών. Αυτό πραγματοποιείται μέσω της χρήσης της οντότητας «λειτουργία εφαρμογής» (Application Function – AF), η οποία διαχειρίζεται την εφαρμογή που εκτελείται στο τερματικό χρήστη (User Equipment – UE) και στο διακομιστή (Application Server – AS) που υποστηρίζει την υπηρεσία. Το AF διαδραματίζει κρίσιμο ρόλο στην παροχή υπηρεσιών υψηλού QoE, καθώς ενημερώνεται από την εφαρμογή και μπορεί να επηρεάσει τις αποφάσεις δρομολόγησης της κυκλοφορίας. Ωστόσο, η ανεξέλεγκτη λειτουργία του AF μπορεί να οδηγήσει σε αστάθεια στο σύστημα. Για την αντιμετώπιση αυτού του ζητήματος σχεδιάζουμε, εφαρμόζουμε και αξιολογούμε θεωρητικά και πειραματικά ένα πλήρως καταναμημένο πλαίσιο λήψης αποφάσεων για την εκχώρηση ροών (flow assignment) στα συστήματα 6G. Το πλαίσιο αυτό αποδεικνύεται ότι, υπό συγκεκριμένες συνθήκες, συγκλίνει σε ένα σταθερό σημείο που παρέχει τη βέλτιστη ισορροπία μεταξύ QoE και αποδοτικότητας κόστους. Οι συναρτήσεις κόστους που χρησιμοποιούνται ενσωματώνουν τόσο το κόστος δικτύου όσο και το υπολογιστικό κόστος, τα οποία προκύπτουν ρεαλιστικά μέσω μιας λεπτομερούς διαδικασίας που διεξάγεται σε μια λειτουργική 5G πλατφόρμα. Αυτή η διαδικασία επιτρέπει τη

μοντελοποίηση της απόδοσης του συστήματος και των απαιτήσεων σε διαφορετικά σενάρια λειτουργίας, τα οποία μπορούν να βοηθήσουν στη βελτιστοποιημένη διαχείριση του κύκλου ζωής των παρεχόμενων υπηρεσιών.

Τέλος, η μελέτη επικεντρώνεται στην πραγματική ανάπτυξη μιας υποδομής 5G που υποστηρίζει τον τεμαχισμό του δικτύου κατά παραγγελία από πολλαπλούς χρήστες. Ο τεμαχισμός του δικτύου επιτρέπει τον διαχωρισμό της φυσικής υποδομής δικτύου σε πολλαπλές λογικές υποδομές που μπορούν να υποστηρίξουν διαφορετικές κατηγορίες υπηρεσιών. Ένα τμήμα δικτύου (slice) έχει τους δικούς του αποκλειστικούς πόρους από το δίκτυο πρόσβασης, μεταφοράς, και πυρήνα, καθώς και στοιχεία από διάφορους τομείς κάτω από τους ίδιους ή διαφορετικούς διαχειριστές. Η κοινή χρήση της υποκείμενης φυσικής υποδομής από τα τμήματα δικτύου περιλαμβάνει την ανάπτυξη κατάλληλων διεπαφών που μπορούν να χρησιμοποιηθούν για την σύνδεση των διαφορετικών δικτυακών στοιχείων, καθώς και τη δημιουργία κατάλληλων περιγραφών (descriptors) για την εικονοποίηση των 5G λειτουργιών (Εικονικές Δικτυακές Λειτουργίες 5G - 5G Virtual Network Functions – VNFs). Η συλλογή και ο κατάλληλος συνδυασμός πολλαπλών VNF δίνει μια 5G υπηρεσία δικτύου (Network Service - NS) από άκρη σε άκρη (End to End - E2E). Μέσω μιας πλατφόρμας διαχείρισης και ενορχήστρωσης (Management and Orchestration Platform - MANO), μπορούμε να συνδυάσουμε αυτές τις υπηρεσίες δικτύου για να δημιουργήσουμε και να διαχειριστούμε ένα 5G τμήμα δικτύου. Για να επιτευχθεί αυτό, στη μελέτη αυτή χρησιμοποιείται ένας ενορχηστρωτής που ονομάζεται Open Source MANO (OSM), ο οποίος είναι συμβατός με το πρότυπο της Εικονικοποίησης Λειτουργιών Δικτύου (NFV). Αναπτύσσονται descriptors τόσο για τις λειτουργίες του επιπέδου ελέγχου του 5G, όσο και για το επίπεδο χρήστη. Συνδυάζοντας αυτούς τους descriptors, επιτυγχάνεται η δυναμική υλοποίηση πολλαπλών τμημάτων δικτύου πάνω σε μια 5G πλατφόρμα που υποστηρίζει πολλαπλούς χρήστες και φιλοξενείται σε μια υποδομή κέντρου δεδομένων. Χρησιμοποιώντας τα δημιουργημένα VNF, μπορούμε να εκτελέσουμε το δίκτυο πυρήνα με το πάτημα ενός κουμπιού και να παρέχουμε πολλαπλά τμήματα δικτύου με διαφορετικά χαρακτηριστικά.

KEYWORDS

Δίκτυα 5^{ης} γενιάς (5G), Δικτύωση Καθορισμένη από Λογισμικό (SDN), Εικονικοποίηση Λειτουργιών Δικτύου (NFV), Ενορχήστρωση Δικτύου, Τεμαχισμός Δικτύου, Επεξεργασία στην Άκρη του Δικτύου, Δίκτυο Μεταφοράς 5G, Κατανεμημένη Διαχείριση Πόρων, Δίκτυο Πυρήνα 5G, Λειτουργία Επιπέδου Χρήστη (UPF), Ποιότητας Υπηρεσίας (QoS).

ABSTRACT

The fifth generation (5G) of wireless and mobile communications is expected to have a far-reaching impact on society and businesses beyond the information and communications technology (ICT) sector. 5G is aligned with the 4th industrial evolution, blurring the lines between the physical, digital, and biological spheres. A common design is necessary to accommodate all service types based on energy and cost efficiency. To address this, this PhD thesis adopts the idea of a universal 5G platform that integrates a variety of networking technologies (wireless and wired), and aims to develop mathematical tools, algorithms and protocols for the energy and operational optimization of this infrastructure and the services it provides. This infrastructure interconnects computing, storage and network components that are placed at different locations, using the concepts of programmable hardware (hardware-HW) and network software (network softwarisation). In this way, it enables the provision of any service by flexibly and efficiently mixing and matching network, computing and storage resources.

The thesis targeted four distinct contributions. All proposed contributions are implemented and investigated experimentally in a 5G open-source lab testbed. The first contribution focused on optimal function and resource allocation adopting the innovative 5G RAN architecture, that splits flexibly the baseband processing function chain between Remote, Distributed and Central Units. This enables access to shared resources provided by micro or large-scale remote data centers, without requiring resource ownership. To support this architecture, networks adopt the Software Defined Networking (SDN) approach, where the control plane is decoupled from the data plane and the associated network devices and is centralized in a software-based controller. In this context, the goal of the proposed approach was to develop effective optimization techniques that identify the optimal functional split, along with the optimal location and size of the SDN controllers. The second contribution concentrated on solving the User Plane Function (UPF) selection problem in 5G core networks. According to the SDN paradigm 5G core control plane functions manage the network, while UPFs are responsible for handling users' data. Depending on the 5G RAN deployment option and the nature of the service, UPF nodes can be placed closer to the network edge, directing traffic to the Multi-access Edge Computing (MEC) servers hence reducing latency, or be placed deeper into the network directing traffic to central cloud facilities. In this context, a framework that selects the optimal UPF nodes to handle user's traffic minimizing total service delay has been proposed. The third contribution pertained to service provisioning in upcoming telecommunication systems. 6G systems require novel architectural Quality of Experience (QoE) models and resource allocation strategies that can differentiate between data streams originating from the same or multiple User Equipment (UEs), respond to changes in the underlying physical infrastructure, and scale with the number of connected devices. Currently, centralized management and network orchestration (MANO) platforms provide this functionality, but they suffer scalability issues. Therefore, future systems are anticipated to operate in a distributed manner, allowing applications to directly intervene in relevant control processes to ensure the required QoE. The proposed approach focused on developing a flow assignment model that supports applications running on UEs. The final contribution of this thesis focused on the deployment of a 5G infrastructure that supports multi-tenant network slicing on demand. Sharing of the underlying physical infrastructure was achieved through the development of suitable interfaces for integrating different network components and the creation of appropriate descriptors for virtual 5G network functions (VNFs). By collecting and combining multiple VNFs, an end-to-end 5G

Network Service (NS) can be obtained. Using a MANO platform, these NSs can be combined to instantiate and manage a 5G network slice.

KEYWORDS

5G, Software Defined Networking, Network Function Virtualization, Orchestration, Slicing, Edge Computing, 5G Transport Network, Distributed Resource Management, 5G Core, User Plane Function, QoS.

CONTRIBUTION STATEMENT

This thesis is based on the following papers:

- I. Alevizaki, Victoria-Maria & Anastasopoulos, Markos & Tzanakaki, Anna & Simeonidou, Dimitra, "Joint Fronthaul Optimization and SDN Controller Placement in Dynamic 5G Networks", in *proc. of Optical Network Design and Modelling 2019*. Available: https://doi.org/10.1007/978-3-030-38085-4_16.
- II. Alevizaki, VM., Anastasopoulos, M., Tzanakaki, A. et al., " Adaptive FH optimization in MEC-assisted 5G environments," *Photonics Network Communications* 40, 209–220, 2020. Available: <https://doi.org/10.1007/s11107-020-00906-8>
- III. V. M. Alevizaki, M. Anastasopoulos, A. Tzanakaki and D. Simeonidou, "Dynamic Selection of User Plane Function in 5G Environments," in *Proc. of Optical Network Design and Modelling*, 2021.
- IV. V. M. Alevizaki, A. I. Manolopoulos, M. Anastasopoulos and A. Tzanakaki, "Dynamic User Plane Function Allocation in 5G Networks enabled by Optical Network Nodes," in *2021 European Conference on Optical Communication (ECOC)*, pp. 1-4, 2021. Available: <https://doi.org/10.1109/ECOC52684.2021.9606154>.
- V. V. -M. Alevizaki, M. Anastasopoulos, A. -I. Manolopoulos and A. Tzanakaki, "Distributed Service Provisioning for Disaggregated 6G Network Infrastructures," in *IEEE Transactions on Network and Service Management*, Vol. 20, Issue 1, pp. 120-137, 2023. Available: <https://doi.org/10.1109/TNSM.2022.3211097>.

ACKNOWLEDGMENTS

First and foremost, I would like to express my sincere gratitude to my supervisor, Associate Professor Anna Tzanakaki for her guidance, motivation and support during my PhD journey.

I would also like to express appreciation to Associate Professor Markos Anastasopoulos for his continuous guidance and support to my Ph.D study.

Over the years I have had enjoyed working with many people in our research group. Some of these deserve a special mention. So, to Alexandros Manolopoulos and Petros Georgiadis thank you for the fruitful discussions, contributions to the co-authored papers, as well as the advice you have provided me with. I am thankful to all the members in our research group for providing a friendly work environment as well as their kind help during my PhD study. It was a pleasant journey studying with them.

I am also grateful to all my friends for their encouragement and support, which makes my life more wonderful.

Last but not least, I would like to express my heartfelt appreciation to my family for their everlasting love and unequivocal support during my whole life.

CONTENTS

DECLARATION OF AUTHORSHIP	II
ΠΕΡΙΛΗΨΗ	III
ABSTRACT	VI
CONTRIBUTION STATEMENT.....	VIII
ACKNOWLEDGMENTS.....	IX
CONTENTS	X
LIST OF FIGURES.....	XIII
LIST OF TABLES	XV
LIST OF ACRONYMS	XVI
CHAPTER 1 INTRODUCTION	1
1.1. MOTIVATION AND PROBLEM STATEMENT	1
1.2. THESIS FOCUS AND CONTRIBUTIONS	5
1.3. THESIS OUTLINE.....	6
REFERENCES.....	7
CHAPTER 2 5G SYSTEM ARCHITECTURE.....	9
2.1. CHAPTER INTRODUCTION.....	9
2.2. NETWORK ARCHITECTURE.....	11
2.2.1. <i>Overall Functional Architecture</i>	11
2.2.2. <i>NG-RAN</i>	16
2.2.3. <i>5G-Core</i>	18
2.2.4. <i>Interfaces</i>	20
2.3. NETWORK SLICING.....	22
2.4. QOS ARCHITECTURE	24
2.5. MANAGEMENT AND NETWORK ORCHESTRATION	26
2.6. SUMMARY	27
REFERENCES.....	27
CHAPTER 3 ADAPTIVE FRONTHAUL OPTIMIZATION IN MEC ASSISTED 5G ENVIRONMENTS	31
3.1. CHAPTER INTRODUCTION.....	31
3.2. THEORETICAL BACKGROUND.....	34
3.2.1. <i>Evolutionary Game Theory: Basic Concepts</i>	34
3.2.2. <i>Dynamics of Multi-Agent Learning</i>	35
3.2.3. <i>Evolutionary Dynamics of Reinforcement Learning</i>	35
3.3. APPLICATION TO WIRELESS NETWORKS IN 5G	36
3.3.1. <i>Problem formulation using EGT</i>	37
3.3.2. <i>Problem formulation using MARL</i>	40
3.4. RESULTS AND DISCUSSION	41
3.4.1. <i>SDN Controller placement using EGT</i>	41
3.4.2. <i>Optimal Split Selection using MARL</i>	45
3.5. SUMMARY	46
REFERENCES.....	47
CHAPTER 4 DYNAMIC USER PLANE FUNCTION ALLOCATION IN 5G NETWORKS	50

4.1.	CHAPTER INTRODUCTION.....	50
4.2.	PROBLEM FORMULATION BASED ON THEORETICAL EGT.....	52
4.2.1.	<i>System Model</i>	52
4.2.2.	<i>Application in 5G networks</i>	54
4.2.3.	<i>Numerical Results and Discussion</i>	56
4.3.	PROBLEM FORMULATION BASED ON LAB MEASUREMENTS.....	56
4.3.1.	<i>System Model</i>	56
4.3.2.	<i>Results and Discussion</i>	58
4.4.	SUMMARY.....	60
	REFERENCES.....	60
CHAPTER 5	DISTRIBUTED SERVICE PROVISIONING FOR 6G NETWORK INFRASTRUCTURES.....	62
5.1.	CHAPTER INTRODUCTION.....	63
5.2.	SERVICE SLICING IN 5G SYSTEMS.....	65
5.2.1.	<i>QoS Architecture</i>	65
5.2.2.	<i>5G Data Collection and Analytics</i>	68
5.2.3.	<i>Problem Statement</i>	69
5.3.	RELATED WORK.....	70
5.4.	THEORETICAL BACKGROUND: EVOLUTIONARY GAME THEORY.....	72
5.5.	PROBLEM FORMULATION.....	73
5.5.1.	<i>Game Formulation</i>	73
5.5.2.	<i>Payoff Function</i>	75
5.5.3.	<i>Dynamics of Adaptation Process</i>	76
5.5.4.	<i>Stability Analysis</i>	77
5.5.5.	<i>From infinite to finite population</i>	77
5.6.	5G SYSTEM PROFILING.....	79
5.6.1.	<i>5G Platform Description</i>	80
5.6.2.	<i>Evaluation Process</i>	81
5.6.3.	<i>Cost and Charging Functions</i>	83
5.7.	NUMERICAL RESULTS.....	84
5.7.1.	<i>Simulation with Experimental Values</i>	84
5.7.2.	<i>Extended simulation with multiple AFs</i>	88
5.8.	SUMMARY.....	89
	APPENDIX.....	91
	<i>A. Derivation of Fokker-Planck Equation</i>	91
	<i>B. Derivation of Replicator Equation</i>	92
	<i>C. Calculation of π_{HPxt} and π_{LPxt} derivatives</i>	92
	REFERENCES.....	93
CHAPTER 6	NETWORK SLICING AND ORCHESTRATION.....	96
6.1.	INTRODUCTION.....	96
6.2.	NFV-MANO FUNDAMENTALS.....	98
6.2.1.	<i>Network Function Virtualization</i>	98
6.2.2.	<i>Network slicing</i>	99
6.3.	PROBLEM STATEMENT.....	101
6.4.	EXPERIMENTAL SETUP.....	102
6.4.1.	<i>5G Platform Description</i>	102
6.4.2.	<i>Orchestration Platform Overview</i>	102
6.5.	IMPLEMENTATION.....	105
6.5.1.	<i>Creation of Network Descriptors</i>	105
6.5.2.	<i>Slice Deployment and Results</i>	108
6.6.	SUMMARY.....	109
	REFERENCES.....	109

CHAPTER 7 CONCLUSIONS AND FUTURE WORK 113

LIST OF FIGURES

FIGURE 1. 1: OVERALL 5G VISION [2].	2
FIGURE 2. 1: KPIS TARGETS FOR 5G AND BEYOND [16][19][20].	10
FIGURE 2. 2: E2E 5G SYSTEM ARCHITECTURE [16].	12
FIGURE 2. 3: CONVERGED TRANSPORT NETWORK. [16].	14
FIGURE 2. 4: FUNCTIONAL SPLIT BETWEEN NG-RAN AND 5G-CORE. 5G SUPPORTS TWO DIMENSIONAL SPLITS: A CP/UP (FUNCTIONAL-VERTICAL) SPLIT AND A CU/DU (GEOGRAPHICAL-HORIZONTAL) SPLIT. [15].	16
FIGURE 2. 5: HORIZONTAL SPLITS OPTIONS FOR THE 5G RAN. [34][37].	17
FIGURE 2. 6: 5G-CORE SERVICE-BASED AND POINT TO POINT ARCHITECTURE.	21
FIGURE 2. 7: 5G NETWORK SLICING OVERVIEW [16].	22
FIGURE 2. 8: EVOLUTION OF THE QoS ARCHITECTURE FROM 4G TO 5G.[15].	25
FIGURE 3. 1: NETWORK ARCHITECTURE. IN THE MEC, A DECISION ABOUT WHICH FUNCTIONS SHOULD BE PROCESSED LOCALLY IS MADE FOR EACH RU. THE REMAINING SET OF FUNCTIONS FOR EACH RU ARE TRANSFERRED THROUGH A COMMON NETWORK INFRASTRUCTURE WITH CENTRALIZED CONTROL TO A DC FOR FURTHER PROCESSING.....	37
FIGURE 3. 2: ASSUMED FH/BH TRANSPORT NETWORK FOR THE SYSTEM DESCRIBED IN FIGURE 3. 1 THE RED CIRCLE REPRESENTS THE POSITION OF THE SDN CONTROLLER, AFTER THE IMPLEMENTATION OF THE HEURISTIC ALGORITHM DESCRIBED IN SECTION IV.A.. THE RED SQUARE REPRESENTS THE OPTIMAL POSITION ESTIMATED ACCORDING TO THE AVERAGE PROPAGATION LATENCY-CASE DESCRIBED IN [11].	43
FIGURE 3. 3: EVOLUTION OF THE PROBABILITIES OF THE THREE SPLIT OPTIONS, WITH THE PARAMETERS DESCRIBED IN TABLE 3. 2, WHEN: (A) THE CONTROLLER IS PLACED IN THE PROPOSED LOCATION (RED CIRCLE IN FIG. 3. 2) BY THE HEURISTIC, (B) THE CONTROLLER IS PLACED IN THE PROPOSED LOCATION (RED SQUARE IN FIG. 3. 2) OF THE AVERAGE PROPAGATION LATENCY-CASE DESCRIBED IN [11].	44
FIGURE 3. 4: THE EVOLUTION OF THE PROBABILITIES OF THE THREE SPLIT OPTION FOR THREE RUs, WITH THE PARAMETERS DESCRIBED IN TABLE 3. 2 USING (A) THE MULTIPOPOPULATION REPLICATOR EQUATION MODEL, (B) THE CROSS LEARNING ALGORITHM, WITH 600 STAGES AND $\theta = 9$.	45
FIGURE 3. 5: THE TOTAL POWER CONSUMPTION OF THE FH SERVICES IN RELATION WITH THE STAGES OF THE CROSS LEARNING ALGORITHM.	46
FIGURE 4. 1: 5G SYSTEM ARCHITECTURE FOR ACCESS TO TWO (E.G. LOCAL AND CENTRAL) DATA NETWORKS [3].	51
FIGURE 4. 2: (A) TRAJECTORIES OF PROPORTIONS OF POPULATION AND (B) CONVERGENCE OF THE ALGORITHM TO THE EQUILIBRIUM (FOR $M1 = 130, M2 = 70, ab = 1, kmeckcc = 10$). IN THE EQUILIBRIUM 16% OF GROUP 1 UES AND 32% OF GROUP 2 UES ARE SERVED BY THEIR LOCAL UPFS, WHILE THE REMAINING ARE SERVED BY THE CENTRAL UPF.	55
FIGURE 4. 3: (A) EVOLUTION OF UES THAT CHOOSE MEC (B) RESOURCE UTILIZATION DURING THE EVOLUTION OF THE UE POPULATION.	59
FIGURE 4. 4: FIXATION PROBABILITIES FOR THREE CC WITH (A) DIFFERENT PROCESSING CAPABILITIES (B) DIFFERENT DISTANCE FROM THE UES.	60
FIGURE 5. 1: 5G SYSTEM ARCHITECTURE.	64
FIGURE 5. 2: 5G QoS SERVICE FLOWS AND PROTOCOL ADAPTATION.	66
FIGURE 5. 3: NWDAF DATA COLLECTION FROM AF VIA NEF.	69
FIGURE 5. 4: UE POPULATION STATE TRANSITION PROBABILITY.	75
FIGURE 5. 5: STRATEGY ADAPTATION PROCESS.	79
FIGURE 5. 6: A) 5G NETWORK DEPLOYMENT AND INFRASTRUCTURE CONNECTIVITY B) SERVERS/SWITCH USED IN THE EXPERIMENTATION.	80
FIGURE 5. 7: UE TRAFFIC AND CPU UTILIZATION FOR THE VM HOSTING A UPF NF	81
FIGURE 5. 8: CPU UTILIZATION FOR 4 VMs WITH DIFFERENT CHARACTERISTICS IN TERMS OF CPU, MEMORY AND DISK.	82
FIGURE 5. 9: AVERAGE THROUGHPUT OF THE HIGH AND LOW PRIORITY QUEUES.	83

FIGURE 5. 10: (A) EVOLUTION OF THE NUMBER OF UEs INSIDE THE POPULATION (B) EVOLUTION OF THE PAYOFF OF THE TWO STRATEGIES IN RELATION WITH THE INITIAL STATE.	85
FIGURE 5. 11: IMPACT OF THE PROPOSED SCHEME ON THE COMPUTE (MEC AND CC) AND NETWORK RESOURCES.	86
FIGURE 5. 12: DEPENDENCE OF THE INTERIOR EQUILIBRIUM FROM THE PROCESSING CAPABILITIES OF THE CC.....	87
FIGURE 5. 13: PERFORMANCE OF THE PROPOSED SCHEME VERSUS CENTRAL BASELINE ALGORITHM.	88
FIGURE 5. 14: PERFORMANCE OF THE PROPOSED SCHEME IN DIFFERENT TRAFFIC SCENARIOS.	89
FIGURE 5. 15: (A) EVOLUTION OF THE UEs POPULATION UNDER THE CONTROL OF ONE OR TWO AFs. (B) EVOLUTION OF EACH AF'S POPULATION OF UEs THAT ARE SERVED AT MEC FACILITY. IMPACT OF THE PROPOSED SCHEME ON (C) THE COMPUTE (MEC AND CC) AND (D) NETWORK RESOURCES.	90
FIGURE 6. 1: NFV ARCHITECTURAL FRAMEWORK.	99
FIGURE 6. 2: NETWORK TOPOLOGY UNDER CONSIDERATION.	101
FIGURE 6. 3: ENVIRONMENT DESCRIPTION.....	102
FIGURE 6. 4: NETWORK SERVICE INSTANTIATION AND CONFIGURATION WITH PROXY CHARMS[31].....	104
FIGURE 6. 5: NETWORK SERVICE INSTANTIATION THROUGH OSM[31]	105
FIGURE 6. 6: GRAPHICAL ILLUSTRATION OF THE VNFDs FOR 5G CP AND UP RESPECTIVELY.....	106
FIGURE 6. 7: GRAPHICAL ILLUSTRATION OF THE NSDs FOR 5G CP AND UP RESPECTIVELY.....	106
FIGURE 6. 8: GRAPHICAL ILLUSTRATION OF THE NSSIs OF EACH NSI	107
FIGURE 6. 9: THE TWO SLICES IMPLEMENTED ON THE LAB	108
FIGURE 6. 10: WIRESHAK TRACES AFTER THE INSTANTIATION OF THE TWO SLICES	108
FIGURE 6. 11: VISUALIZATION PLATFORM SCREENSHOTS FOR UE TRAFFIC OF EACH SLICE	109

LIST OF TABLES

TABLE 2. 1: KPIs AND TARGETS FOR 5G AND BEYOND TERMS [16][19][20].....	11
TABLE 2. 2: INDICATIVE, FUTURE SLA STRUCTURE [15][47].....	23
TABLE 3. 1: NETWORK AND PROCESSING DEMANDS OF EACH FUNCTIONAL SPLIT.....	38
TABLE 3. 2: PARAMETERS OF THE SYSTEM CONFIGURATION	42
TABLE 4. 1: CPU UTILIZATION FOR 4 VMs WITH DIFFERENT CHARACTERISTICS IN TERMS OF CPU, MEMORY AND DISK	58
TABLE 5. 2: STANDARDIZED 5QI TO QOS CHARACTERISTICS MAPPING [2]	67
TABLE 5. 3: VM CONFIGURATIONS FOR THE VM HOSTING UPF.....	81
TABLE 5. 4: CPU COST FUNCTION APPROXIMATION RESULTS OF FIGURE 5. 8.....	83
TABLE 5. 5: NETWORK COST FUNCTION APPROXIMATION RESULTS OF FIGURE 5. 9.....	84

LIST OF ACRONYMS

ACRONYM	DESCRIPTION
3GPP	3rd Generation Partnership Project
4G	Fourth Generation Telecommunication systems
5G	Fifth Generation Telecommunication systems
5GMM	5G Mobility Management
5GN	5G Networks
5G-NR/NG-RAN	5G New Radio/Next Generation RAN
5G-PPP	European Union Public-Private Partnership for 5G
5G-PPP	European Union Public-Private Partnership
5GS	5G System
5QI	5G QoS Identifier
6G	Sixth Generation Telecommunication systems
AF	Application Function
AI	Artificial Intelligence
AMF	Access and Mobility Function
AN	Access Network
API	Application Programming Interface
AQM	Active Queue Management
AR/VR	Augmented/Virtual Reality
ARP	Allocation and Retention Policy
AS	Application Server
AUSF	Authentication Server Function
B5G	Beyond 5G
BBU	Baseband Unit
BH	Backhaul
BS	Base Station
BTN	Backhaul Transport Network
CAPEX/OPEX	Capital/Operational Expenditures
CC	Central Cloud
CHF	Core Charging Function
CL-AS	closely located AS
CN	Core Network
CP	Control Plane
C-RAN	Cloud RAN
CU	Central Unit
CUPS	Control and User Plane Separation
DC	Data Center
DL	Downlink
DN	Data Network
D-RAN	Distributed RAN
DRB	Data Radio Bearer
DSCP	Differentiated Services Code Point
DU	Distributed Unit
E2E	End to end
eCPRI	enhanced Common Public Radio Interface
EGT	Evolutionary Game Theory
EH	Extended Header
eMBB	enhanced Mobile Broadband

EPC	Evolved Packet Core
EPS	Evolved Packet System
ERAB	EPS Radio Access Bearer
ESS	Evolutionary Stable Strategies
FEC	Forward Error Correction
FFS	Flexible Functional Splits
FFT	Fast Fourier Transform
FH	Fronthaul
FTN	Fronthaul Transport Network
GMB	guaranteed minimum bandwidth
gNB	Next Generation Node B
GTP-U	General Packet Radio Service (GPRS) Tunnelling Protocol for the user plane
HARQ	Hybrid Automatic Repeat Request
HNF	Hybrid Network Service
HTTP	Hypertext Transfer Protocol
HW	Hardware
ICT	Information and Communications Technology
IDFT	Inverse Discrete Fourier Transform
iFFT	inverse Fast Fourier Transform
IM	Information Model
IoT	Internet of Things
IP	Internet Protocol
IQ	Quadrature Signals
I-UPF	Intermediate UPF
KPI	Key Performance Indicator
LCM	LifeCycle Management
LEVO	Long-term Evolution Period
LOS/NLOS	Line of Sight/ Non-Line of Sight
LP/HP	Low/High priority
LTE	Long Term Evolution
MAC	Media Access Control
MANO	Management and Orchestration
MARL	Multi-Agent Reinforcement Learning
MAS	Multi Agent System
MEC	Multi-access Edge Computing
MEVO	Medium-term Evolution Period
MH	Midhaul
MIMO	Multiple Input Multiple Output
ML	Machine Learning
mMTC	massive Machine Type Communications
MNO	Mobile Network Operator
MOP	Multi-objective Optimization
NaaS	Network as a Service
NAS	Non-Access Stratum
NBI	Northbound Interface
NE	Nash Equilibrium
NEF	Network Exposure Function
NF	Network Function
NFV	Network Function Virtualization
NFVI	NFV Infrastructure

NGAP	New Generation Application Protocol
NGMN	Next Generation Mobile Network
NG-RAN	New Generation RAN
NIC	Network Interface Card
NRF	Network Repository Function
NS	Network Service
NSD	NS Descriptor
NSI	the 5G Network Slice Instance
NSSF	Network Slice Selection Function
NSSI	Network Slice Subnet Instance
NST	Network Slice Template
NWDAF	Network Data Analytics Function
ODL	OpenDayLight SDN Controller
OSM	Open Source MANO
PCF	Policy Control Function
PDCCP	Packet Data Convergence Protocol
PDF	Probability Density Function
PDU	Protocol Data Unit
PFCP	General Packet Radio Service
PHY	Physical
PNF	Physical Network Service
PON	Passive Optical Network
PRACH	Physical Random-Access Channel
PSA	PDU Session Anchor
QAM	Quadrature amplitude modulation
QFI	QoS Flow ID
QoE	Quality of Experience
QoS	Quality of Service
RAN	Radio Access Network
RAT	Radio Access Technology
RE	Replicator Equation
ReL	Reinforcement Learning
REST	Representational State Transfer
RF	Radio Frequency
RFSP	Rat Frequency Selection Priority
RL	5G Resource Layer
RL-AS	remotely located AS
RLC	Radio Link Control
RMSE	Root-Mean-Square Error
RoE	Radio over Ethernet
RRM	Radio Resource Management
RTT	Round Trip Time
RU	Remote Unit
SBA	Service-Based Architecture
SBI	Service-Based Interface
SD	Service Descriptor
SDAP	Service Data Adaptation Protocol
SDF	Service Data Flow
SDN	Software Defined Networking
SDO	Standard Developing Organization
SEVO	Short-term Evolution Period

SI	Service Instance
SIL	5G Service Instance Layer
SLA	Service Level Agreement
SMF	Session Management Function
SST	Slice Service Type
SW	Software
TDM	Time Division Multiplexing
TN	Transport Network
TSON	Time Shared Optical Network
UC	Use Case
UDM	Unified Data Management Function
UDR	Unified Data Repository
UE	User Equipment
UL	Uplink
UP	User Plane
UPF	User Plane Function
URLLC	Ultra-Reliable Low Latency Communication
VCA	VNF Configuration and Abstraction
VIM	Virtualized Infrastructure Manager
VM	Virtual Machine
VNF	Virtual Network Function
VNFD	VNF Descriptor
WAN	Wide Area Network
WDM	Wavelength Division Multiplexing

Chapter 1

Introduction

Contents

1.1.	Motivation and Problem Statement	1
1.2.	Thesis Focus and Contributions	5
1.3.	Thesis Outline.....	6
	References	7

1.1. Motivation and Problem Statement

It is anticipated that the fifth generation (5G) of wireless and mobile communications will have a significant impact on society and business that will go far beyond the information and communications technology (ICT) sector. 5G will offer much higher average data rates than prior cellular generations, enabling improved mobile broadband services. Although mobile broadband services are already widely used, 5G is anticipated to enable the next level of human and machine connectivity and interaction, for example through the widespread use of virtual or augmented reality, free-viewpoint video, and telepresence [1].

5G is in line with the 4th industrial evolution [2], that is a fusion of technologies, blurring the lines between the physical, digital, and biological spheres. Ubiquitous mobile coverage will facilitate the advancement of a variety of business sectors outside the ICT domain, that are referred to as vertical sectors [3]. 5G characteristics will have a great impact on the automotive field and transportation in general, by offering advanced forms of collaborative driving, protection mechanism for road users, increased efficiency in railroad transportation and many other developments. Furthermore, mining and construction sectors can leverage the automation that 5G provides for the remote control of vehicles or machines in dangerous or inaccessible areas. 5G can revolutionize health care through the possibility of wirelessly enabled smart pharmaceuticals or remote surgery with haptic feedback. Even the agricultural sector can benefit from 5G, by employing big data analytics and the Internet of Things (IoT) for tracking, monitoring, automating and analyzing its operations, leading to improvements on products' quality and quantity while reducing the overall waste and production cost. In general, the whole quality of living can be optimized through the so-

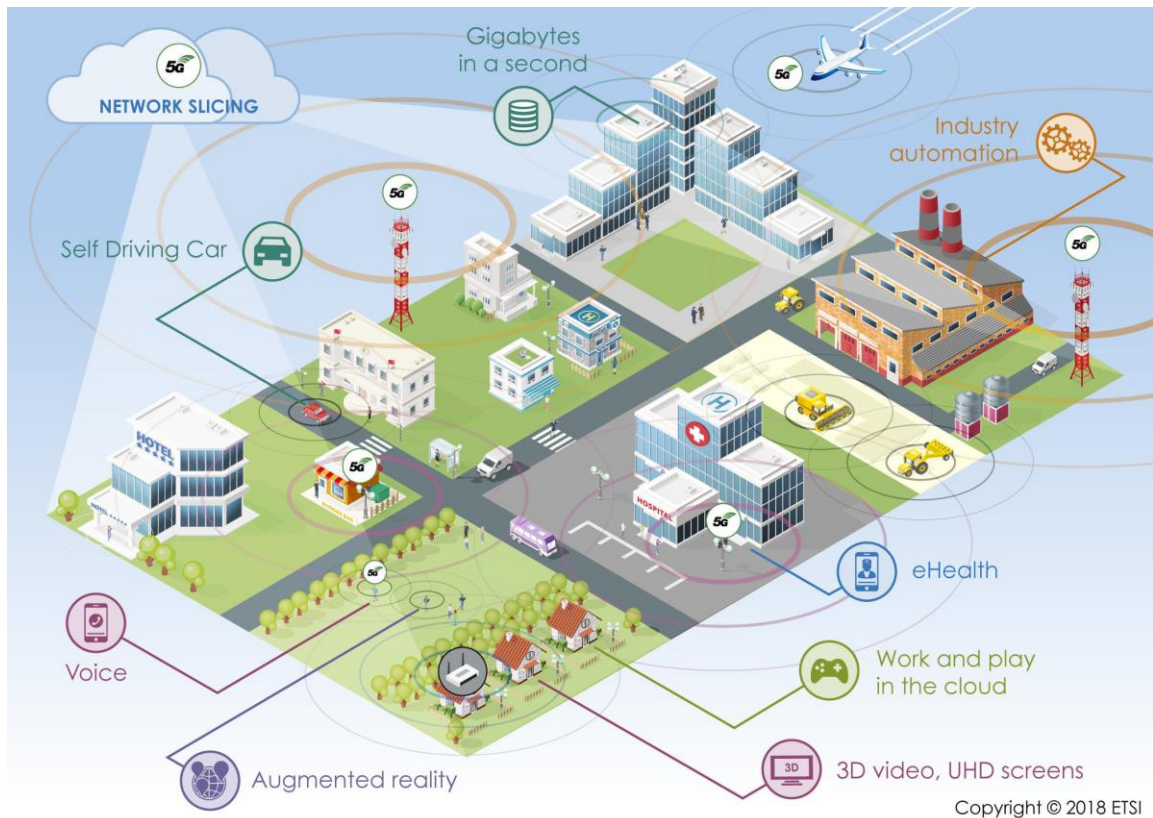


Figure 1. 1: Overall 5G Vision [4]

called 5G smart-cities that give emphasis to better energy, transportation, environment and waste management. **Figure 1. 1** shows the overall envisaged 5G system [4].

A distinguishing feature of 5G, compared to past generations of wireless communications, is the great diversity of technology drivers and use cases that are highlighted [5]. More specifically, prior generations have always been defined by a single monolithic system architecture that most of the times have been targeting a single requirement and a single business ecosystem, such as mobile broadband in the case of Long-Term Evolution (LTE). However, addressing this limitation, 5G was designed to offer multi-service and multi-tenancy capabilities. Taking that into consideration, 5G systems are not just the next evolution phase of 4G networks with new spectrum bands, higher spectral efficiencies, and higher peak throughput, but target to support a large variety of new services and business models. In this direction, the following main categories of 5G services have been defined [1] [6]:

- Enhanced mobile broadband (eMBB) is a term used to describe enhanced access to data, services, and multi-media content with better user experience and performance. This service type, is human-centric and can be thought of as an evolution of the services currently offered by 4G networks. The main goal is to offer seamless radio coverage practically anywhere and anytime with noticeably improved user data rates. It comprises of Use Cases (UCs) with a wide range of requirements, such as hotspot UCs with a high user density, very high traffic capacity, and low user mobility, to wide area coverage cases with medium to high user mobility.

- UCs with strict criteria for features like mobility (up to 100 km/h), latency (as low as 1 ms), reliability and availability (higher than 99.999% or "five nines"), belong to the category of 5G Ultra-Reliable and Low-Latency Communications (URLLC). These services are predicted to play a significant role to sectors beyond the ICT. Wireless management of production and manufacturing processes in the industrial sector (Industry 4.0), remote medical operation, automation in a smart grid, safety in transportation, are typical examples of such services.
- Massive machine-type communications (mMTC) are services that are characterized by the transmission of non-delay-sensitive data by a very large number of connected devices, usually at a very low rates. The main issue with these services is that devices must typically be inexpensive and have a very extended battery life. Agricultural applications, smart metering, and logistics applications (such as those that involve the tracking of tagged objects) are some prominent examples of this service type. In agricultural applications, for example, small, inexpensive, and low-power sensors are dispersed over large areas to measure ground humidity, fertility, etc.

However, concerns have emerged on the limitations of present vertical services and the requirements of future ones, in order to achieve total automation and digitization of vertical industries. Many 5G-envisioned services cannot readily be mapped to one of the three basic service types outlined above, since they relate to more than one service types. A typical example is augmented reality, a service that overlays information on the real world for the purposes of education/safety/training/gaming, which demands high values of throughput while maintaining low values for latency (10s and 100s of Gbps per communication link). Furthermore, the new use cases emerging from the Factory of the Future concept [7], where items in a factory environment will be connected wirelessly, cannot identify as one of the aforementioned service types, because they expect that both energy consumption and latency levels (as low as 250 μ s) are maintained at a minimum level.

These types of services pose the main challenges for the deployment of a unified 5G system. Developing a 5G network by taking into account each service type, or even individual UCs, separately, would lead to very different 5G system designs and architectures. However, based on energy and cost efficiency considerations, only a common design that accepts all service types is seen as a viable solution [8]. In this context, the concept of network slicing is in the center of attention in 5G systems. It basically divides the physical network infrastructure into multiple independent logical ones, each designed to support a certain class of service, by interconnecting a subset of the radio access, transport, and core network parts. To do this, 5G utilizes the concepts of Network Function Virtualization (NFV) [10] and Software Defined Networking (SDN) [11]. NFV enables the migration from the notion of network elements to Network Functions (NFs). Multiple network operations can be hosted as virtual machines (VMs) on general purpose servers, with the benefit of faster service delivery and cost reduction. This facilitates independence between technologies, offering increased granularity in how resources are committed and provided, without the requirement of owning and installing specific hardware (HW) or software (SW) [9]. SDN, on the other hand, refers to a network architectural approach, where the control and forwarding layer functions of the network are decoupled. In this architecture the network is controlled through a logically centralized entity, the controller, which controls a number of forwarding devices. Networks become programmable and manageable, while resource allocation,

scalability in distributed data centers, and device virtualization that is necessary for efficient resource sharing are facilitated. However, due to the communication requirements between the controller and the physical devices, great attention must be paid at the location and size of the controller, in order to minimize the total end-to-end (E2E) latency.

The concepts of NFV and SDN are adopted by the 5G-core architecture. Each 5G core function is softwarized and isolated, thus being able to be deployed over general purpose hardware. At the same time these functions are separated to Control and User Plane functions (CUPS), according to the SDN paradigm. The 5G-core Control Plane (CP) functions are accountable for decision-making and managing the network. The 5G-core User Plane (UP) consists of the User Plane Function (UPF) entity, which is responsible for the actual transmission of the data and the connectivity of the 5G network with the external IP network. The number and location of the UPFs may vary and depends on the 5G system implementation.

5G revolutionizes the way radio access networks (RAN) operate. 5G RAN (5G New Radio-5G NR or Next Generation RAN- NG-RAN) integrates multiple network technologies (wired and wireless), in order to interconnect a wide variety of end devices with computing and storage resources. The wireless access solutions that 5G networks will offer are expected to support a heterogeneous set of integrated interfaces and coexist with 2-4G technologies. For improved spectral efficiency and throughput, cells with limited coverage, also known as small cells, can be deployed either by adopting the traditional model of Distributed Radio Access Networks (D-RAN), where digital processing units, referred to as Baseband Units (BBUs), and antennas are co-located, or the more recently proposed concept of Cloud Radio Access Networks (C-RAN). In C-RAN, Remote Units (RUs) are connected to the Central Unit (CU), where the centralized BBU is located, through high-bandwidth links known as fronthaul (FH)[12]. The centralization of processing and coordination offered by C-RAN makes it well suited to address the increased operational and investment costs, as well as the limited scalability and flexibility of D-RAN. However, C-RAN requires enormous transmission bandwidth and imposes severe transport network delays constraints [12][13]. To address these challenges, the introduction of processing equipment in Distributed Units (DUs), that are closer to the edge of the access network seems to be a promising solution.

To this end, the concept of "disaggregation of resources", that is the idea of disconnecting network components and placing them at district and in some cases remote locations, is expected to play a key role. A new architecture is proposed that splits the baseband functional processing chain between the CU and the DUs, which are located closer to the RUs. FH connects the RUs with the DUs, while the links that connect the DUs with the CUs are referred as midhaul (MH). The choice of optimal split can be made based on factors such as transmission network and service characteristics, with significant benefits in terms of resources and energy efficiency [14]. Flexible splits can be provided by programmable digital HW used to support BBUs. The shared resources offered by this architectural model are provided by either micro data centers (Multi-access Edge Computing-MEC) or large-scale remote data centers (DCs), without the need for resource ownership. This alternative approach to access networks introduces the need to develop new optimization frameworks for the design and operation of 5G infrastructures with improved performance in efficiency, flexibility and density.

To ensure that the features of the proposed architecture are properly exploited and offer the requested scalability and efficiency, it is important to develop tools that will enable

the optimization of infrastructure design and operation. In this direction, in this thesis different parts of the 5G infrastructure will be examined and optimized through the adoption of new, scalable, stochastic, multi-objective optimization (MOP) algorithms that are based on different mathematical approaches. The overall architecture of the proposed infrastructure and its technical objectives are aligned with the performance parameters as described by the European Union Public-Private Partnership (5G PPP)[15].

1.2. Thesis Focus and Contributions

This PhD thesis, adopts the concept of a new generation universal 5G platform which integrates a multitude of networking technologies (wireless and wired) that are jointly optimized. This infrastructure will interconnect "disaggregated" compute/storage and network components and adopt the concepts of programmable HW and network softwarisation. It will enable the provision of any service by flexibly and efficiently mixing and matching network, computing and storage resources. The aim of this thesis is the development of mathematical tools, algorithms and protocols for the optimization of the infrastructure and the services it provides. To this end, a thorough study and application of important concepts of 5G networks, such as NFV, SDN and installation, operation and maintenance of a network cloud was performed. In the following we discuss the overall contributions of this thesis.

The first contribution of this thesis focuses on the new architectural approach of RAN flexible functional splits. In order to support this novel architecture seamlessly, the traditional networks need to be transformed into open and dynamic infrastructures that can effectively manage the disaggregated resources. To this end, the SDN architecture was examined, which aim at the programmable and dynamic management of network resources through centralized control. The interest was focused on determining delays associated with the control related communication (signaling) required by this architecture. The development of effective optimization technics that identify the optimal functional split along with the optimal location and size of the SDN controllers, taking into consideration the overall system's stability was one of the main goals of this analysis. An outcome of this work can be found in [16]. The analysis was extended with advanced optimization schemes, in the interest of approaching a more realistic 5G environment, where the rapid increase in traffic implies the need for greater scaling capabilities to manage its spatial and temporal variations, as well as terminals with different quality requirements, and presented in [17].

The second contribution of this thesis focuses on the UPF selection problem of 5G-core networks. Based on the 5G-RAN deployment option and the type of service that needs to be provided, UPF nodes can be placed closer or further away from the 5G-RAN. Through this approach, UPF elements placed close to the network edge can redirect traffic to the MEC servers reducing latency, whereas UPF nodes placed deeper into the network can send traffic to central cloud facilities. In this thesis, the problem of optimal UPF selection in a 5G network, with the objective of minimizing the total service delay was examined. In this direction, in collaboration with colleagues, a 5G infrastructure was implemented in the laboratory using open-source tools. Parts of this work can be found at [18] and [19].

The third contribution concentrates on the service provisioning in future telecommunication systems. The analysis covered aspects of 5G- and 6G-systems supporting a variety of services through a common infrastructure that is effectively

shared through network slicing. The proposed approach involved new architectural Quality of Experience (QoE) models and resource allocation schemes that are capable of a) differentiating Data Streams originating from the same or multiple user equipment (UEs), b) reacting to changes of the underlying physical infrastructure, and c) scaling with the number of connected devices. Currently, this functionality is provided by centralized management and network orchestration platforms that suffer scalability issues. Therefore, future systems are expected to operate in a distributed manner allowing applications to directly intervene in the relevant control processes to guarantee the required QoE. The proposed approach focused on developing a flow assignment model that supports applications running on UEs. The work was published in [20].

The last contribution of this thesis concentrates on the actual deployment of a 5G infrastructure that supports on-demand multi-tenant network slicing. A network slice is formed utilizing dedicated access, transport, core and edge network resources as well as cross-domain components from various domains under the same or different administrations. Sharing of the underlying physical infrastructure is achieved using the concepts of hardware programmability and NFV. This involves the development of appropriate interfaces that can be used to integrate the different network components, as well as the creation of appropriate descriptors for the virtual 5G NFs (VNFs). The collection and appropriate combination of multiple VNFs gives an E2E 5G Network Service (NS). Through a Management and Orchestration (MANO) platform, we can combine these NSs in order to instantiate and manage a 5G network slice. All of our contributions are implemented and investigated experimentally in our 5G open-source lab testbed. Parts of this work were presented in the 5G-COMPLETE project [21].

1.3. Thesis Outline

This thesis is organized as follows:

Chapter 2 describes the overall 5G System Architecture. First, the overall 5G functional architecture is introduced. We proceed with analyzing the radio, transport and core network parts of the infrastructure. The innovative concept of slicing in 5G systems is then presented, followed by a description of the Quality of Service (QoS) architectural scheme, that is currently seen at the heart of 5G. Finally, the chapter concludes with a discussion about the Management and Network Orchestration of 5G systems.

Chapter 3 focuses on the selection of the optimal functional split in 5G-RAN. It proposes a hybrid centralized/distributed 5G network management solution that focuses on the optimization of FH flows. This problem is analytically solved adopting a novel mathematical model that allows RUs to dynamically adjust their FH split options with the objective to minimize their total operational expenditures. In this environment, the controller placement problem is also investigated, as this decision has a direct impact on the stability of the whole system. The stability of the proposed scheme depends on network latency, thus a metric for sizing the SDN transport network is proposed. Finally, the model was extended in order to approximate real-life scenarios where multiple RUs interact.

Chapter 4 concentrates on the UPF selection problem in 5G networks. UPF nodes have to be designed to support challenging 5G services with very tight performance requirements. According to the type of service that needs to be implemented, the number of UPF elements in the network varies. In this environment, the selection of the appropriate UPF to carry the traffic is of vital importance. To address this problem, we

propose a novel scheme that allows dynamic selection of the optimal UPF elements. The problem is formulated considering a specific optical node implementation and using accurate modeling of the delays introduced when this programmable optical node is adopted to act as UPF. Realistic cost functions for the model have been calculated using lab measurements derived from an open source 5G platform hosted in an optical datacenter cloud environment.

Chapter 5 proposes the introduction of network analytics functionalities in 5G systems, and studies how the users in a converged 5G network can be optimally served, based on their service QoS requirements. To address this problem, we propose a novel scheme that allows the UEs to dynamically select their service strategy, in order to optimize their service experience without uncontrollably exploiting the existing resources. For the evaluation of our solution, we considered an optical transport network and a UPF solution adopting a programmable optical node implementation. Real lab measurements derived from an open source 5G platform hosted in an optical datacenter cloud environment were used.

Chapter 6 presents the implementation of the 5G network slicing concept. The proposed solution concentrates on automating the deployment of a sophisticated 5G multi-operator network supporting E2E provisioning of multiple slices. For this purpose, we used Open Source MANO (OSM), an NFV-MANO compliant orchestrator and created network descriptors for the core and the user plane of a 5G network. Combining those descriptors, we successfully deployed dynamic 5G network slices on top of a softwarized multi-operator 5G platform hosted in a containerized data centre infrastructure.

Chapter 7 concludes the thesis and highlights the possible future works.

References

- [1] Marsch, P. and Bulakci, O. and Queseth, O. and Boldi, M., *5G System Design: Architectural and Functional Considerations and Long Term Research*, ISBN: 9781119425120, Wiley, 2018. Available: <https://books.google.gr/books?id=QFxDwAAQBAJ>.
- [2] Rao, S.K., Prasad, R., “Impact of 5G Technologies on Industry 4.0.,” *Wireless Personal Communications 100*, pp. 145–159, 2018. Available: <https://doi.org/10.1007/s11277-018-5615-7>
- [3] 5G PPP (2022). *5G and Verticals*. [Online]. Available: <https://5g-ppp.eu/verticals/>
- [4] ETSI (2022). *Why we need 5G?* [Online]. Available: <https://www.etsi.org/technologies/5G>
- [5] Dohler, M., & Nakamura, T., *5G Mobile and Wireless Communications Technology (A. Osseiran, J. Monserrat, & P. Marsch, Eds.)*, Cambridge: Cambridge University Press, 2016. Available: <https://doi.org/10.1017/CBO9781316417744>
- [6] Reply (2022). *5G Technology: Mastering the magic triangle*. [Online]. Available: <https://www.reply.com/en/industries/telco-and-media/5g-mastering-the-magic-triangle>
- [7] International Electrotechnical Commission, “Factory of the future”, White Paper 2015 [Online]. Available: <https://www.iec.ch/basecamp/factory-future>
- [8] 5G-PPP & European Commission (2016). *5G Vision - The 5G Infrastructure Public Private Partnership: the next generation of communication networks and services*. [Online]. Available:

- <https://espas.secure.europarl.europa.eu/orbis/sites/default/files/generated/document/en/5G-Vision-Brochure-v1.pdf>
- [9] S. Han et al., "Network support for resource disaggregation in next-generation datacenters", In *Proceedings of HotNets-XII. ACM*, New York, USA, 2013.
- [10] B. Han, V. Gopalakrishnan, L. Ji and S. Lee, "Network function virtualization: Challenges and opportunities for innovations," in *IEEE Communications Magazine*, vol. 53, no. 2, pp. 90-97, Feb. 2015, doi: 10.1109/MCOM.2015.7045396.
- [11] Open Networking Foundation (2019). *Software-Defined Networking (SDN) Definition - Open Networking Foundation*. [Online]. Available: <https://www.opennetworking.org/sdn-definition/>
- [12] A. Tzanakaki et al., "5G infrastructures supporting end-user and operational services: The 5G-XHaul architectural perspective," *IEEE ICC*, Kuala Lumpur, Malaysia, 2016.
- [13] M. Ruffini, "Multi-Dimensional Convergence in Future 5G Networks," *IEEE/OSA Journal of Lightwave technology*, Vol. 35, No. 3, March 2017.
- [14] U. Dötsch et al., "Quantitative Analysis of Split Base Station Processing and Determination of Advantageous Architectures for LTE," *Bell*
- [15] 5G-PPP (2022). *KPIS*. [Online]. Available: <https://5g-ppp.eu/kpis/>
- [16] Alevizaki, Victoria-Maria & Anastasopoulos, Markos & Tzanakaki, Anna & Simeonidou, Dimitra, "Joint Fronthaul Optimization and SDN Controller Placement in Dynamic 5G Networks", in *proc. of Optical Network Design and Modelling* 2019. Available: https://doi.org/10.1007/978-3-030-38085-4_16.
- [17] Alevizaki, VM., Anastasopoulos, M., Tzanakaki, A. et al., "Adaptive FH optimization in MEC-assisted 5G environments," *Photonics Network Communications* 40, 209–220, 2020. Available: <https://doi.org/10.1007/s11107-020-00906-8>
- [18] V. M. Alevizaki, M. Anastasopoulos, A. Tzanakaki and D. Simeonidou, "Dynamic Selection of User Plane Function in 5G Environments," in *Proc. of Optical Network Design and Modelling*, 2021.
- [19] V. M. Alevizaki, A. I. Manolopoulos, M. Anastasopoulos and A. Tzanakaki, "Dynamic User Plane Function Allocation in 5G Networks enabled by Optical Network Nodes," in *2021 European Conference on Optical Communication (ECOC)*, pp. 1-4, 2021. Available: <https://doi.org/10.1109/ECOC52684.2021.9606154>.
- [20] V. -M. Alevizaki, M. Anastasopoulos, A. -I. Manolopoulos and A. Tzanakaki, "Distributed Service Provisioning for Disaggregated 6G Network Infrastructures," in *IEEE Transactions on Network and Service Management*, Vol. 20, Issue 1, pp. 120-137, 2023. Available: <https://doi.org/10.1109/TNSM.2022.3211097>.
- [21] H2020 Project 5G-COMPLETE, Deliverable D6.2: "Report on the integration of 5G-COMPLETE technologies". [Online]

Chapter 2

5G System Architecture

Contents

2.2.	Network Architecture.....	11
2.2.1.	Overall Functional Architecture	11
2.2.2.	NG-RAN	16
2.2.3.	5G-Core	18
2.2.4.	Interfaces.....	20
2.3.	Network slicing.....	22
2.4.	QoS Architecture	24
2.5.	Management And Network Orchestration.....	26
2.6.	Summary.....	27
	References	27

2.1. Chapter Introduction

In order to meet the performance requirements of future services, which go far beyond the capabilities of the current 4G technologies, new approaches need to be implemented in the 5G and beyond network architectures. The current 4G network service provisioning model, which is based on abstract and rigid network requirements for generic user classes is inadequate to cope with the already observed and upcoming explosive growth of mobile Internet traffic. Thus, 5G and beyond networks are moving away from the traditional network element-driven architecture to service-driven approaches that allow to satisfy more flexible and sophisticated requirements on a per service basis, as determined by the various stakeholders, applications, and services.

From a technological standpoint, 5G and beyond networks evolve towards softwarised, distributed, multi-domain and technology ecosystems so that they can satisfy the strict objectives linked to services and performance. From a commercial standpoint, in-depth public discussions involving established and new parties have led to the introduction of

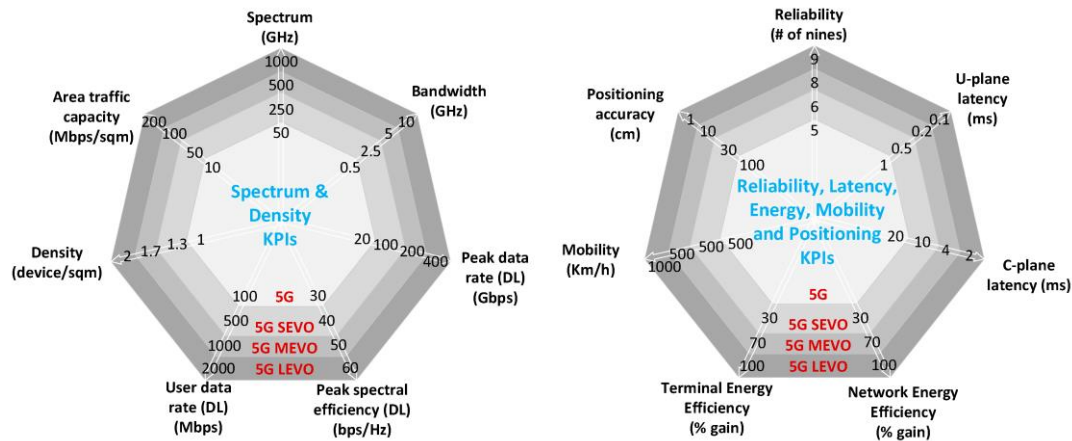


Figure 2. 1: KPIs Targets for 5G and beyond [15][18][19]

a variety of vertical applications, industries, and stakeholders that will be offered or enabled by 5G. Their description and the way they interact and fit in the 5G architecture are available in the literature ([1] - [7]). Except from the main vertical applications (m/eMTC, URLLC, eMBB) that will benefit from this architecture, 5G is being developed with the intention of enabling a wide range of future uses beyond those that are currently fully understood.

Supporting these services entails more than merely reducing latency or increasing bandwidth for specific users; [8]. To that end, 3rd Generation Partnership Project (3GPP) has established the crucial network service-related Key Performance Indicators (KPIs) and target values that must be met by 5G equipment and network deployments ([9] - [14]). Apparently, these goals won't be achieved immediately, but that's to be expected considering that building a new mobile network is a long-term process. At the same time, as new use cases emerge, these KPIs will continuously evolve to new target values in order to meet the requirements. **Figure 2. 1** illustrates the most important KPIs, by dividing the 5G evolution in phases, namely the today's 5G NR [16],[17], the short (SEVO), medium (MEVO) and long (LEVO)-term evolution period [18],[19]. Overall, the KPIs for 5G networks reflect the evolution and growth of the technology over time, with each stage bringing new challenges and opportunities for improving network performance and expanding connectivity. In the SEVO, the focus is on faster data rates, low latency, and spectrum efficiency. MEVO focuses on supporting mission-critical applications, mMTC, and energy efficiency. In the LEVO, the aim is on achieving ubiquitous connectivity, and developing Artificial Intelligence (AI) and Machine Learning (ML) capabilities. The SEVO KPIs are important for delivering eMBB applications, while the MEVO KPIs are crucial for supporting IoT applications. The LEVO KPIs are essential for creating a fully connected world where everything is connected to the internet. **Table 2. 1**[15] depicts the specific target values for each KPI throughout the four 5G evolution periods.

The rest of this chapter is structured as follows. First an overview of the overall 5G functional architecture is presented. Then an extensive representation of the main components of the architecture, that is the radio, core and transport network of 5G is described. The innovative principles of 5G slicing are then presented, followed by an

Table 2. 1:

KPIs AND TARGETS FOR 5G AND BEYOND TERMS [15][18][19]

KPI	5G	5G-SEVO	5G-MEVO	5G-LEVO
Spectrum	<56 GHz	<250 GHz	<500 GHz	<1000 GHz
Bandwidth	<0.5 GHz	<2.5 GHz	<5 GHz	<10 GHz
Peak Data Rate	DL > 20Gbps, UL>10Gbps	DL > 100Gbps, UL>50Gbps	DL > 200Gbps, UL>100Gbps	DL > 400Gbps, UL>200Gbps
User Data Rate	DL > 100Mbps, UL>50Mbps	DL > 500Mbps, UL>250Mbps	DL > 1Gbps, UL>0.5Gbps	DL > 2 Gbps, UL>1 Gbps
Peak Spectral Efficiency	DL > 30bps/Hz, UL>15bps/Hz	DL > 40bps/Hz, UL>20bps/Hz	DL > 50bps/Hz, UL>25bps/Hz	DL > 60bps/Hz, UL>30bps/Hz
Density	>1 Device/sqm	>1.3 Device/sqm	>1.7 Device/sqm	>2 Device/sqm
Area Traffic Capacity	>10 Mbps/sqm	>50 Mbps/sqm	>100 Mbps/sqm	>200 Mbps/sqm
Reliability	URLLC >5 nines	>6 nines	>8 nines	>9 nines
U-Plane Latency	URLLC < 1 ms	< 0.5 ms	< 0.2 ms	< 0.1 ms
C-Plane Latency	<20 ms	<10 ms	<4 ms	<2 ms
Network Energy Efficiency	Qualitative	>30% gain	>70% gain	>100% gain
Terminal Energy Efficiency	Qualitative	>30% gain	>70% gain	>100% gain
Mobility	<500 Km/h	<500 Km/h	<500 Km/h	<1000 Km/h
Positioning Accuracy	NA(<1m)	<30cm	<10cm	<1cm

introduction to 5G QoS architecture, that is currently seen at the heart of the 5G. Finally, the Management and Network Orchestration is discussed.

2.2. Network Architecture

2.2.1. Overall Functional Architecture

The overall E2E architecture of the 5G vision is shown in **Figure 2. 2**. 5G networks feature a diverse physical deployment that includes a variety of frequency bands, cell sizes and processing sites [20]. Given the overall ‘softwarization’ and ‘virtualization’ principles of the 5G vision, the core and access network functions are separated. The new logical point-to-point interface between the RAN and the Core Network (CN) will provide control and user plane separation. In order to be independent from the potential RAN deployment variants, it will need to be open and future-proof. The intention is for the new interface to support access-agnostic CN functionalities and permit autonomous development of the core and access networks.

To support the processing requirements of the various 5G services, the network architecture described above needs to offer access to compute and storage resources. The concepts of SDN and NFV open the road for the physical deployment of processing sites that will be used to support various network functions and services in accordance with the orchestrating services [21]. Those sites are divided into classes, depending on their distance in relation with the access network, and their computational capabilities [22]:

- **Centralized Data Center (Central Cloud- CC):** This class of resources includes high-performance servers and storage systems that are located in a centralized

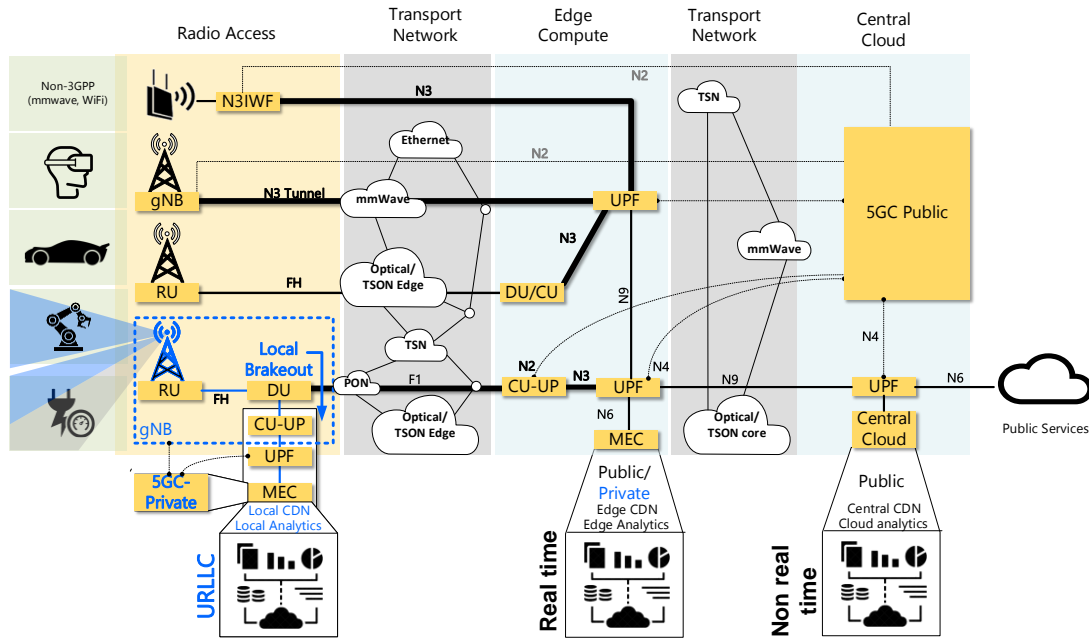


Figure 2. 2: E2E 5G System Architecture [15]

data center. These resources are used to support core network functions such as authentication, billing, and network management.

- MEC [23]: Edge computing resources include smaller data centers that are located near the network edge, closer to end-users, and devices. Its proximity to the access network may be able to reduce latency and jitter [24]. Depending on their proximity to the end-users, these resources can be used to support services that require low latency, such as augmented and virtual reality (AR/VR) applications, or even applications that require real-time data processing, such as autonomous vehicles, and industrial automation. MEC can also be used to support RAN functions, by hosting virtualized BBUs (vBBUs), decoupling them from the RUs. FH links connect the edge cloud with the RUs and require low latency and high-capacity attributes [15][25][26]. The FH connection speed is usually between 1 and 10 Gbps per antenna in order to meet the strict timing requirements of baseband processing. Specifically, since the delay for centralizing radio access protocol layers should be between 100 and 200s or less, a dedicated point-to-point connection between the edge cloud the and antenna sites, not greater than a few 10s of kilometers, should be employed as FH. It is important to multiplex (i.e., wave or time division multiplex) a number of these connections into a single fiber in order to increase efficiency. Ethernet switching needs to be studied further, because it adds more latency and delay jitter.

The antenna site can also accommodate the installation of dedicated Base Stations (BSs), which can be macro or small cells [27]. Macrocells are large, high-power cellular BSs that are used to provide wide-area coverage to a large number of users and are the backbone of the cellular network. On the other hand, small cells are low-power, short-range cellular BSs that are used to extend coverage and improve network capacity. 5G networks rely on a dense network of small cells to provide high-speed, low-latency connectivity to end-users. One of the primary benefits of small cells is that they can be deployed in areas where traditional macrocellular networks are unable to provide adequate coverage, such as indoor locations or densely populated urban areas. By

providing additional network capacity, small cells can help reduce network congestion and improve overall network performance. They also offer lower latency, which is essential for supporting real-time applications such as gaming and AR/VR. Small cells are typically lower in power consumption than traditional macrocells, which can help network operators reduce their energy costs and improve network efficiency. However, deploying a dense network of small cells can be expensive, as it requires a significant investment in infrastructure and equipment. Small cells also require ongoing maintenance to ensure that they are operating at optimal levels, which can be time-consuming and costly. Additionally, small cells operate on the same frequency bands as traditional macrocells, which can lead to interference and reduced network performance if they are not properly designed and deployed. Finding suitable locations to deploy small cells can also be challenging, as it requires access to public and private property and permission from local authorities. Despite these challenges, as 5G networks continue to evolve, small cells are expected to become even more important and essential for providing high-quality connectivity to the end-users.

Both macrocells and small cells in 5G networks are connected to the backbone network through a combination of wired and wireless backhaul (BH) technologies. The amount of user plane traffic traveling through a BS determines partly the required data rate of the BH link, while the acceptable latency is imposed by the requested radio functionality. Macrocells are typically connected to the backbone network through fiber optic cables, which provide high-speed, low-latency connectivity. These cables are typically buried underground or mounted on utility poles and are connected to the cellular BS via a wireless link. In small cell networks, the BH can be provided by a variety of technologies depending on the location and deployment scenario. For example, small cells deployed in outdoor environments may be connected to the core network through fiber optic cables, microwave or mmWave links, while small cells deployed in indoor environments may be connected to the core network through wired Ethernet or Wi-Fi links. In some cases, small cells may also be connected to other small cells in a mesh network configuration, where each small cell acts as a node in the network and can communicate with other nodes to provide connectivity to end-users. Such base stations typically use optical fibers for their BH.

Due to the diversity of the location and the capabilities of the processing sites, 5G Base Stations, namely the next-generation Node Bs (gNBs), are composed of three logical units: the CU, the DU and the RU. These units can be physically separated and placed at different network locations [28]. While their colocation offers backward compatibility with the LTE Enhanced Packet Core (EPC), the disaggregated architecture offers more advantages like the separation of the CP and UP, thus allowing further optimization for the location of different RAN entities.

Due to the disaggregated nature of gNBs, the transport network must be able to offer operational services in addition to BH services with the aid of FH/MH links [29]. FH links offer connectivity services between densely distributed RUs and regional data centers hosting CUs with extremely strict synchronization and latency requirements. A shared network infrastructure to jointly support BH and FH (and MH when necessary) functions is crucial, in order to maximize sharing benefits, offer enhanced efficiency in resource utilization and provide cost, scalability, and sustainability benefits. In such a converged transport network, latency still needs to be addressed as a primary criterion. To support this, the 5G architecture incorporates a variety of cutting-edge wireless and

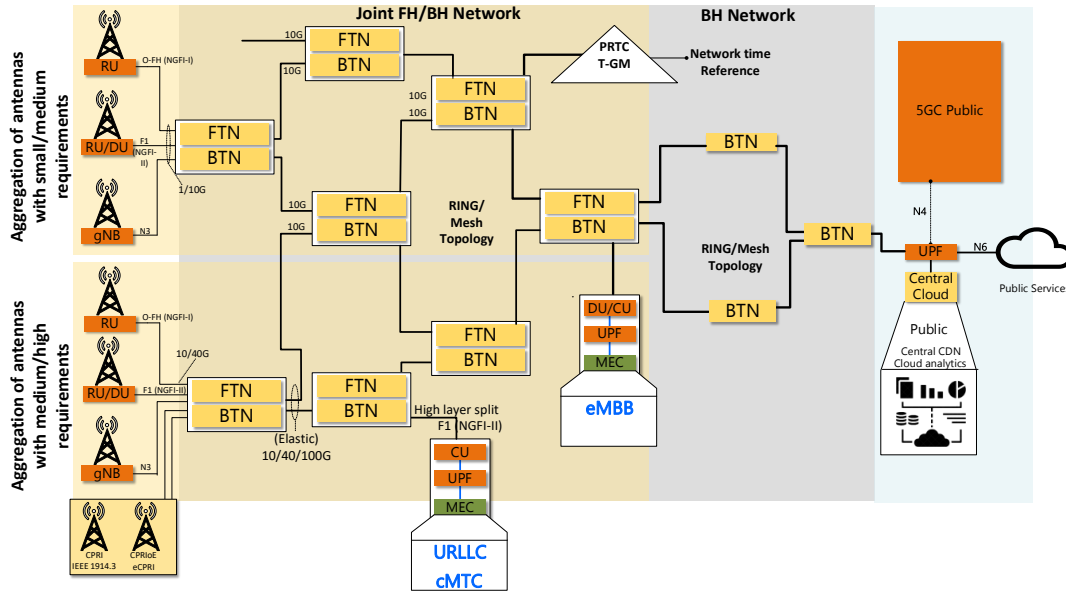


Figure 2. 3: Converged transport network. It includes both FH and BH transport nodes (FTN/BTN) capable of handling a range of connectivity options using wired and wireless technologies with capacities from 1Gbps to 100Gbps. [15]

wired access and transport network technologies [26]. The implementation of a high-capacity, flexible optical transport network combined with mmWave links serves as a promising solution. Although deploying mmWave links is less expensive than building optical fibers, the achievable capacity of mmWave links is in general lower. They are therefore only appropriate for backhauling sites with a small number of cells and modest data rates at the air interface, i.e., single small cells. Demanding capacity and flexibility requirements for traffic aggregation and transport can be supported by hybrid approaches including passive optical networks (PONs) offering enhanced capacity through wavelength division multiplexing (WDM) and dynamic, spectrally elastic, and frame-based optical network solutions [30]. **Figure 2. 3** provides an example of one such optical network, proposed by the 5G PPP project 5G-COMPLETE [15]. It is a multi-technology transport network that can handle FH and BH traffic using wired and wireless technologies with capacities ranging from 1Gbps to 100Gbps. Low-cost nodes can be installed near RUs to offer connectivity to radio transmission points and collect and aggregate traffic from various cell sites to a central location. The architecture addresses different cell sizes and can use higher capacity nodes to aggregate traffic from hub sites to edge sites. In the transport infrastructure different options for the integration of RU, DU, and CU entities exist, depending on the split options that is being adopted, the scale of the network and the latency requirements of mobile services. The mobile core entities can be located closer to the end user for low-latency services, while less latency-sensitive services can be handled at the core site. In terms of granularity, capacity, and flexibility, this transport network will be able to accommodate the high transport needs of 5G environments. mmWave technology, which operates in the sub-6 GHz and 60 GHz bands, will provide high mobility in diverse environments through dynamic beamforming and programmable beam steering techniques, while also enabling high bandwidth connection for both non-line-of-sight (NLOS) and line-of-sight (LOS) scenarios.

A key feature offered by 5G as already discussed is decoupling the control and user plane functionalities. This is enabled through two fundamental 5G concepts, namely the architectural microservices [24] and the network function decomposition. The concepts refer to the service-based architecture (SBA) approach where an application consists of small independent services that communicate through well-defined Application Programming Interfaces (APIs). This architecture allows for each service to scale or update without disrupting other services in the application. 5G Microservices enable the design of logical architectures customized according to the performance and functional requirements of each use case [25]. The basic idea is to split the traditionally monolithic NFs, which frequently correspond to physical network elements in previous cellular systems, into basic modules or NFs defined with adequate granularity. The modularization of the 5G network functions separates the NFs that relate to the access network (AN) with those of the CN. This decreases the dependency of the 5G core on the access (and vice versa), and provides new ways of connectivity other than cellular radio [26]. Furthermore, it separates the CP from the UP functions. This would enable the definition of various logical architectures through the interconnection of various subsets of CP and UP NFs. Along with the flexibility in network deployment and operation, CUPS provides efficient and cost-effective traffic management. The UP is responsible for the transfer of the user data and comprises the end-to-end path between the UE, the RAN, the UPF and the Data Network (DN). The UPF is responsible for the actual transmission of the data. It also performs packet inspection and routing, as well as UP QoS handling. The CP, on the other hand, comprises of multiple NFs that are communicating through the SBA and is accountable for decision-making and managing the network. It interacts with the NG-RAN through two 5G-core functions, namely the Access and Mobility Management Function (AMF), and Session Management Function (SMF). AMF is responsible for the user subscription and authentication as well as for the security and mobility management. SMF is in charge of calculating the appropriate routes and communicating its policies to the UP elements. QoS management is also under the control of the SMF. NG-RAN, on its part, is responsible for routing the CP and UP data, managing the radio resources (RRM) and the QoS Flows, creating new network slices, and setting up the UE connection. The communication of the UP elements with the CP elements is achieved through point-to-point interfaces [31].

The key differentiation between 5G and 4G lies in its QoS architecture, which is built on Protocol Data Unit (PDU) sessions. This allows for the creation of several ‘QoS flows’ between the UEs and the destination DN via the 5G system's UPF. The PDU session comprises a group of end-to-end packet flows, referred to as Service Data Flows (SDFs), specific to the application [32]. The initial step towards establishing the PDU session is to register the UE with the 5G core AMF. Then, the AMF verifies if the UE request matches the user's subscription by contacting the SMF, which manages PDU sessions. The SMF seeks assistance from two other 5G-core network functionalities, the Unified Data Manager (UDM) and the Policy Control Function (PCF), to perform this task [15][33]. While the UDM provides information on the user subscription and authentication, the PCF supplies policy rules for QoS Flows and flow-based charging to the SMF. Once the user subscription is authenticated, the SMF determines the appropriate user plane path to support the requested service and shares necessary information, such as PDU Session ID and QoS characteristics, with the UPFs. In the final step, the AMF notifies the RAN and the UE with the QoS information retrieved from the SMF [7]. In the final step, the AMF notifies the RAN and the UE about the QoS information obtained from the SMF. Once all entities involved are aware of the QoS-related details, the PDU session is established for both Uplink (UL) and Downlink (DL) data.

2.2.2. NG-RAN

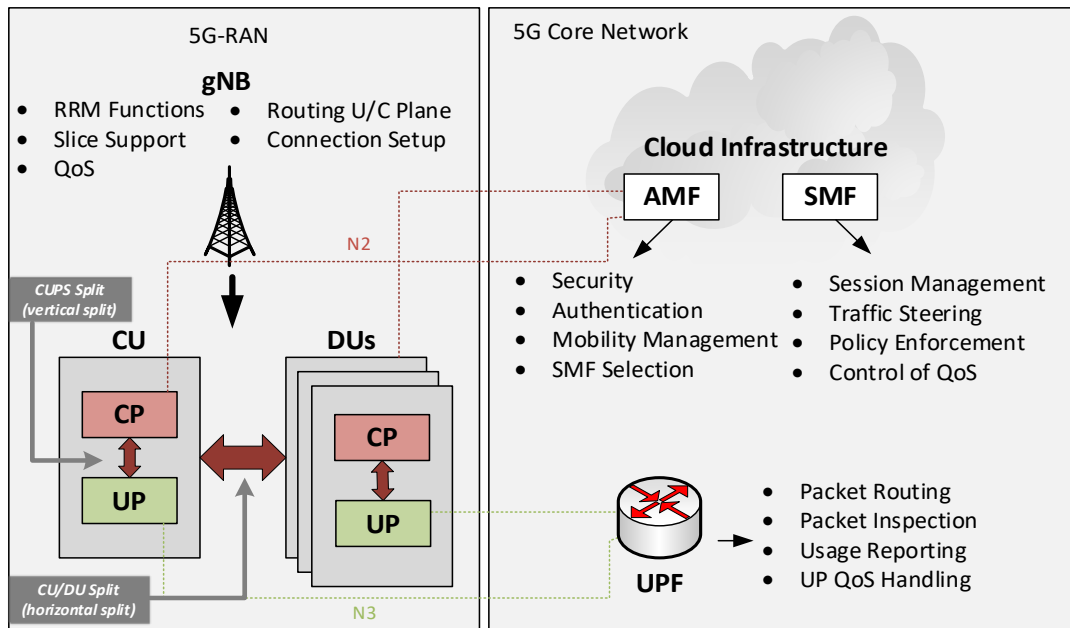


Figure 2. 4: Functional split between NG-RAN and 5G-CORE. 5G supports two dimensional splits: a CP/UP (functional-vertical) split and a CU/DU (geographical-horizontal) split. [22]

The introduction of multiple processing sites in 5G opens the road for different RAN deployments. RAN in 5G systems is redesigned to support various functional splits along two dimensions (**Figure 2. 4**): A “vertical split” referring to CUPS, and a “horizontal” split of functions between processing sites depending on their time requirements. This architecture enables a more flexible mapping of NFs to physical network entities, depending on the use case and deployment constraints [34].

The vertical split opens the road for the employment of SDN in the RAN [35]. It uncouples the CP and UP operation, making possible for them to be scaled and optimized independently. Each function can be upgraded separately, without requiring the modification of other NFs. Furthermore, the standardization of a global interface for the CP will allow a consistent control over network entities and NFs from different vendors. Despite the benefits of CUPS, there are some concerns that should be taken into consideration. Whenever a new feature is introduced, the CP/UP interface must be newly standardized, and that can cause significant delays in the delivery of these new features. Apart from this, the most important impediment of the vertical split is the tight interconnection of CP and UP functions, especially in the lower protocol stack layers. For example, a BH connection between CP and UP is optimal for the proper employment of the Media Access Control (MAC) scheduler.

On the other hand, horizontal split refers to the split of computational resources required for baseband processing from the BSs (the RUs) and their placement at a co-location facility (DU or CU), according to the service requirements, thus minimizing the energy footprint and maximizing efficiency [26][29]. The connection between the RUs and the DUs/CUs is performed through FH links. In the case where the DUs and the CU are located in different sites, their connection is performed through MH links. The

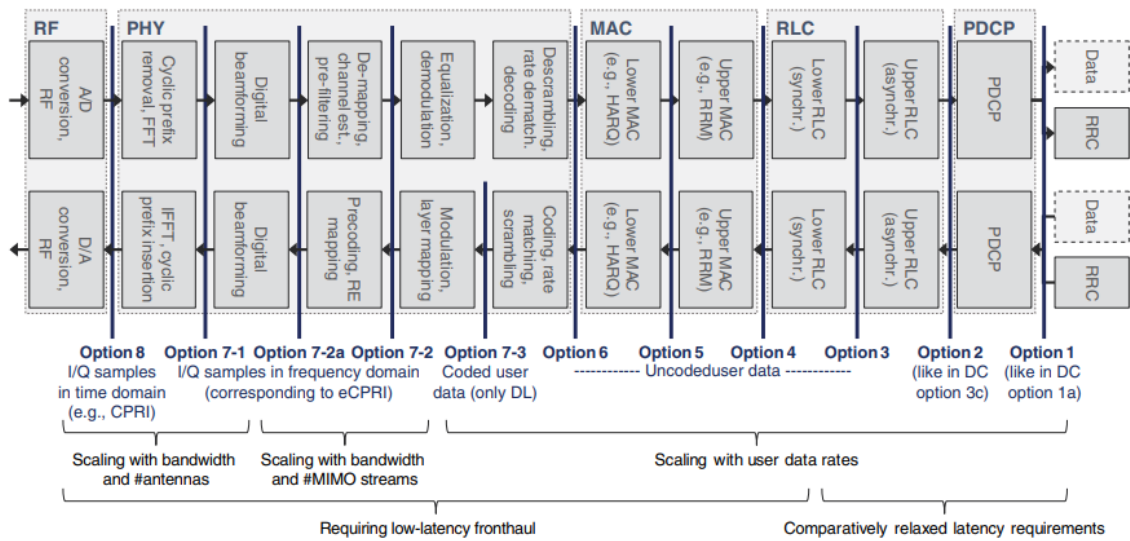


Figure 2. 5: Horizontal splits options for the 5G RAN. [34][36]

architecture of Flexible Functional Splits (FFSs) of RAN functions can offer centralized resource management and multi-cell processing, leading to significant performance and cost gains. The location of the NFs is based on the time requirements of the service, for example the whole protocol stack can be placed at the edge for time-sensitive applications, whereas for delay-tolerant applications a set of the processing functions can be allocated at the central cloud. This way the RAN processing can be customized according to different deployment characteristics, like local processing power of the infrastructure, BH or FH characteristics etc. In general, the DU mainly hosts some or all of the lower layer-2 protocols and physical layer processing such as Lower Radio Link Control (RLC), MAC, and Physical (PHY), while the CU functions include non-real-time upper layer processing (Upper RLC), Packet Data Convergence Protocol (PDCP), and Service Data Adaptation Protocol (SDAP)- a new protocol that is introduced for the QoS Flow management between the UE and the RAN - as well as core network functions that have been moved to the edge of the network in order to enable MEC services. Additionally, in an effort to minimize the transport capacity and latency requirements between RU and DU(s), some of the PHY processing can be moved to the RUs. In order to enable connectivity in the proposed architecture, new logical interfaces are introduced for the gNB and specifically F1 has been defined to support connectivity between CUs and DUs [15].

Figure 2. 5 shows the main horizontal splits that have been proposed by the 5G research area [36]. Each functional split has its own characteristics in terms of acceptable data rate and latency. In general, the greater the number of NFs that are placed close to the radio, the lower the latency and data requirements for the related interfaces. However, this causes reduced centralization benefits. Functional splits can be different for the UL and DL. Depending on the position in the processing chain where the split happens, there are six options for high-layers splits, that are splits above the physical layer, and 3 options for splits within the physical layer. In terms of simplicity for the description we will consider the existence of only remote and central units (DUs are collocated with the CUs). Options 1 and 2, shown in **Figure 2. 5**, are already utilized in LTE systems. Option 3 distinguishes the processing between asynchronous and synchronous functionalities, by placing the entire RLC functionality in the CU, apart from the

aggregation functionality. Split option 3 is mostly pursued, because it offers the highest centralization gains and enables improved flow control. Options 5 and 6 locate some or all of the MAC processing in the CU. These splits are less favored by the 3GPP, because they require sub-frame-level timing interactions between CU and RU, and FH delay would affect timing and scheduling of Hybrid Automatic Repeat Request (HARQ). Split 7 has three alterations, according to where the split within the physical layer happens. Option 7-1 performs the Inverse Fast Fourier Transform (iFFT) and cyclic prefix insertion/removal at the RU, and Quadrature Signals (IQ) in frequency domain are exchanged over the interface. In this way only samples related to occupied sub-carriers need to be exchanged, instead of time domain samples reflecting the whole system bandwidth. Also, less quantization bits per symbol may be needed when quantizing in frequency domain. In the UL, some Physical Random Access Channel (PRACH) processing may also be performed at the RU. Options 7-2 and 7-23 perform the same functions with option 7.1 at the RU with the addition of the pre-coding and digital beamforming. Thus, the FH requirements scale with the number of Multiple-Input Multiple-Output (MIMO) layers, and not with the number of antenna ports as in the case of options 7-1 and 8. Options 7-2 and 7-2a differ in the extent of precoding happening locally or centrally, or the location where channel estimation is performed in the uplink etc. Option 7.3 is only considered for DL, and further reduces bandwidth requirements on the interface, as coded user data is exchanged before modulation. As a downside, such split would likely strongly increase the complexity of the local unit. Finally, option 8 is the fully centralized option, where all functions are performed at the central cloud. Although this split offers all the benefits of centralization, and a simple local unit deployment, the capacity requirements of the FH are significantly raised, in relation to other splits.

The different RAN functional split render the distinction between traditional BH and FH more ambiguous. On one hand, higher layer split points, such as the separation of PDCP and RLC, have bandwidth and latency requirements that are more in line with conventional BH than FH. Even for such split points, it is obvious that effective methods for high-bit-rate transport and multiplexing will be needed due to the significant increase in the volume of traffic that needs to be aggregated. On the other hand, latency and jitter are becoming more of a critical limitation for the lower RAN layer split points. Thus, traditional dimensioning approaches for the transport infrastructure will no longer be applicable with 5G. A more relevant way to categorize the new FH requirements is the “loose” or “tight” latency requirements. The new 5G transport network will need to support these varying bandwidth needs and aggregation of RAN functional splits, which may demand that many, various split point interfaces be delivered over the same physical infrastructure. The heterogeneity of transport network equipment must also be addressed in order to reduce costs by integrating data, control, and management planes across all technologies as much as possible. A novel redesign of FH/BH network segment is necessary, since the capacity needed for next-generation radio interfaces, which use 100MHz channels and squeeze the bit-per-MHz ratio utilizing massive MIMO or even full-duplex radios, cannot be fulfilled simply by advancing existing technology.

2.2.3. 5G-Core

To perform network control functionalities the 5G architecture relies on a number of core NFs. These include [7][4][22]:

- **Access and Mobility Management Function (AMF):** The AMF is involved in most of the signaling call flows in a 5G network, and it supports encrypted signaling towards device, allowing them registration, authentication, and moving between cells. It interacts with Radio Network using a new protocol, namely the New Generation Application Protocol (NGAP), and with the devices with Non-Access Stratum (NAS) 5G Mobility Management (5GMM) messages.
- **Session Management Function (SMF):** The SMF is on the one hand responsible for setting up the connectivity between the UE and the Data Networks, and on the other hand for managing the user plane functionality for that connectivity. It manages the end device sessions such as establishment, modification and release of individual sessions, and allocation of IP addresses per session. The SMF is designed with the flexibility to support various types of end-user protocols, different options to ensure service continuity, as well as a flexible user plane architecture. Additionally, it interacts with the PCF to retrieve policies for the PDU sessions which it then passes to the UPF, and it is responsible for collecting charging data.
- **User Plane Function (UPF):** Its main task is to process and forward user data and its functionality is controlled from the SMF. It connects with external IP networks and acts as a stable anchor point for the devices toward external networks, hiding the mobility. This means that IP packets with a destination address belonging to a specific device is always routable from the Internet to the specific UPF that is serving this device even as the device is moving around in the network. Moreover, it performs different types of processing if the forwarded data. It also generates traffic usage reports for the SMF, which then includes them in charging reports to other NFs. The UPF can also apply packet inspection, to analyze the content of user data packets, which it can use either as input to policy decisions, or as a basis for traffic reporting. Additionally, it can redirect traffic or apply different data rate limitations. It also acts as buffer when a device is in idle state or unreachable from the network. It applies QoS marking of packets towards the radio network or to external networks. This can be used from the transport network to prioritize packets in case of congestion inside the network.
- **Unified Data Management Function (UDM):** It executes functions requested from the AMF. It also generates the data used to authenticate attaching devices. Moreover, it authorizes access to specific users based on subscription data, for example, applying different access rules for roaming subscribers and home subscribers. The UDM also keeps track of which instance is serving which device, in case of more than one AMF and SMF existing in the network.
- **Unified Data Repository (UDR):** It is basically a database that stores various types of data such as, subscription data and data defining various types of network or user policies. The data stored in the UDR are commonly offered as services to other NFs, namely, UDM, PCF and NEF.
- **Authentication Server Function (AUSF):** It has a limited but important functionality which is to provide the service of authenticating a specific device, in that process utilizing the authentication credentials created by the UDM. Moreover, it generates cryptographical material to allow for secure updates of roaming information and other parameters in the device.

- **Policy Control Function (PCF):** The PCF is a key element in the 5G Core and it interacts with various NFs. It is in charge of providing policy control of functionality for the SMF, AMF, UE access selection and PDU session selection. It can also support Negotiation of future data transfers. In relation to the SMF, the PCF provides QoS and charging control for SDFs, as well as policy control and event reporting for the PDU sessions. In relation with the AMF, it offers access and mobility policy control, including management of service area restrictions and of the Rat Frequency Selection Priority (RFSP) index. Finally, the PCF interacts with the UE via the AMF to provide policy information such as discovery, session continuity mode selection, network slice selection and more.
- **Network Repository Function (NRF):** The NRF is a repository that keeps the profiles of all NFs available to the network. Each time a NF (Service Consumer) needs to receive some service from another NF (Service Provider), will simply have to look the NRF to find the most suitable NF. Newly deployed or changed NFs have to report their new profile information to the NRF. Profile information updates can either be triggered by the NF itself or by another entity on behalf of the NF. NF profiles in the NRF include information such as the NF type, address, capacity etc.
- **Network Exposure Function (NEF):** The NEF is responsible for the exposure of events and capabilities from the 5G System (5GS) to the applications and NFs of the operator's network or a third-party network. Events such as the UE location, reachability, roaming status and loss of connectivity can be monitored by the NEF and then be made available to specific applications and NFs. Authorized applications by the network can use the NEF for specific requests, such as QoS and charging policies.
- **Application Function (AF):** Applications that are considered trusted, either inside the operator's network or outside of it, can communicate directly with the other Core Network Functions and influence some of their aspects, e.g. an application that runs on a MEC server may influence the traffic routing decisions. AFs can interact with the Core NFs either directly or via the NEF.
- **Network Slice Selection Function (NSSF):** In 5GS, different networks can be virtually divided to isolated slices. For example each use case (eMBB, URLLC, mMTC) can be served by a different slice. The NSSF is aware of the existence of all network slices and the AMF(s) that are dedicated to each slice, and selects the set of slices and AMF(s) that should serve the UE.

2.2.4. Interfaces

5G Networks (5GNs) introduce new ways of interconnecting the involved network elements compared to the previous LTE EPC technology. Utilizing the notion of architectural microservices [37] and network function decomposition [38], the evolution of mobile communication systems towards 5G was specifically designed to provide a network architecture that allows for each service to scale or update without disrupting other services in the network.

The SBA was adopted for the connection of the 5G core functions [39][40]. Through specifically designed Service Based Interfaces (SBIs), each core NF offers one or more services to other NFs in the network. These services are made available over NF interfaces that are connected to the common Service-Based Architecture. In practice

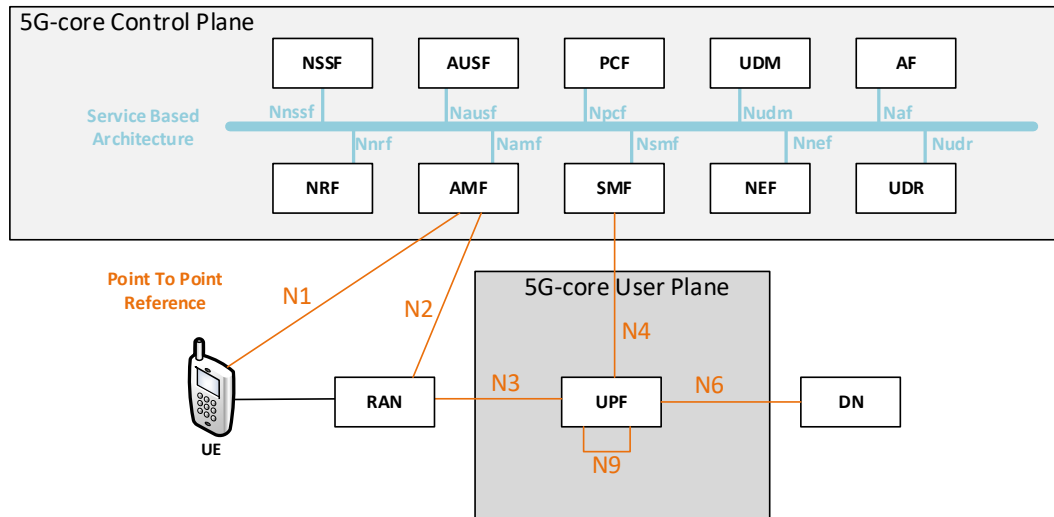


Figure 2. 6: 5G-CORE Service-Based and Point to Point architecture.

this means that Network Functions are accessible by other Network Functions over an API. The defined communication method for 5G Core is the HTTP Representational State Transfer (REST) API. It is important to note that SBA is used for signaling functionality and not for the transfer of actual user data. The SBI concept is expected to make network capabilities easier to extend, compared to the traditional point-to-point architecture which relies on detailed and extensive protocol specification efforts.

When two NFs interact over the 3GPP SBA, a role is assigned to each of them. The NF that requests a service has the role of a Service Consumer, while the NF offering the service is assigned the role of a Service Producer. Service Discovery, i.e. when a Service Consumer needs to locate a certain Service Producer, is realized by making use of the NRF. The NRF is, as it was mentioned, a repository where all available services of all NFs are registered and it is in charge of connecting the NFs in order for a service to be delivered.

Interaction between the UP elements (UE-RAN-UPF-DN) as well as between some CP NFs (SMF, AMF) and the UP components, communication is achieved through their point-to-point interfaces. The most commonly used point to point interfaces are [42][31][7]:

- N1: Connects the UE with the AMF entity of the CN for NAS messages
- N2: Used for the signaling between the RAN and the CN (AMF) through the NGAP protocol
- N3: Used for data connectivity between the RAN and the UPF through a General Packet Radio Service Tunnelling Protocol (GTP) tunnel for the user plane (GTP-U protocol)
- N4: Used for communication between the SMF and the UPF over the Packet Forward Control Packet (PFCP) protocol
- N6: Interconnects the Anchor UPF with the external DN
- N9: Interconnects different UPF entities across a given CN.

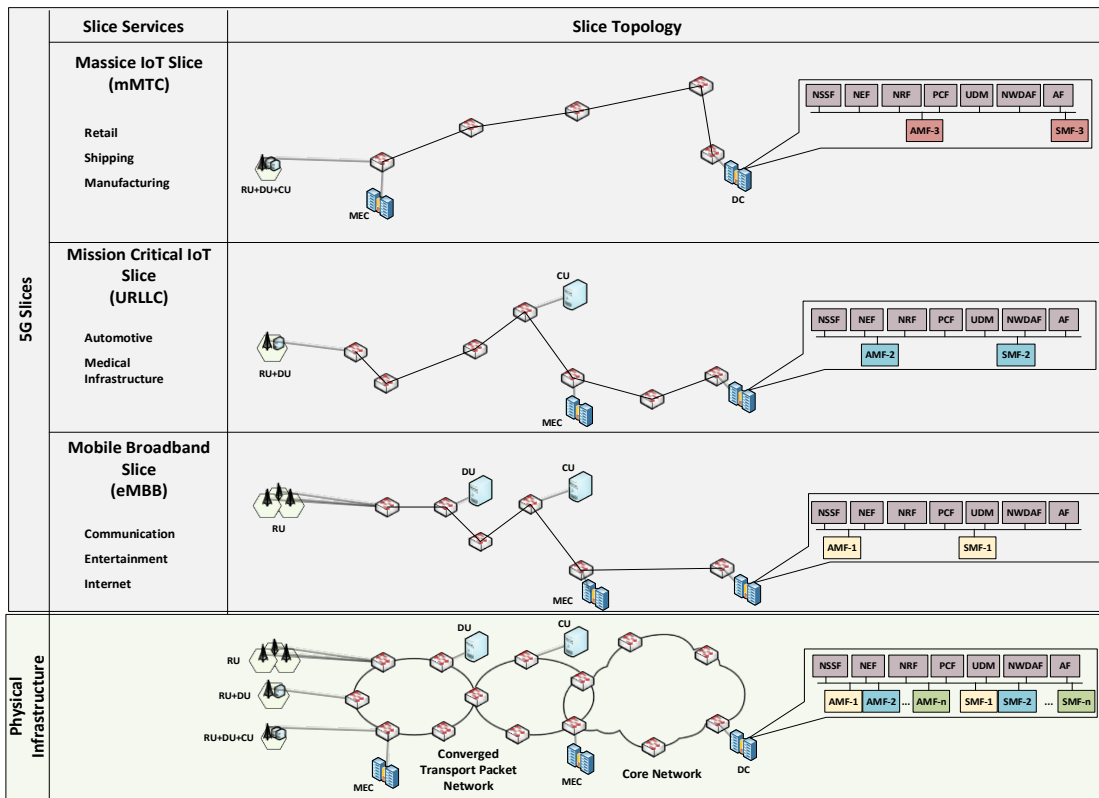


Figure 2. 7: 5G Network slicing overview [15]

Figure 2. 6 shows the SBA architecture of the 5G core, as well as the point-to-point interfaces of the basic NFs of 5G core.

2.3. Network slicing

5G networks are expected to support a wide range of end-devices with various system needs and features, including tablets, smartwatches, massive IoT devices, mission-critical IoT devices, and other gadgets. This will create a heterogeneity in terms of KPIs among different services, that cannot be fulfilled by current network infrastructures. To achieve the desired KPIs and criteria at technical and commercial levels, Beyond 5G (B5G) technological solutions will need to include not just sophisticated network technologies disaggregated at the hardware-software level, or at the control-data plane, but also architectural solutions and technologies that provide flexible deployments on a per use case basis. A possible solution would be the establishment of multiple physical networks, one for each service or business entity. Therefore, services would be isolated, using their own resources and without the need for laborious hardware and network reconfigurations. It is obvious that this method cannot be used in real networks in a cost-effective way, hence a solution that enables effective resource sharing and multi-tenancy infrastructure use is required.

This diversity of services to be addressed by the operators along with the need for efficient resource utilization of the 5G infrastructure led to the idea of “network slicing”, which gives network operators the option to divide the physical network into various end-to-end virtualized networks [41][43]. Figure 2. 7 shows an overview of the network

Table 2. 2:

INDICATIVE, FUTURE SLA STRUCTURE [15][47]

SLA PARAMETERS	DEFINITION
TIME PERIOD	The actual time over which the service is provided, incl. the start and end date (e.g. 1/1/2020 - 30/4/2020)
PERIODICITY	The periodicity of the service offering; e.g. “Continuous”, “Every weekend”, 8:00am – 17:00pm, etc.
LOCATION	The location of the ANNs to be connected (where the end-users are expected to reside) (e.g. at a specific stadium area/ municipality/ building block(s)/campus/ hotspot (shopping mall, park, etc.), Nationwide, etc.
CLOUD SERVICE RESOURCES	The Cloud Service Resources in terms of vCPUs, Memory-RAM, Storage Space, etc.
AVAILABILITY	The time within a specific time-period in which the service is up and running, or inversely the max. time within a specific time-period in which the service is unavailable.
RELIABILITY	Frequency (number of times) of non-availability of the service per specific time period.
MONTHLY AVAILABILITY	The percentage (%) of time within a specific time-period in which the service is up and running.
SCALING RULES	The set of pre-defined rules based on which the service can be scaled up.
OPERATIONAL RULES	The set of pre-defined rules based on which specific operations are triggered, e.g. events' warning, monitoring initiation, etc.

slice concept. The radio access, transport, and core networks of each of the sliced networks are each specifically designed to support a certain class of device. There is no interference between the other slices, since the slices are provisioned with independent network resources. Mission-critical IoT services benefit greatly from network slicing, which also accomplishes dynamic network resource allocation for various traffic conditions.

As it was mentioned, the main idea behind network slicing is the partition of traffic in multiple logical networks that are all executed on and share a common physical infrastructure [22]. Those logical networks are independent from each other and offer distinct Service Qualities, named Service Level Agreement (SLA). This way the 5G system can be optimized to match specific requirements of the different services or operators/tenants [44]. A variety of users can create their own (virtualized) network instance by using network resources and functions. This will enable the monetization of so-called vertical sectors. With network slicing, tenants will be able to reconfigure the behavior of individual NFs and adapt the network topology according to their needs [45]. The end-to-end nature of network slicing in the context of 5G, as well as the requirement to express a service through a high-level description and to flexibly map it to the appropriate infrastructural elements and NFs, are perhaps the most significant factors that set 5G slicing apart from other types of slicing that have been considered in the past (such as cloud computing).[46]

In the business realm, the introduction of network slicing evolves the SLAs and requirements from straightforward connectivity QoS guarantees (for retail or wholesale purposes) to complicated service descriptions, which may include or combine factors like [15]:

- Versatile Connectivity (Both data rate and latency, and in some cases jitter) QoS guarantees over multiple links.
- Compute resource requests.
- Specific reliability/ availability/ security thresholds.
- Specific time constraints and allowances, even requested ad-hoc.
- Scaling capabilities in terms of resources/ number of links etc.
- Operational aspects and additional functionalities.

Indicatively, **Table 2. 2** [47] presents a few examples of how SLAs can evolve for vertical and network operator services in the 5G context. Even with the finest granularity, these requirements may not always be achieved by basic end-user classification. To accommodate the variety of vertical and network operator services, an overall slicing architecture would be required. Therefore, slicing at three distinguished layers is being proposed, namely the network layer slice, the infrastructure/resource layer slice, and the network management layer slice ([48], [49]). Each layer, must provide suitable interfaces, physical resources, and virtualization capabilities for the correct establishment of the end-to-end logical network.

2.4. QoS Architecture

As was previously established, 5G will need to simultaneously support a wide range of use cases. For this purpose, versatile and extremely granular E2E tools are required to handle and prioritize various traffic types, or even different packets that belong to the same traffic type. Previous cellular systems cannot cope with this challenge, since traditionally in LTE, the differentiation of mobile data happens at the radio bearers (left part of **Figure 2. 8**). All packets are handled equally inside a single bearer, which can be configured to reflect guaranteed or non-guaranteed bit rate traffic and is identified by a QoS class identifier (reflecting priority, tolerable latency, and packet loss rate). A one-to-one mapping between radio bearers and Evolved Packet System (EPS) bearers is also present. Since in 5G a single wireless communication link may carry data related to very different services, for the separation of the traffic one should set up multiple individual radio bearers, thus rendering the system very inefficient. Additionally, packets within a PDU session cannot be treated differently by LTE. Therefore, it is widely agreed that 5G must offer a considerably more granular approach to QoS management, enabling more flexible and independent QoS handling in the CN and RAN, and in particular permitting a packet-specific traffic differentiation when necessary.

The 3GPP-approved E2E QoS management architecture for 5G is based on the idea of QoS flows [42][22]. Each QoS flow has a unique ID (QFI), and is managed by the SMF entity of the CN. A key differentiation to previous systems is that a PDU session can contain multiple QoS flows. The routing of the packets inside a QoS flows is performed by the SMF. After the calculation of the appropriate route, the SMF informs the access networks in terms of a QoS profile, that contains information about the flow (e.g. 5G QoS Identifier -5QI, Allocation and Retention Policy -ARP, if it is guaranteed or not bit rate traffic etc.). The SMF also provides QoS rules to the UE, and SDF classification to the UPF, that is a template dictating the routing of the packets inside a QoS Flow [34].

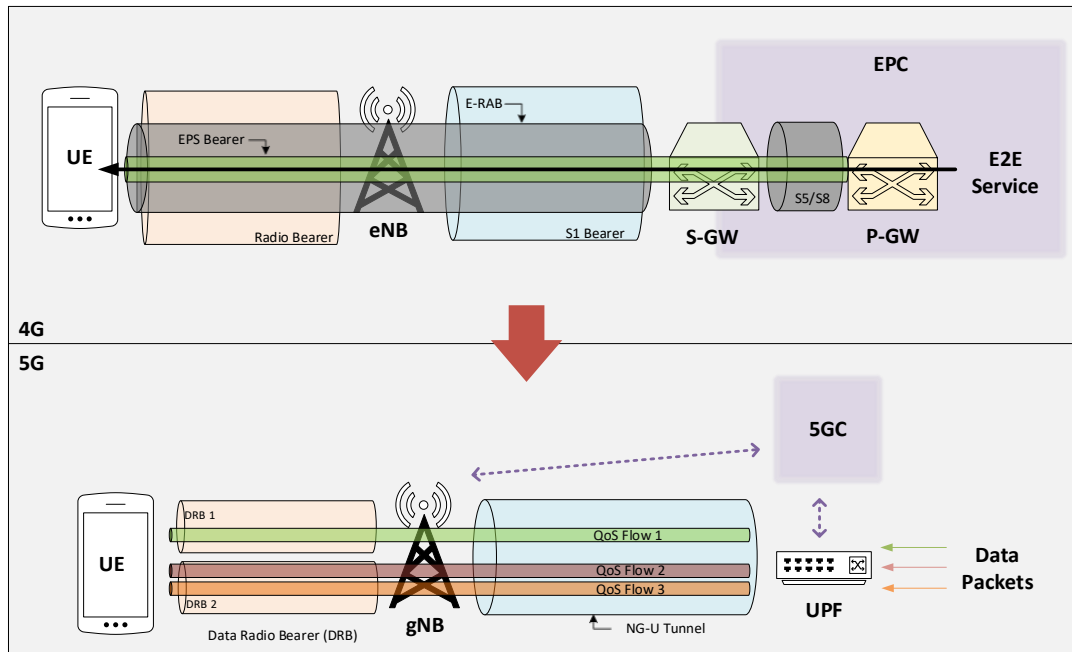


Figure 2. 8: Evolution of the QoS architecture from 4G to 5G.[22]

Specifically, the following processing happens for the DL and UL respectively. In order to assign specific data packets to the defined QoS flows in the DL, the UPF in the CN uses the SDF classification rules. Then, QoS flows can be arbitrarily assigned to various Data Radio Bearers (DRBs) in the access network, according to their QFI value and the QoS profile that is provided to the RAN. In contrast to the LTE approach, where there is a precise mapping of EPS bearers to radio bearers, in 5G systems this can be handled separately from the CN. With this new architecture design, there exist many optimization opportunities in both the assignment of flows to radio bearers and the assignment of packets to QoS flows. For example, urgent QoS flows could be assigned by the access network to lower carrier frequencies, with additional methods for greater reliability, such as multi-connectivity, while less critical QoS flows could use mmWave frequencies. In the UL, the UE classifies packets to QoS flows, according to the QoS rules that it received from the AMF. The QoS Flows are then mapped to available data radio bearers [32].

A comparison between LTE and 5G QoS architecture is illustrated in **Figure 2. 8**. In the left the LTE’s “barrier oriented” approach is shown, with a strict one-to-one mapping between EPS bearers, EPS radio access bearers (ERAB), S1 and radio bearers. On the other hand, 5G uses a “flow-based” approach as depicted in the right of **Figure 2. 8**. The assignment of packets to flows is decoupled from the assignment of flows to DRBs. The former is a function performed by the CN, while the latter is carried out by the RAN. In this new architecture, the SDAP protocol is introduced on top of PDCP, for the QoS Flow management between the UE and the RAN [31]. In the downlink it maps a specific QoS Flow inside a PDU session to a suitable DRB. In the uplink it marks the packets with the appropriate QFI, in order to be treated appropriately by the core network.

2.5. Management And Network Orchestration

It is extremely challenging to handle customized service creation and rapid delivery in very short periods, as it is envisaged in 5G networks, when networks are constructed in the traditional fashion. A core prerequisite for 5G systems has been the support of flexible and configurable network architectures, so that they can adapt to any use case and service requirements [50]. With the introduction of SDN [51] and NFV [52], 5G realizes network services by dynamically deploying functions and programming communication channels among them rather than re-architecting the network. In order to guarantee an effective utilization of the infrastructure while meeting the performance and functional requirements of heterogeneous services, the MANO plane is crucial [53]. Services will no longer be deployed, configured and management on a node-by-node basis, but in an integrated and coordinated way. This approach is related to the removal of individual device configuration in favor of a more robust management mechanism that can offer network-wide service design, configuration, deployment, and monitoring. Complex networking systems, resources, and services can be automated in their arrangement and coordination thanks to orchestration. Such a process requires implicit autonomic control over all systems, resources, and services as well as inherent intelligence.

This network-wide orchestration has the primary benefit of providing a single point of integration and a centralized representation of the distributed network, regardless of the number of resources involved or the locations of those resources. Any manageable component with a set of features (capacity, connectivity etc.), which is related to a physical or virtual network (such as an optical or packet network), or to a data center (e.g., compute or storage) can be referred to as a resource. This opens up the potential for the implementation of advanced services across all network domains and offers significant prospects for achieving the required degree of automation and targeted KPIs. By automating the underpinning configuration and monitoring processes from the orchestration module, it is possible to reduce the inherent complexity of delivering and administering sophisticated and multi-featured services [54].

To manage and configure every component of the entire network service simultaneously, higher-level abstractions and automated methods are required [55]. Abstraction enables representation of an entity in terms of chosen attributes, which are shared by similar resources and that are handled and controlled, while concealing or summarizing characteristics unrelated to the selection criteria. The administration of resources can be generalized and made simpler by the abstraction, removing the initial impediments caused by manufacturing differences, in particular the technology used to create them, or the resources' physical realization.

Taking into consideration 5G network slicing, there are two main levels of network orchestration [55]. The first one corresponds to the inter-slice orchestration that deals with the orchestration of resources to accommodate different network slices in the network. The second corresponds to the intra-slice orchestration, that refers to the orchestration of resources within a network slice. Furthermore, different levels of orchestration are needed when delivering a service. Initially, the resources required to support a certain service should be correctly assigned and configured in accordance with the requirements of the service that will be supported. This is a task allocated to the resource orchestration level. It is not necessary for resource orchestrators to comprehend the service logic or the topology that establishes the relationship between the NFs that make up the service; instead, they solely deal with resource level

abstraction. The service logic as requested by the customer concerns the service level orchestration. It specifies the functions necessary to satisfy the customer request as well as the way in which these functions interact to offer the whole service. The NFs in the underlying infrastructure will be dynamically instantiated by the service orchestrator.

2.6. Summary

5G promises to provide higher data rates, lower latency, and better reliability compared to its predecessor, 4G. This chapter discusses various aspects of the overall 5G system architecture. First, the 5G functional architecture, that is designed to support a wide range of use cases and applications, from eMBB and mMTC to URLLC, was examined. The main components of this architecture were thoroughly investigated, namely NG-RAN and the 5G-CORE. The NG-RAN is responsible for providing radio access to the end-users, while the 5G-CORE is responsible for providing 5G services connectivity and control. Following the examination of the 5G architectural components, we focused on a key concept of 5G, referred to as slicing. 5G slicing enables the network to be divided into multiple virtual networks, each with its own specific requirements and characteristics. This allows the network physical infrastructure to be segmented and each network segment to be independently customized in support of different use cases. To enable 5G slicing, the novel QoS mechanisms that were introduced in 5G systems were also analyzed. Finally, the chapter concludes with a discussion on 5G MANO that provides the ability to monitor and control network resources in order to ensure efficient operation of the 5G network.

References

- [1] 5g-PPP (2022). *5G PPP TB & 5G IA Verticals TF, "Empowering Vertical Industries through 5G Networks - Current Status and Future Trends, Version 1.0.* [Online]. Available: <https://5g-ppp.eu/wp-content/uploads/2020/09/5GPPP-VerticalsWhitePaper-2020-Final.pdf>.
- [2] 5G PPP (2022). *5G Empowering Vertical Industries.* [Online]. Available: https://5g-ppp.eu/wp-content/uploads/2016/02/BROCHURE_5PPP_BAT2_PL.pdf.
- [3] *IMT vision—Framework and overall objectives of the future development of IMT for 2020 and beyond*, Recommendation ITU-R M.2083-0 09/2015, Geneva, Switzerland, 2015. [Online]. Available: https://www.itu.int/dms_pubrec/itu-r/rec/m/R-REC-M.2083-0-201509-1!!PDF-E.pdf
- [4] 5G-PICTURE, EU funded project, Deliverable D7.5: "Commercial exploitation and market Impact document". [Online]. Available: https://www.5g-picture-project.eu/download/5g-picture_D7.5.pdf
- [5] I. Mesogiti, E. Theodoropoulou, G. Lyberopoulos, F. Setaki, A. Ramos, P. Gouvas, A. Zafeiropoulos and R. Bruschi, "A Framework to Support the Role of Telecommunication Service Providers in Evolving 5G Business Models," *IFIP Advances in Information and Communication Technology*, vol. 560.
- [6] 5G-VINNI, EU funded project, Deliverable D5.1: "Ecosystem analysis and specification of B&E KPIs". [Online]. Available: https://zenodo.org/record/3345665#.YA-7_ugzaUk
- [7] *View on 5G Architecture*, 5G PPP Architecture Working Group, Version 3.0, 06-2019. [Online]. Available: https://5g-ppp.eu/wp-content/uploads/2019/07/5G-PPP-5G-Architecture-White-Paper_v3.0_PublicConsultation.pdf

- [8] 5G-PPP (2020). *5G and Verticals*. [Online]. Available: <https://5g-ppp.eu/verticals/>
- [9] NGMN 5G Initiative Team, "A Deliverable by the NGMN Alliance: NGMN 5G White Paper". [Online]. Available: <https://www.ngmn.org/5g-white-paper/5g-white-paper.html>
- [10] *Feasibility Study on New Services and Markets Technology Enablers for massive Internet of Things; Stage 1*, 3GPP TR 22.861. [Online].
- [11] *Feasibility study on new services and markets technology enablers for critical communications; Stage 1*, 3GPP TR 22.862. [Online].
- [12] *Feasibility study on new services and markets technology enablers for enhanced mobile broad-band; Stage 1*, 3GPP, TR 22.863. [Online].
- [13] *Feasibility study on new services and markets technology enablers; Stage 1*, 3GPP TR 22.891. [Online].
- [14] *Feasibility study on new services and markets technology enablers for network operation; Stage 1*, 3GPP TR 22.864. [Online].
- [15] H2020 Project 5G-COMPLETE, Deliverable D2.1: "Initial report on 5G-COMPLETE network architecture, interfaces and supported functions". [Online]
- [16] *Study on Self-Evaluation Towards IMT-2020 Submission*, 3GPP TR 37.910. [Online]. Available: <https://itectec.com/archive/3gpp-specification-tr-37-910/>
- [17] *Minimum Technical Performance Requirements for IMT2020 Radio Interfaces*, ITU-R IMT-2020. [Online]. Available: https://www.itu.int/en/ITU-R/study-groups/rsg5/rwp5d/imt-2020/Documents/SO1-1_Requirements%20for%20IMT-2020_Rev.pdf
- [18] Alain Mourad, Rui Yang ,Per Hjalmar Lehne and Antonio De La Oliva, "A Baseline Roadmap for Advanced Wireless Research Beyond 5G", *Electronics* 2020, 9(2), 351;; [Online]. Available: <https://doi.org/10.3390/electronics9020351>.
- [19] Joint EU-US programme, EMPOWER, Deliverable D2.2: "First technology roadmap for advanced wireless," [Online]. Available: https://www.advancedwireless.eu/wp-content/uploads/Deliverables/EMPOWER_deliverable_D2_2_final.pdf.
- [20] The O-RAN alliance (2021). [Online]. Available: <https://www.o-ran.org/>.
- [21] Y. Mao, C. You, J. Zhang, K. Huang and K. B. Letaief, "A Survey on Mobile Edge Computing: The Communication Perspective," in *IEEE Communications Surveys & Tutorials*, vol. 19, no. 4, pp. 2322-2358, Fourthquarter 2017. Available: <https://doi.org/10.1109/COMST.2017.2745201>.
- [22] Marsch, P. and Bulakci, O. and Queseth, O. and Boldi, M., *5G System Design: Architectural and Functional Considerations and Long Term Research*, ISBN: 9781119425120, Wiley, 2018. Available: <https://books.google.gr/books?id=QFxDwAAQBAJ>.
- [23] H. Liu, F. Eldarrat, H. Alqahtani, A. Reznik, X. de Foy and Y. Zhang, "Mobile Edge Cloud System: Architectures, Challenges, and Approaches," *IEEE Systems Journal*, vol. 12, no. 3, pp. 2495-2508, 2018. Available: <https://doi.org/10.1109/jsyst.2017.2654119>
- [24] R. A. Addad, D. L. C. Dutra, M. Bagaia, T. Taleb and H. Flinck, "Fast Service Migration in 5G Trends and Scenarios," in *IEEE Network*, vol. 34, no. 2, pp. 92-98, March/April 2020. Available: <https://doi.org/10.1109/MNET.001.1800289>

- [25] M. Satyanarayanan et al., “An open ecosystem for mobile-cloud convergence,” in *IEEE Communications Magazine*, vol. 53, no. 3, pp. 63-70, March 2015. Available: <https://doi.org/10.1109/MCOM.2015.7060484>
- [26] Tzanakaki et al., “Wireless-Optical Network Convergence: Enabling the 5G Architecture to Support Operational and End-User Services,” in *IEEE Communications Magazine*, vol. 55, no. 10, pp. 184-192, Oct. 2017. Available: <https://doi.org/10.1109/MCOM.2017.1600643>
- [27] M. Kamel, W. Hamouda and A. Youssef, “Ultra-Dense Networks: A Survey”, in *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, pp. 2522-2545, 2016. Available: <https://doi.org/10.1109/comst.2016.2571730>
- [28] M. A. Adedoyin and O. E. Falowo, “Combination of Ultra-Dense Networks and Other 5G Enabling Technologies: A Survey,” in *IEEE Access*, vol. 8, pp. 22893-22932, 2020. Available: <https://doi.org/10.1109/ACCESS.2020.2969980>
- [29] M. Peng, Y. Li, Z. Zhao and C. Wang, “System architecture and key technologies for 5G heterogeneous cloud radio access networks,” in *IEEE Network*, vol. 29, no. 2, pp. 6-14, March-April 2015. Available: <https://doi.org/10.1109/MNET.2015.7064897>
- [30] Tzanakaki, A. Manolopoulos, M. Anastasopoulos, and D. Simenidou, “Optical Networking in Support of User Plane Functions in 5G Systems and Beyond,” in *Photonics in Switching and Computing 2021*, pp. W2B.3, January 2021. Available: <https://doi.org/10.1364/PSC.2021.W2B.3>
- [31] Rommer, S., Hedman, P., Olsson, M., Frid, L., Sultana, S., Mulligan, C., *5G Core Networks: Powering Digitalization*, ISBN: 9780081030097, Elsevier Science, 2019.
- [32] Devopedia (2021). *5G Quality of Service*, Version 3. [Online]. Available: <https://devopedia.org/5g-quality-of-service>
- [33] *5G; 5G System; Policy and Charging Control signalling flows and QoS parameter mapping; Stage 3*, ETSI TS 129 513 V16.6.0 (01/2021d). [Online]
- [34] *5G; NR; NR and NG-RAN Overall description; Stage-2*, ETSI TS 138 300 V16.4.0 (2021c). [Online]
- [35] Open Networking Foundation (2019), *Software-Defined Networking (SDN) Definition*. [Online]. Available: <https://www.opennetworking.org/sdn-definition/>
- [36] Cpri.info (2019). *eCPRI Specification V1.1*. [Online]. Available: http://www.cpri.info/downloads/eCPRI_v_1_1_2018_01_10.pdf
- [37] Microservices.io (2020). *What are microservices?* [Online]. Available: <https://microservices.io/>
- [38] Ian F. Akyildiz, Shuai Nie, Shih-Chun Lin, Manoj Chandrasekaran, “5G roadmap: 10 key enabling technologies”, *Computer Networks*, Volume 106, 2016, Pages 17-48, ISSN 1389-1286. Available: <https://doi.org/10.1016/j.comnet.2016.06.010>.
- [39] B. Chatras, “Applying a Service-Based Architecture Design Style to Network Functions Virtualization,” in *IEEE Conference on Standards for Communications and Networking (CSCN)*, Paris, 2018.
- [40] F. T. Kuhn, F. Schnicke and P. O. Antonino, “Service-Based Architectures in Production Systems: Challenges, Solutions & Experiences,” in *ITU Kaleidoscope: Industry-Driven Digital Transformation (ITU K)*, Ha Noi, Vietnam, 2020.

- [41] K. Katsalis, N. Nikaein, E. Schiller, A. Ksentini and T. Braun, "Network Slices toward 5G Communications: Slicing the LTE Network," in *IEEE Communications Magazine*, vol. 55, no. 8, pp. 146-154, Aug. 2017, Available: <https://doi.org/10.1109/MCOM.2017.1600936>
- [42] *System architecture for the 5G System (5GS)*, 3GPP TS 23.501 version 16.6.0 Release 16. [Online]
- [43] NGMN (2020). *Description of Network Slicing Concept*. [Online]. Available: <https://www.ngmn.org/publications/description-of-network-slicing-concept.html>.
- [44] Huawei.com (2020). *5G Network Slicing for Vertical Industries*. [Online]. Available: <https://www.huawei.com/minisite/5g/img/5g-network-slicing-for-vertical-industries-en.pdf>
- [45] Alcardo Alex Barakabitz, Arslan Ahmad, Rashid Mijumbi, Andrew Hines, "5G network slicing using SDN and NFV: A survey of taxonomy, architectures and future challenges," *Computer Networks*, Volume 167, 2020, 106984, ISSN 1389-1286. Available: <https://doi.org/10.1016/j.comnet.2019.106984>
- [46] X. Foukas, G. Patounas, A. Elmokashfi and M. K. Marina, "Network Slicing in 5G: Survey and Challenges," in *IEEE Communications Magazine*, vol. 55, no. 5, pp. 94-100, May 2017. Available: <https://doi.org/10.1109/MCOM.2017.1600951>
- [47] I. Mesogiti et al., "Network Services SLAs over 5G Infrastructure Converging Disaggregated Network and Compute Resources," *2018 IEEE 23rd International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, Barcelona, Spain, 2018, pp. 1-5. Available: <https://doi.org/10.1109/CAMAD.2018.8514989>
- [48] *5G Service-Guaranteed Network Slicing White Paper*, China Mobile Communications Corporation, Issue V1.0, February 2017. [Online]. Available: <https://www-file.huawei.com/-/media/corporate/pdf/white%20paper/5g-service-guaranteed-network-slicing-whitepaper.pdf?la=en>
- [49] ITU-news (2021). *Why end-to-end network slicing will be important for 5G*. [Online]. Available: <https://news.itu.int/why-end-to-end-network-slicing-will-be-important-for-5g/>.
- [50] L.M. Contreras, P. Doolan, H. Lønsethagen and D.R. López, "Operation, organization and business challenges for network operators in the context of SDN and NFV", in *Elsevier Computer Networks*, Vol. 92, pp. 211–217, 2015
- [51] ONF, "SDN Architecture", Issue 1, 2016. [Online]. Available: <https://opennetworking.org/wp-content/uploads/2014/10/TR-521-SDN-Architecture-issue-1.1.pdf>
- [52] etsi.org (2020). *Network Function Virtualization (NFV)*. [Online]. Available: <http://www.etsi.org/technologies-clusters/technologies/nfv>
- [53] *Networks Functions Virtualization (NFV); Management and Orchestration*, ETSI GS NFV MAN 001, V1.1.1, Dec. 2014
- [54] *Management and Orchestration; Provisioning*, 3GPP TS 28.531 January 2020. [Online]
- [55] *Management and orchestration; 5G Network Resource Model (NRM) (Release 16)*, 3GPP TS 28.541 5G, March 2020. [Online]. Available: <https://portal.3gpp.org/desktopmodules/Specifications/SpecificationDetails.aspx?specificationId=3400>.

Chapter 3

Adaptive Fronthaul Optimization in MEC Assisted 5G environments

Contents

3.1.	Chapter Introduction	31
3.2.	Theoretical Background.....	34
3.2.1.	Evolutionary Game Theory: Basic Concepts	34
3.2.2.	Dynamics of Multi-Agent Learning	35
3.2.3.	Evolutionary Dynamics of Reinforcement Learning	35
3.3.	Application to Wireless Networks in 5G	36
3.3.1.	Problem formulation using EGT.....	37
3.3.2.	Problem formulation using MARL	40
3.4.	Results and Discussion	41
3.4.1.	SDN Controller placement using EGT.....	41
3.4.2.	Optimal Split Selection using MARL.....	45
3.5.	Summary.....	46

3.1. Chapter Introduction

The digital evolution observed in the modern world, renders existing technologies inefficient, since an increased number of devices needs to connect to the network, thus creating huge demands in terms of capacity and latency [1]. In view of this, 5G network technologies are aiming at a new open and flexible network paradigm that will satisfy the needs of various vertical operators (transport, media, automotive, manufacturing, healthcare etc.) in a cost and energy efficient manner [2], [3].

The huge increase of data traffic that is foreseen in the next years, introduces the need for higher network densification, a concept that is being implemented employing a large number of cells with limited coverage, also known as small cells. This way, network's capacity can be increased, while the end-to end delay decreases [4]. However, this solution comes at the expense of increased Capital and Operational Expenditures (CAPEX and OPEX) as a large number of new BSs needs to be adopted and operated contributing also to the increase of the CO₂ footprint of the infrastructure.

One way to overcome the aforementioned limitations is by decoupling the computational resources (the BBUs) from the BSs (the RUs) and place them in a co-location facility (the CU), minimizing the footprint and maximizing efficiency. The RUs are connected with the CU through high bandwidth and low latency FH links. Through centralized signal processing and management, significant benefits are expected including improved planning of the shared wireless access network, reduced intercell interference between adjacent cells and increased spectral efficiency, better mobility management and faster handovers, quick and easy network upgrades. The centralized RAN architecture can be also combined with Cloud Computing, a concept known as Cloud-RAN/C-RAN) [5]. Despite the benefits of C-RAN, this architecture can be challenging for Mobile Network Operators (MNOs). This is mainly due to the need for high capacity transport links to support fronthaul services. Furthermore, moving away from optimized-RAN dedicated hardware to virtualized RAN functionalities can be costly inefficient in terms of installation, testing, integration, operation and maintenance [1].

The adoption of MEC along with C-RAN can make the above-mentioned investment much more efficient [6]. The placement of processing equipment closer to the edge of the access network is necessary for the realization of the envisioned 5G vertical services[7]. The vision of ubiquitous broadband access with high quality user experience for both mobile and stationary users (eMBB), in combination with mMTC, will immensely augment the amount of data volume. In addition, with the increased network traffic, the strict delay and capacity requirements of 5G URLLC, render necessary the processing of time sensitive applications, such as autonomous driving, remote robotic surgery and augmented reality etc., closer to the user [5], [8].

The challenges and cost associated with the transport network requirements of C-RAN can be balanced by splitting the processing of 5G protocol stack functions between the CU and the RUs, using FFSs [1], [9]. As it was mentioned in the previous chapter, with FFS a subset of the RAN functions is performed remotely at the RU taking advantage of the processing capabilities of the MEC to which the RU is connected, and the remaining functions are executed at the CU. Hence, the FFS technique can be implemented adopting an architecture able to assign the BBU processing functions dynamically between MEC servers and large-scale DCs placed at the optical access and metro domains, respectively. Efficient management and orchestration of this architecture can be achieved applying novel network designs that are aligned with the SDN open reference architecture [10]. SDN separates the control and data planes and as such it migrates the switch control externally of the physical devices to a logical entity the controller. The controller is in charge of populating the forwarding table of the switches. The communication between the two entities is carried out through a secure channel. This centralized structure enables the controller to perform network management functions, while allowing easy modification of the network behavior through the centralized control layer. However, this type of architecture suffers increased end-to-end latency due to the communication requirements between the controller and the physical devices, particularly when the scale of the underlying infrastructure increases. In addition, as the size of the infrastructure grows, the number of variables and parameters that should be considered by the SDN controller during calculation of the optimal routing policies (i.e. forwarding rules) increases exponentially. This results in an increase of the flow set up time having a negative effect on the quality of the offered services. To address this problem two key solutions have been proposed. The first is associated with the identification of the optimal location and size of the SDN controllers so that the data-to-control plane latencies are kept below an acceptable level [11], [12].

The second approach is associated with the development of hybrid traffic management policies according to which some of the decisions are taken centrally by the SDN controller, while some other are taken locally [13]. This approach reduces the number of variables that should be handled by the centralized SDN control providing the optimal share between performance (in terms of flow setup time) and overhead (controller synchronization and management).

Taking into consideration this approach, in the present study we propose a hybrid centralized/distributed 5G network management solution that focuses on the optimization of FH flows. The proposed solution relies on the following key technologies, including

1. the concept of FFS, according to which RUs are able to individually select their preferred split option policy, to unilaterally optimize their performance. In order to minimize control plane overheads, FFS decisions are taken at the RUs in a non-cooperative/self-optimizing manner.
2. the employment of MEC resources in the form of specific purpose low processing power servers embedded in the wireless access network (also known as cloudlets). These resources are used to provide the necessary processing power for lower level physical functions and,
3. an SDN controlled FH/BH transport network connecting the MEC domains with medium to large-scale DCs hosting general purpose servers placed at the optical access and metro domains. These SDN controllers are responsible to calculate the necessary forwarding rules and populate the associated forwarding tables at the network switches.

The above, result in a hybrid control scheme according to which high level FH connectivity decisions are taken by the centralized controller, while local decisions associated with the optimal FFS are taken by the RUs in a non-cooperative manner. RUs periodically evaluate the selected FFS scheme and adopt a new policy if a better performing split option has been identified. This problem is analytically solved adopting a novel mathematical model based on Evolutionary Game Theory (EGT) that allows RUs to dynamically adjust their FH split options with the objective to minimize their total operational expenditures. In this environment, the controller placement problem is also investigated, as this decision has a direct impact on the stability of the whole system. The stability of the proposed scheme depends on network latency, thus a metric for sizing the SDN FH/BH network is proposed. Finally, Multi-Agent Reinforcement Learning (MARL) was used for approximating real-life scenarios where the pairwise interactions between the RUs do no longer hold.

The rest of the chapter is organized as follows. After a brief overview of EGT and MARL in Section 3.2 the problem under consideration is analyzed in Section 3.3. Then, its application to the proposed network model is presented in Section 3.4, the controller placement problem is addressed and the optimal split option is identified in simple and more complexed schemes. Finally, conclusions are drawn in Section 3.5.

3.2. Theoretical Background

3.2.1. Evolutionary Game Theory: Basic Concepts

EGT studies the interactions of non-cooperative players that play repeatedly strategic games [14]. Contrary to classic Game Theory that examines the behavior of rational players, EGT focuses on how the strategies can "survive" through evolution and how they can help players who choose them to "strengthen" and better meet their needs.

Evolutionary processes are described by three main components: the population, the game and a dynamic set of equations that is used to model the processes. The most common dynamics is called the Replicator Equation (RE) and can be expressed as:

$$\dot{x}_i(t) = x_i(t) \left(F_i(\mathbf{x}(t)) - \bar{F}(\mathbf{x}(t)) \right), \quad i \in S \quad (1)$$

where S is the set of strategies that are available to the population, $\mathbf{x}(t) = [x_1(t) \ x_2(t) \ \dots \ x_i(t) \ \dots]^T$ is the population state at time t with $x_i(t)$ being the proportion of the population that uses strategy i at time t . $F_i(\mathbf{x}(t))$, $\bar{F}(\mathbf{x}(t))$ are used to describe the expected payoff of strategy i and the mean payoff, respectively [14]. According to this equation, the growth rate \dot{x}_i/x_i of the strategies that are currently used is equal to the excess of the current payoff $F_i(\mathbf{x}(t))$ over the average population's payoff, $\bar{F}(\mathbf{x}(t))$. This means that a selected strategy will either survive or eliminated in the long run depending on whether its payoff is better or worse than the average payoff of all strategies.

The replicator equation can be extended to multiple populations to cover scenarios where the interacting agents are drawn from different populations. In this case, each population is characterized by a state vector \mathbf{x}_i that depicts the distribution of the available pure strategies inside the population. At each interaction a player is drawn from each population to play the game. The multipopulation dynamics equation can be described as follows:

$$\dot{x}_{ih}(t) = x_{ih}(t) \left(F_i(e_i^h, \mathbf{x}_{-i}(t)) - \bar{F}(\mathbf{x}(t)) \right), \quad i \in N, h \in S \quad (2)$$

where N is the number of populations, S is the set of strategies that are available to the populations, $\mathbf{x}_i(t) = [x_{i1}(t) \ x_{i2}(t) \ \dots \ x_{ih}(t) \ \dots]^T$ is population's i state at time t with $x_{ih}(t)$ representing the proportion of the population i that uses strategy h at time t , and $F_i(e_i^h, \mathbf{x}_{-i}(t))$, $\bar{F}(\mathbf{x}(t))$ is the expected payoff of strategy i and the mean payoff of all strategies, respectively [14].

In the above, the time duration of the interaction between individuals is assumed to be infinitesimal and the corresponding payoffs are awarded immediately. However, this assumption is not valid under realistic problem settings. Specifically, in communication networks, the impact of a selected action may be applied to the involved entities with delay due to network latency. Thus, it is more realistic to consider a system where the strategies evolve considering the fitness values perceived in a past moment. The adjusted RE when delay τ is introduced is given below [15],[16]:

$$\dot{x}(t) = x_i(t) \cdot \left(f_i(\mathbf{x}(t - \tau)) - \sum_{j \in S} x_j(t) \cdot f_j(\mathbf{x}(t - \tau)) \right) \quad (3)$$

3.2.2. Dynamics of Multi-Agent Learning

A Multi Agent System (MAS) can be seen as a loosely connected network of independent agents that can work together to solve problems that go beyond the scope, resources and capabilities of each individual agent. Many sophisticated problems of modern society can be approximated using MAS, such as urban and air traffic control, multi-robot coordination, distributed sensing, energy distribution and load balancing, finance and auditing [18]-[19].

MAS systems are decentralized. This means that there is no central control, and no agent has complete information about the system. Decision making in MAS is a collective process and the actions of other agents can interfere with one's action plan. Each agent makes decisions that are always relevant to him in the context of overall MAS coordination. The agents can act optimally only when they know the environment with which they interact. However, such a knowledge is often unattainable, because the environment changes according the joint action of all the agents (non-stationarity). Therefore, each agent is faced with a moving-target learning problem: the best policy changes as the other agents' policies change [20], [21]. Thus, learning in these systems is imperative. The agents must interact repeatedly with the environment in order to learn its dynamics.

This kind of process is the main objective of MARL. It combines the technics used by classical Reinforcement Learning (ReLe), with the concepts used in game theory and economics. Classical ReL refers to systems with only one agent that tries to learn and optimize his strategy. However, due to the non-stationary nature of a MAS environment, ReL methods more than often fail to guarantee the conversion to optimal policy [22]. To confront this challenge, concepts from EGT are employed to describe learning in MAS. Just like MAS, EGT deals with dynamic and uncertain environments in which, agents with incomplete information try to optimize their behavior [19]. The combination of MAS with EGT enables the detailed analysis of the learning dynamics in MAS, thus facilitating the comparison and parameter tuning of different learning algorithms.

3.2.3. Evolutionary Dynamics of Reinforcement Learning

The first link between the two fields was established by Borgers and Sarin [23], who proved the convergence of a basic ReL algorithm - Cross Learning - to the most common population dynamics of EGT - the replicator equation - in the continuous time limit.

In this model, a large number of agents is considered to play the same normal-form game repeatedly in discrete time. At each repetition each agent chooses a strategy from his strategy with some probability and receives a payoff for that strategy. Players' choices are described as random because they are affected by some unmodelled psychological factors [23]. The payoffs represent reinforcement experiences and in the Cross model are always positive. The higher the payoff of the selected strategy, the higher is the likelihood that the agent will select this strategy in the next repetition. In this way, the strategy probabilities of each agent adjust over time in response to experience [23].

The strategy probabilities change as follows. When the i^{th} agent plays the n^{th} repetition (stage) of the normal form game, he selects a strategy j (according to his strategy

probability distribution) and receives a payoff U_j^i ; then he updates the probability $P_j(n + 1)$ to select strategy j at stage $n + 1$ using the following set of equations:

$$P_j(n + 1) = \theta \cdot U_j^i + (1 - \theta \cdot U_j^i) \cdot P_j(n) \quad (4)$$

$$P_{j'}(n + 1) = (1 - \theta \cdot U_j^i) \cdot P_{j'}(n) \text{ for all } j' \neq j \quad (5)$$

In (4)-(5), θ represents the time interval between the repetitions of the game. It is proved that if $\theta \rightarrow 0$, $n\theta \rightarrow t$ the above equations yield to a continuous time system modeling the evolution of the probability p_{ij} an agent i to select strategy j that can be expressed by the multipopulation RE:

$$\dot{p}_{ij}(t) = p_{ij}(t) \left(F_i(e_i^j, \mathbf{p}_{-i}(t)) - \bar{F}(\mathbf{p}(t)) \right) \quad (6)$$

This means that the Cross-learning model converges to the RE in the continuous time limit, each time interval sees “many” iterations of the game, and that the adjustments which players make between two iterations of the game are “very small”. In this case, a law of large numbers can be applied, and the (stochastic) learning process becomes in the limit deterministic [23]. Based on these results, researchers have made a connection between evolutionary models and various ReL algorithms [24], [25], [26].

3.3. Application to Wireless Networks in 5G

We consider the 5G network topology shown in **Figure 3. 1**. In this scenario, the RUs are installed, managed and operated by coexisting MNOs. The RUs share a set of computational resources that are located both at the edge of the access network (in a MEC server) and at the metro/core network (in the Cloud). The interconnection between the MEC server and the central cloud servers is carried out by an SDN-controlled optical FH/BH transport network. In this environment, SDN controller is responsible to calculate the forwarding rules and install these rules in the network devices, interconnecting the RUs with the MEC and the centralized computational resources. These decisions are taken in a centralized manner.

To reduce the computational complexity, functional split related decisions are taken locally at the RUs. These elements decide where to perform the processing of the low layer functions of the LTE protocol stack. Aligned with the enhanced Common Public Radio Interface (eCPRI) specification, a set $S = \{1,2,3\}$ containing three possible functional split options are considered [9]:

1. Split E (split 1 for simplicity) in which MEC is responsible for the radio frequency (RF) processing of the received signals and the Cloud performs the entire baseband processing.

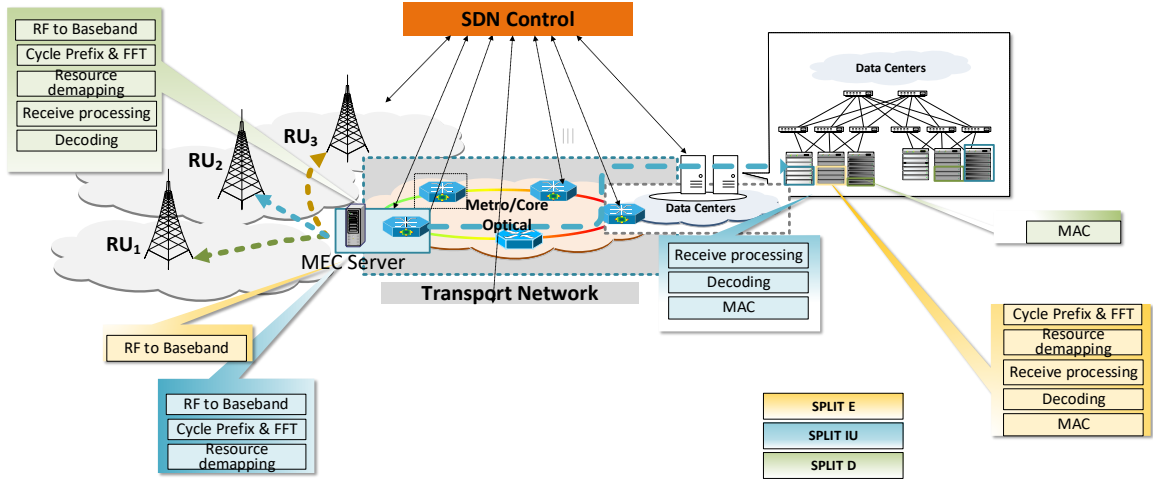


Figure 3. 1:Network architecture. In the MEC, a decision about which functions should be processed locally is made for each RU. The remaining set of functions for each RU are transferred through a common network infrastructure with centralized control to a DC for further processing

2. Split IU (split 2) in which MEC handles the per cell processing (RF processing, cyclic prefix elimination, frequency domain transformation (FFT) and resource demapping), while the remaining functions are performed at the Cloud (Equalization, IDFT, QAM, multi-antenna processing, Forward Error Correction (FEC), higher level operations (MAC, RLC, PDCP) and,
3. Split D (split 3) where the entire lower layer function chain is performed at the MEC server, and the higher lever functions in the Cloud. One can conclude that as the split is placed lower in the 5G protocol stack, the required transport capacity increases [27].

3.3.1. Problem formulation using EGT

Each RU periodically selects one of the three possible functional splits with probability x_i , $i=1,\dots,3$. The decisions are sent to the SDN controller, who is responsible for the application of the policies (i.e. set the appropriate rule forwarding policies at network switches in order establish the necessary connectivity between the RUs, the MEC devices and the central cloud facilities). We consider the scenario in which all the necessary resources are available. When the policies have been applied, the payoffs are calculated and the RUs are reviewing their split option strategy. Specifically, if a better (lower) payoff is observed, then the probability of an RU to select the specific split option increases (decreases). The new policies are sent to the controller and the same procedure is repeated. The time between each repetition is referred to as revision time. To address this scenario, EGT can provide a suitable optimization framework that can be used to support energy-aware FH service provisioning over a common infrastructure.

Let $\mathbf{x}(t) = [x_1(t) \ x_2(t) \ x_3(t)]^T$ be the state vector of the RU, where $x_i(t)$, $i \in S$ refers to the RU's probability of choosing split i . If the RU revise its strategy with a time rate $r_i(\mathbf{x})$, the change of the proportion of the probabilities is described by the following dynamical equation:

$$\dot{x}_i(t) = \sum_{j \in S} x_j(t) r_j(\mathbf{x}(t)) p_j^i(\mathbf{x}(t)) - \sum_{j \in S} x_i(t) r_i(\mathbf{x}(t)) p_i^j(\mathbf{x}(t)) \quad (7)$$

Table 3. 1

NETWORK AND PROCESSING DEMANDS OF EACH FUNCTIONAL SPLIT

Split	Network Rate	Processing functions	
		Local	Remote
1 (E)	R_1 ($N_o \cdot f_s \cdot 2 \cdot N_Q \cdot N_R$)	RF	FFT, RE Demapping, Rx Processing, DEC, MAC
2 (IU)	R_2 ($N_{sc} \cdot T_s^{-1} \cdot 2 \cdot N_Q \cdot N_R \cdot \eta$)	RF, FFT, RE Demapping	Rx Processing, DEC, MAC
3 (D)	R_3 ($N_{sc} \cdot T_s^{-1} \cdot \eta \cdot E_{spectral}$)	RF, FFT, RE Demapping, Rx Processing, DEC	MAC

where $p_i^j(\mathbf{x})$ is the rule of change in the probability of choosing split i when the RU samples split j and can be expressed as:

$$p_i^j(\mathbf{x}(t)) = \begin{cases} x_j(t)(u(j, t) - u(i, t)) & j \neq i \\ 1 - \sum_{j \neq i} x_j(t)(u(j, t) - u(i, t)) & \text{otherwise} \end{cases}, \quad i, j \in S \quad (8)$$

with $u(i, t)$ denoting the payoff of split i at time t . Substituting Eq. (8) to Eq. (7) and making the assumption that all review rates are constantly equal to one ($r_i(\mathbf{x}) \equiv 1 \frac{\text{revision}}{\text{time unit}}$), the following differential equation yields:

$$\dot{x}_i(t) = x_i(t)[u(i, t) - \sum_{j \in S} x_j(t)u(j, t)] \quad (9)$$

which satisfies the replicator dynamics model introduced in Eq. (1).

The objective of each MNO is to minimize their own service power consumption requirements and, hence, the service operational costs. Thus, the payoff function per operator is formed by summing up the power consumption of the network and compute elements required to support FH services. **Table 3. 1** summarizes the network and processing demands of each functional split.

For this problem setting, the payoff of an RU belonging to an MNO that chooses split i against another RU operated by a different MNO who chooses split j is described by the payoff matrix \mathbf{A} , with elements:

$$a_{ij} = -\left(P_{PROCESSING_{ij}} + P_{NET_{ij}}\right) + b, \quad i, j \in S \quad (10)$$

where $P_{PROCESSING}$ and $P_{NET_{ij}}$ refer to the total compute and network energy consumption respectively, when split i competes with split j and b is a positive constant that guarantees the robustness of the system. Technical parameters like the oversampling factor (N_o), the sampling frequency (f_s), the quantization bits per I/Q (N_Q), the number of receiving antennas (N_R), the number of subcarriers used (N_{sc}), the

percentage of used resource elements (η), and the spectral efficiency ($E_{spectral}$) affect the required capacity and the power consumption of each processing function [27], [28].

Due to latencies lays that are introduced in the transport network, the payoff values are provided with some delay. It is evident, that this kind of procedure indicates that the strategies will evolve based on information related to a past moment. This will be reflected to the expected payoff of the strategies. This delay mainly composed of propagation, serialization, switching/routing and queuing delay across the network. Although propagation and switching/routing delays are constant, the rest are highly affected by the network traffic. Given that that network delays can be modeled as random variables following a specific probability density function (PDF) $P(t)$ [29], we expect that RUs will also receive the corresponding payoffs with delay τ following the same PDF. The expected payoff u of an RU using strategy i as well as the average payoff are determined by [30]:

$$u(i, t) = \int_0^\infty P(\tau) (\mathbf{A}\mathbf{x}(t - \tau))_i \quad (11)$$

$$\bar{u} = \sum_{j \in S} x_j(t) u(j, t) \quad (12)$$

Substituting Eq. (11)-(12) in Eq. (9) we get a nonlinear system of differential equations. Since this system cannot be easily solved by analytical methods it is important to examine its qualitative behavior without actually solving it. We concentrate on finding the stability of a solution exploiting the Lyapunov stability theorem. This method is based on the expansion of the right part of the dynamical system as a Taylor series about an equilibrium point \mathbf{x}^0 . If the initial condition $\mathbf{x}(0) = \mathbf{x}_0$ is close enough to \mathbf{x}^0 , then \mathbf{x} will be a small perturbation for some time interval extending from zero. Thus, it is acceptable to neglect the higher-order terms, and approximate the nonlinear system by the linear system [16]:

$$\dot{\mathbf{x}}(t) = \mathbf{J}_0 \mathbf{x}(t) + \mathbf{J}_1 \mathbf{x}(t - \tau), \quad (13)$$

where $\mathbf{J}_0 \in \mathbb{R}^{2 \times 2}$ and $\mathbf{J}_1 \in \mathbb{R}^{2 \times 2}$ are respectively, the Jacobian matrix, and the delayed Jacobian matrix evaluated at equilibrium at \mathbf{x}^0 .

The stability of the system requires that all roots of its characteristic equation have a negative real part. The characteristic equation can be expressed as:

$$\det(\mathbf{I}\lambda - \mathbf{J}_0 - \mathbf{J}_1 Q) = 0 \Rightarrow \lambda^2 + D\lambda + E\lambda Q + FQ^2 + GQ + H = 0 \quad (14)$$

where $\lambda \in \mathbb{C}$, \mathbf{I} is the $N \times N$ identity matrix, $Q = \int_0^\infty P(\tau) e^{-\lambda\tau}$ corresponds to the Laplace transform of the delayed term in Eq. (13) and D, E, F, G, H are parameters that depend on the Jacobian matrices' elements.

The system admits to seven equilibrium points: three corner points, one interior and three corner side points. The linearization about each of the three corner critical points produces an ordinary differential equation that is independent of the delayed variables as in the non-delayed three strategies game.

At the interior critical point all the payoffs are equal. The differential system that emerges depends only on the delayed variables, thus one should anticipate that the

distributed delay will affect its stability. The parameters D, G, H are eliminated and the characteristic equation is formed as:

$$u^2 + E \cdot u + F = 0, u = \frac{\lambda}{Q} \quad (15)$$

The last three critical points are equilibriums where only two of the three strategies survive (corner side points). Their characteristic equation can be written as:

$$(\lambda - l_1) \cdot (\lambda - l_2 Q) = 0 \quad (16)$$

where l_1 and l_2 are parameters that depend on the corner side equilibrium point.

As we can conclude from the above, our analysis can be restricted for finding the solution of the equation

$$\lambda - C \int_0^\infty P(\tau) e^{-\lambda\tau} = 0 \quad (17)$$

where C is a variable that depends on the equilibrium point (e.g. $C = l_2$ if it is a corner side equilibrium point). The above equation is the characteristic equation of the linear differential equation

$$\dot{x}(t) = C \int_0^\infty x(t - \tau) f(t) d\tau \quad (18)$$

Thus, the conclusions derived for the stability of Eq. (13) can be expanded to our case. Based on [31], we derive the following necessary and sufficient condition for the asymptotic stability of the equilibriums:

Proposition 1: If $C < 0$ and the expected value (E) of the delay's probability density satisfies the condition:

$$E(\tau) < \frac{\pi}{\gamma \cdot |C|} \quad (19)$$

where $\gamma = 2$ when the pdf is symmetrical, or else $\gamma = \sup \{ \gamma | \cos w = 1 - \frac{\gamma w}{\pi}, w > 0 \}$, then the equilibrium point is stable [31]. As far as the variance of the distribution is concerned, the stability of the system increases as the variance grows [31].

3.3.2. Problem formulation using MARL

In the previous formulation of the problem the interactions of the RUs have been modelled as symmetric and pairwise random matching with the RUs being identical. However, this scenario only represents a very restricted fraction of the 5G ecosystem cases. In most scenarios, the RUs that interact in each time interval are more than two, and in many cases with different technical parameters. This environment falls in the category of multiagent systems, where the optimal behavior of its agents can be learned through the methods of MARL.

The RUs are characterized by a probability state vector, which indicates how likely he is to play any of his strategies. We assume that at each time interval, each RU has information only about her strategy, the information being the payoff that she receives. The RUs cannot observe the strategy of each other. After obtaining the payoff, the RUs update their probability vector, with the adjustment formulas given in Eq (4)-(5). This procedure can be described by the multipopulation replicator equation (Eq (6)) in the continuous time limit[23].

The fact that each RU can only observe her strategy and respective payoff means that the information about the effect of the joint action of all RUs on the environment is given through the received payoff value. In order to construct the appropriate payoff function for this purpose, we introduce a map $K: \mathfrak{R}^{1 \times 3^{N-1}} \times \mathfrak{R} \rightarrow \mathfrak{R}$ that takes a vector matrix \mathbf{A} and a summation as arguments and returns a new summation. The mathematical formula that describes the mapping is:

$$K(\mathbf{A}, \sum_{m=1}^M a_m) = \sum_{m=1}^M A_m a_m \quad (20)$$

The mapping is valid only when the size of the vector \mathbf{A} is the same as the number of the summed terms M .

For this problem setting, the payoff value of the i^{th} RU who chooses split j is formed as:

$$u_{ij} = K(A^j, \prod_{n \neq i}^{N-1} \sum_{s=1}^3 x_{ns}) \quad (21)$$

where N is the number of the RUs that are interacting, A^j is a vector of size 3^{N-1} that describes the payoff of every possible combination of the other RUs' strategies, assuming the $i - th$ RU plays the $j - th$ strategy and x_{ns} is the probability of the $n - th$ RU to play the $j - th$ strategy. The payoff value per RU is formed as in the previous model, namely by summing up the power consumption of the network and compute elements required to support FH services (shown in **Table 3. 1**).

Substituting Eq. (21) in Eq. (6) we get a nonlinear system of differential equations, that is the dynamics in the continuous limit of the learning process. The stability analysis of the system can facilitate the tuning of the learning parameters (number of iterations, time interval θ).

3.4. Results and Discussion

In the previous section, we discussed a model for finding the optimal split option, adopting at first a simple model based on EGT, and then a more sophisticated scheme using MARL. The parameters that define the optimal split depend on the processing and network power requirements. The network power consumption depends on the required connection rate between the antenna and the remote processing platform, the evolution of traffic growth over time and network technology improvements. On the other hand, the processing power depends on the technical parameters described in **Table 3. 1**. In this section, we will first examine through a numerical the position of the SDN controller that guarantees the stability of the system using EGT. Based on the result, we will tune the parameters of our learning model, in order to include the case of asymmetrical, multi-interaction system.

3.4.1. SDN Controller placement using EGT

As it was mentioned earlier, the SDN controller is responsible for collecting and providing to the MNOs of the RUs the required information from all controlled devices. The maximum delay corresponds to the delay of the most distant node to the controller path plus the delay of the controller-MEC path. Thus, assuming that each controlled device may host a MEC, the stability of the system is achieved only when the round-trip time (RTT) of the controller's path to the most remote device is less than the limit

Table 3. 2

PARAMETERS OF THE SYSTEM CONFIGURATION

Symbol	Quantity	Value
B	bandwidth	20 MHz
Ant	number of the rx antennas	2
M	modulation	6 bits/symbol
R	coding rate	1/2
dt	time-domain duty-cycling	100%
f_s	sampling frequency	30.72 MHz
N_o	oversampling factor	2
N_{sc}	number of used subcarriers	1200
T_s	symbol duration	66.6 μ s
N_Q	quantization bits per I/Q	10
S	spectral efficiency	3 bit/cu
η	assumed RB utilization	70%

imposed by Eq (19). Based on this limit, we propose a heuristic algorithm that tries to identify the minimum number and associated position of SDN controllers with the aim to guarantee the stability of the 5G infrastructure. This is performed with low computational complexity.

At first, the heuristic algorithm finds the maximum network radius, that is the number of hops of the longest end-to-end path. Then, for each node it calculates the maximum RTT to all the other nodes inside the network radius. If the result of all nodes is a number higher than the limit imposed in Eq (19), the network radius is reduced by one, and the same procedure is repeated, until a case is found where the RTT from a node to all other nodes within the network radius meets the condition of Eq (19). The nodes that meet this requirement, are marked as possible controller candidates. From this set, the algorithm chooses as the first controller the one that is connected to the largest number of devices within the network radius. These devices and the first controller are removed from the network, and the whole procedure is repeated for the downscaled network. The algorithm ends when the downscaled network has no network nodes.

In order to see the effectiveness of the proposed EGT-model we consider the system depicted in **Figure 3. 1**, with the parameters shown in **Table 3. 2**. We assumed that large scale DCs provide superior performance per Watt, compared to cloudlets, so the cost of remote processing is higher than the cost of the central one. The cost ratio (remote/central processing) was assumed to be equal to two. Furthermore, the relationship of the transport network's energy consumption with the required capacity for the support of the FH services was assumed to be nonlinear, since the non-linear

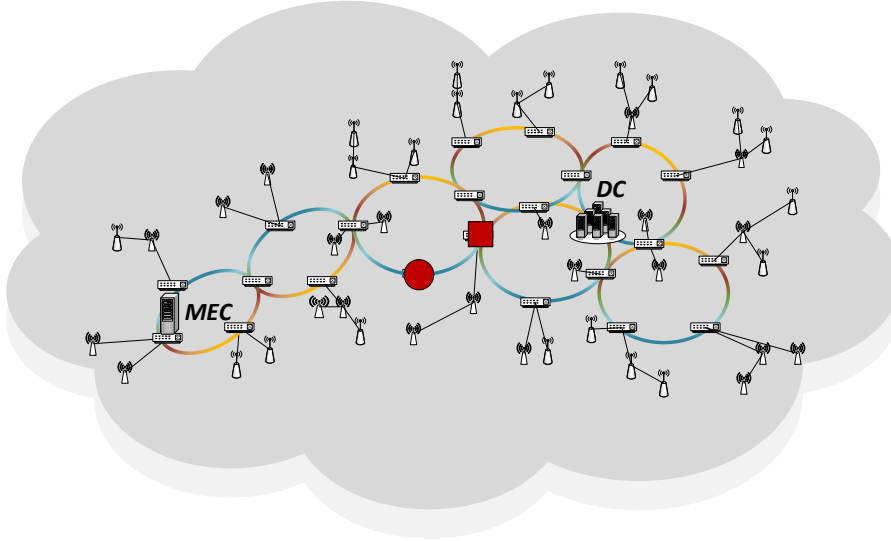


Figure 3. 2: Assumed FH/BH transport network for the system described in **Figure 3. 1**. The red circle represents the position of the SDN controller, after the implementation of the heuristic algorithm described in section IV.A.. The red square represents the optimal position estimated according to the average propagation latency-case described in [11].

model is best to describe the technology advancements in terms of energy efficiency of network devices [32].

The stability analysis of system (9) indicates that the equilibrium point in such a scenario is $x_1^* = 0.2957$, $x_2^* = 0.7043$, $x_3^* = 0$. This means that in the non-delayed system the optimal split choice is split 2. However, as it was stated previously the SDN transport network introduces additional delay to the system. This delay can be divided to two main components, namely the processing delay of the SDN controller and the transport delay.

The SDN controller chosen for the implementation is the Opendaylight controller (ODL), that is a scalable controller infrastructure that supports SDN implementations in modern heterogeneous networks of different vendors [33]. For measuring the processing delay of the ODL controller, we developed an application that communicates externally with the controller. For evaluations, a linear network topology with Out of Band control plane was emulated in Mininet, a tool that can emulate and perform the functions of network devices in a single physical host or VM [34]. Both Mininet and Opendaylight controller were implemented on the same machine (Intel® Core™ i5-7400U CPU @ 3.00GHZ (4 cores)) to overcome the Ethernet interface speed limitations. 7.7 GiB of memory was available. The system was running Ubuntu 16.04 LTS-64 bit. The application implements at first step a mechanism for collecting data on the network topology and at second step a mechanism for sending echo messages to all switches simultaneously, and measuring the time elapsed for receiving a reply. The time response of ODL is measured by averaging the delay results of a large number of tests, in order to achieve higher accuracy. The results showed an exponential relation between the controller’s processing delay and the network devices.

Regarding transport delay, we used monthly delay measurements extracted from GRNET [35], in order to find the dependence of the E2E transport delay on the E2E hops. Our analysis concluded that this relationship can be well approximated with a linear function. Furthermore, the best pdf that fitted the end to end delay was the generalized t-student distribution [36].

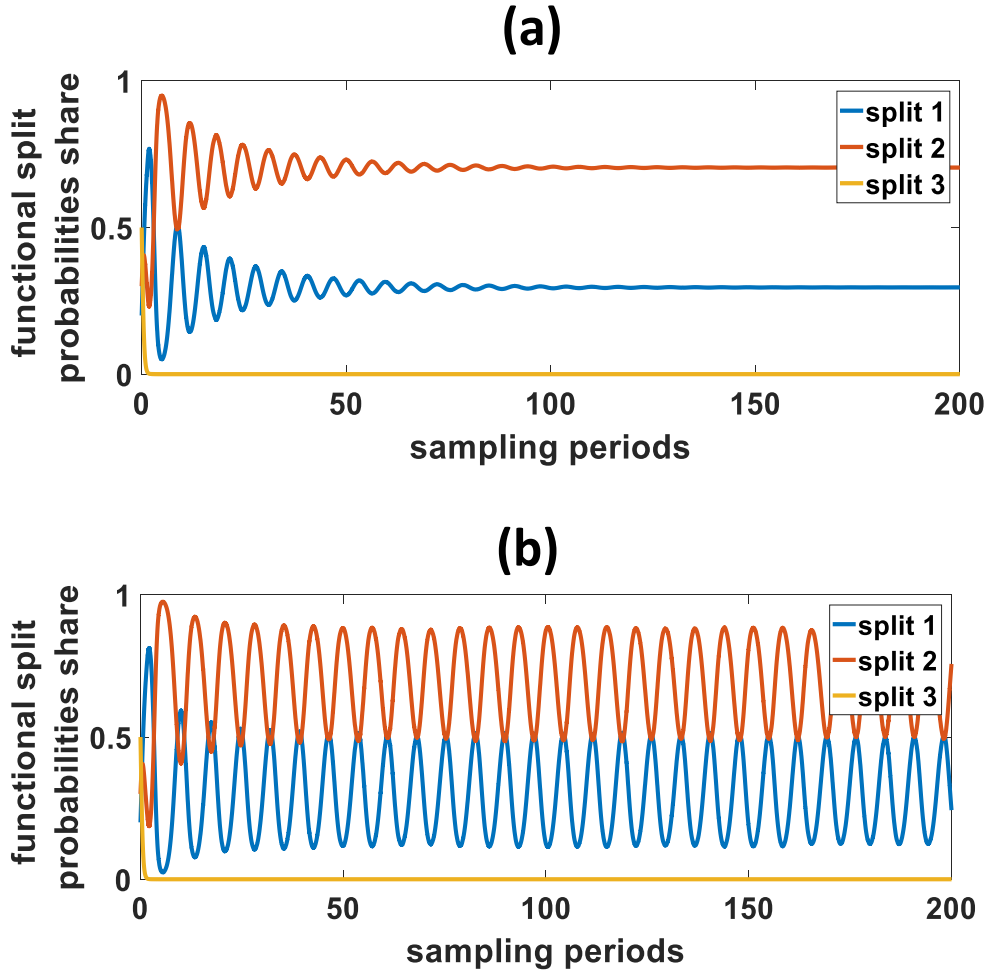


Figure 3. 3: Evolution of the probabilities of the three split options, with the parameters described in Table 3. 2, when: (a) the controller is placed in the proposed location (red circle in Figure 3. 2) by the heuristic, b) the controller is placed in the proposed location (red square in Figure 3. 2) of the average propagation latency-case described in [11]

Taking these into consideration, we expect that the induced SDN-transport network's delay will be a random variable that is characterized by the generalized t-student distribution, with expected value that depends on the size of the transport network and the hops between two network nodes. Thus, the upper delay limit for our example is given by Eq (19) as: $E_{max} = 1.6449$ time units.

The assumed FH/BH transport network's topology for our example is depicted in **Figure 3. 2**. The figure also shows the possible controller placements after implementing the heuristic algorithm described in the previous section. In order to test the validity of the heuristic, we investigate the evolution of strategies in two cases: 1) when the controller is placed in one of the proposed locations by the heuristic, 2) when the controller is placed in the location identified by the average propagation latency optimization technique described in [11]. **Figure 3. 3** illustrates the evolution of split option selection probability under the proposed EGT based approach and the average latency minimization scheme described above. As can be seen in the former case (**Figure 3. 3. (a)**) after few sampling periods the scheme converges to a mixed solution where all antennas operate under a single split option mode that will be either split 1 or

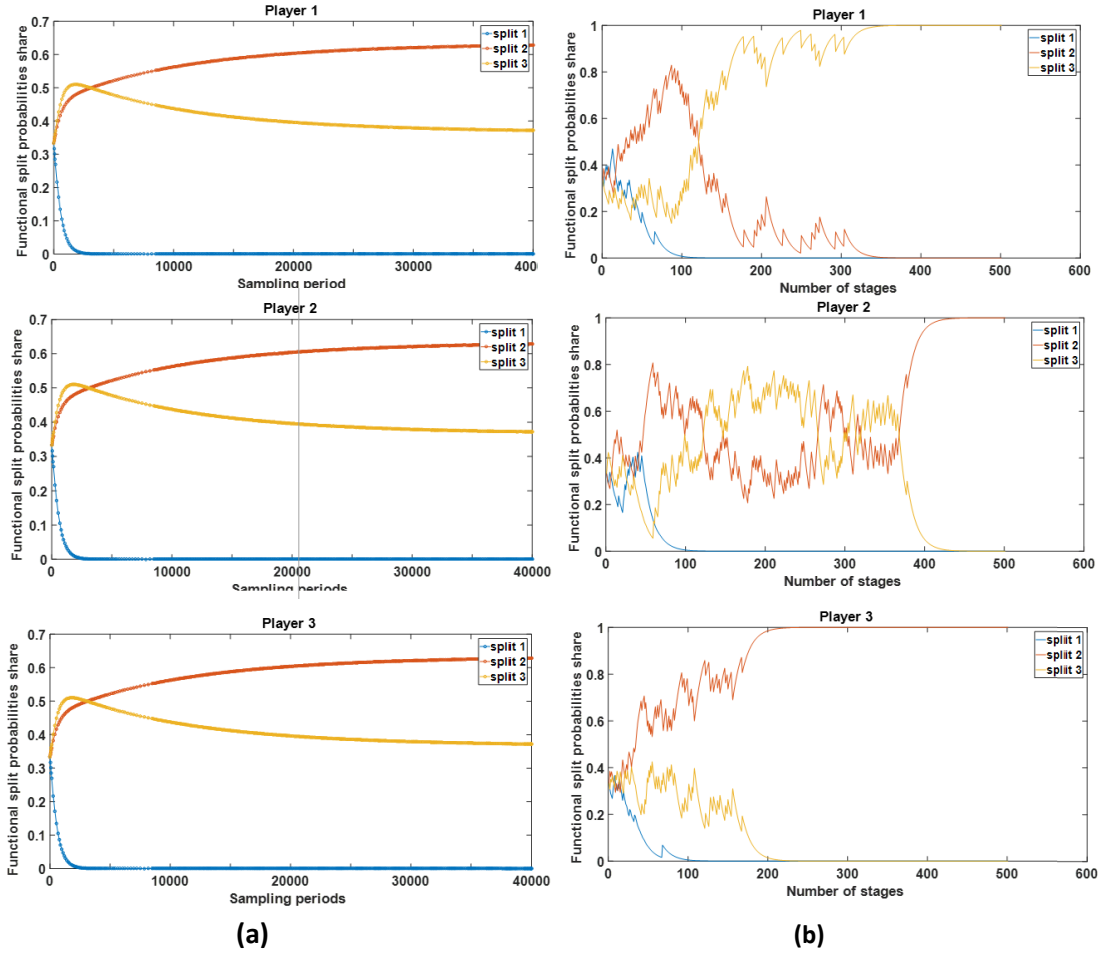


Figure 3.4: The evolution of the probabilities of the three split option for three RUs, with the parameters described in Table 3. 2 using (a) the multipopulation replicator equation model, (b) the cross learning algorithm, with 600 stages and $\theta = 9$.

split 2. However, in the second case (**Figure 3.3 (b)**), the placement of the SDN controller at a node that does not satisfy the stability threshold imposed by equation (19) leading to an unstable operational mode for the 5G network. The reason behind this is that the increased control plane delay in this case introduces inaccurate information of the network status at the controller. Therefore, decision making is performed with outdated information that leads to an oscillation around the optimal operating point preventing it from converging to a stable solution.

3.4.2. Optimal Split Selection using MARL

In this section we will consider a more sophisticated scenario, where the RUs don't imitate each other's strategy, but are trying to optimize their behavior based on the impact of their action in the environment. Since more than one RU are trying to optimize their behavior in the same time, the problem falls to the category of MARL.

First, we will consider the simple case of three interacting RUs. The RUs are identical and characterized by a probability distribution over their set strategy. After they execute an action, they receive the analogous payoff and update their probability vector according to the Cross-Learning algorithm. The received payoff is the total power consumption of the network and compute elements required to support FH services. As

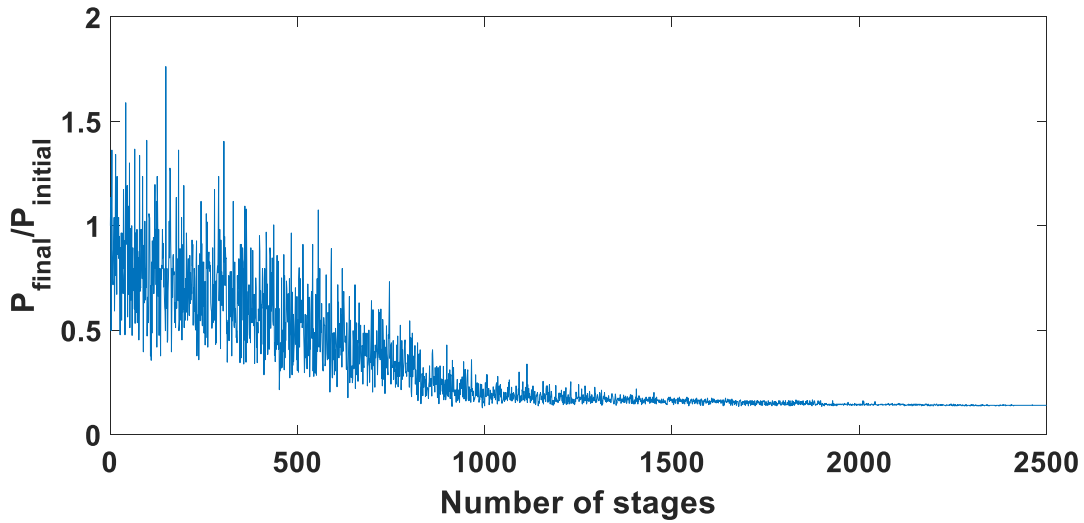


Figure 3. 5: The total power consumption of the FH services in relation with the stages of the cross learning algorithm

it was mentioned earlier the multipopulation replicator equation can describe the behavior of the system in the continuous time limit. In **Figure 3. 4** we can see the evolution of the probabilities of the three RUs, using the RD and a simple cross learning algorithm. The RD model shows that the splits that have a probability to occur are split 2 and 3, a fact that is confirmed in the cross-learning model. For the cross learning model, the parameters were chosen as $n = 600$ stages and $\theta = 9$ time units between each stage, after many trials, with the objective to achieve the minimization of the total energy consumption.

The next step is to consider the scenario of **Figure 3. 2**. The controller is placed at the position indicated in section IV-A. The system consists of 43 identical RUs in terms of technical parameters, that update their probability distribution vector according to the Cross-Learning model. When more than three RUs are interacting the stability analysis of the multipopulation RD dynamics becomes rapidly unfeasible. In this situation, we can tune the parameters of the learning model by combining the delay limit with the minimization of the total power consumption of the FH services. It should be noted that the RUs are identical for simplicity purposes, however the same procedure can be applied for RUs with different technical parameters. The learning algorithm shows that the minimization of the total energy occurs when 30% of the RUs operate at split option 2 and 70% at split option 3. **Figure 3. 5** shows the power consumption of the FH services with respect to the algorithm stages, where it is clear that the targeted goal is achieved. Taking into consideration that the time interval between the stages is equal to the delay limit measured in milliseconds (calculated in the previous section), the minimization of the power consumption is achieved after 4 seconds.

3.5. Summary

To address the limitations of current RANs, centralized-RANs adopting the concept of flexible splits of the BBU functions between RUs and the CU have been proposed. Further efficiency gains in terms of cost and energy consumption are achieved through the implementation of this architectural model exploiting compute resources, required

for the BBU function processing, located both at the MEC and relatively large-scale centralized DCs. This architecture adopts high bandwidth/low latency SDN controlled optical transport networks. In this scenario, and with the aim to dynamically identify the optimal split option that minimizes infrastructure operational costs, in terms of power consumption, we have proposed a novel mathematical model based on EGT. In addition, optimal placement of the transport network SDN controllers is determined by a heuristic algorithm with the objective to guarantee the stability of the whole system. Finally, for more complex scenarios of asymmetrical, multi-interaction systems, we used the relation of a basic learning algorithm with EGT to identify the optimal split option for the RUs that minimizes the total FH operational costs.

References

- [1] Tzanakaki et al., "Wireless-Optical Network Convergence: Enabling the 5G Architecture to Support Operational and End-User Services", *IEEE Communications Magazine*, vol. 55, no. 10, pp. 184-192, 2017. Available: 10.1109/mcom.2017.1600643.
- [2] "5G and Verticals < 5G-PPP", *5g-ppp.eu*, 2019. [Online]. Available: <https://5g-ppp.eu/verticals/>. [Accessed: 30- Oct- 2019].
- [3] "5G Network Slicing for Vertical Industries", *Huawei.com*, 2019. [Online]. Available: <https://www.huawei.com/minisite/5g/img/5g-network-slicing-for-vertical-industries-en.pdf>. [Accessed: 30- Oct- 2019].
- [4] M. Kamel, W. Hamouda and A. Youssef, "Ultra-Dense Networks: A Survey", *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, pp. 2522-2545, 2016. Available: 10.1109/comst.2016.2571730.
- [5] "View on 5G Architecture", *5g-ppp.eu*, 2019. [Online]. Available: <https://5g-ppp.eu/wp-content/uploads/2018/01/5G-PPP-5G-Architecture-White-Paper-Jan-2018-v2.0.pdf>. [Accessed: 30- Oct- 2019].
- [6] H. Liu, F. Eldarrat, H. Alqahtani, A. Reznik, X. de Foy and Y. Zhang, "Mobile Edge Cloud System: Architectures, Challenges, and Approaches", *IEEE Systems Journal*, vol. 12, no. 3, pp. 2495-2508, 2018. Available: 10.1109/jsyst.2017.2654119.
- [7] "Cloud RAN and MEC: A Perfect Pairing", *Etsi.org*, 2019. [Online]. Available: https://www.etsi.org/images/files/ETSIWhitePapers/etsi_wp23_MEC_and_CRAN_e_d1_FINAL.pdf. [Accessed: 30- Oct- 2019].
- [8] "What is enhanced Mobile Broadband (eMBB)", *5g.co.uk*, 2019. [Online]. Available: <https://5g.co.uk/guides/what-is-enhanced-mobile-broadband-emb/>. [Accessed: 30- Oct- 2019].
- [9] "eCPRI Specification V1.1", *Cpri.info*, 2019. [Online]. Available: http://www.cpri.info/downloads/eCPRI_v_1_1_2018_01_10.pdf. [Accessed: 30- Oct- 2019].
- [10] "Software-Defined Networking (SDN) Definition - Open Networking Foundation", *Open Networking Foundation*, 2019. [Online]. Available: <https://www.opennetworking.org/sdn-definition/>. [Accessed: 30- Oct- 2019].
- [11] B. Heller, R. Sherwood and N. McKeown, "The controller placement problem", *ACM SIGCOMM Computer Communication Review*, vol. 42, no. 4, p. 473, 2012. Available: 10.1145/2377677.2377767.

- [12] D. Hock et.al., "Pareto-Optimal Resilient Controller Placement in SDN-based Core Networks", *IEEE Proceedings of the 2013 25th International Teletraffic Congress (ITC)*, 10-12 Sept. 2013, Shanghai, China. Available: 10.1109/ITC.2013.6662939
- [13] M. Noormohammadpour and C. S. Raghavendra, "Datacenter Traffic Control: Understanding Techniques and Tradeoffs," in *IEEE Communications Surveys & Tutorials*, vol. 20, no. 2, pp. 1492-1525, Second quarter 2018
- [14] J. Weibull, *Evolutionary game theory*. Cambridge, Mass.: MIT Press, 2004.
- [15] T. Yi and W. Zuwang, "Effect of Time Delay and Evolutionarily Stable Strategy", *Journal of Theoretical Biology*, vol. 187, no. 1, pp. 111-116, 1997. Available: 10.1006/jtbi.1997.0427.
- [16] G. Obando, J. Poveda and N. Quijano, "Replicator dynamics under perturbations and time delays", *Mathematics of Control, Signals, and Systems*, vol. 28, no. 3, 2016. Available: 10.1007/s00498-016-0170-9.
- [17] N. Anastasopoulos, D. Asteriou, "Optimal dynamic auditing based on game theory." *Oper Res Int J*, 2019. <https://doi.org/10.1007/s12351-019-00491-3>
- [18] N. P. Anastasopoulos, M.P., Anastasopoulos, "The evolutionary dynamics of audit. *Eur. J. Oper. Res.* pp. 469–476, 2012.
- [19] D. Bloembergen, K. Tuyls, D. Hennes and M. Kaisers, "Evolutionary Dynamics of Multi-Agent Learning: A Survey", *Journal of Artificial Intelligence Research*, vol. 53, pp. 659-697, 2015. Available: 10.1613/jair.4818.
- [20] M. Wiering and M. Otterlo, *Reinforcement Learning*. Berlin: Springer Berlin, 2014.
- [21] P. Hernandez-Leal, M. Kaisers, T. Baarslag and E. Munoz de Cote, "A Survey of Learning in Multiagent Environments: Dealing with Non-Stationarity", *arXiv*, vol. 170709183, 2017. [Accessed 5 November 2019].
- [22] L. Busoniu, R. Babuska and B. De Schutter, "A Comprehensive Survey of Multiagent Reinforcement Learning," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 2, pp. 156-172, March 2008. doi: 10.1109/TSMCC.2007.913919
- [23] T. Börgers and R. Sarin, "Learning Through Reinforcement and Replicator Dynamics", *Journal of Economic Theory*, vol. 77, no. 1, pp. 1-14, 1997. Available: 10.1006/jeth.1997.2319.
- [24] Karl Tuyls, P.J 't Hoen, and B. Vanschoenwinkel, "An Evolutionary Dynamical Analysis of Multi-Agent Learning in Iterated Games", *Journal of Autonomous Agents and Multi Agent Systems*, Vol. 12, No. 1, pp. 115–153, 2006.
- [25] L. Panait, K. Tuyls, and S. Luke, "Theoretical Advantages of Lenient Learners: An Evolutionary Game Theoretic Perspective", *Journal of Machine Learning Research*, Vol. 9, pp. 423–457, 2008.
- [26] T. Klos, G.J.V Ahee and K. Tuyls, "Evolutionary Dynamics of Regret Minimization", *Technical report*, 2010.
- [27] Wübben et.al., "Benefits and Impact of Cloud Computing on 5G Signal Processing", *IEEE Signal Processing Magazine*, pp. 35-44, November 2014.
- [28] C. Desset et.al., "Flexible power modeling of LTE base stations," *IEEE Wireless Communications and Networking Conference (WCNC)*, Shanghai, China, April, 2012.

- [29] Y. Xia and D. Tse, "Inference of Link Delay in Communication Networks," in *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 12, pp. 2235-2248, Dec. 2006
- [30] N. Ben Khalifa et.al. "Random time delays in evolutionary game dynamics", in *Proceedings of IEEE CDC*, Osaka, Japan, pp 3840–3845.
- [31] Samuel Bernard et.al.. "Sufficient conditions for stability of linear differential equations with distributed delay", *Discrete & Continuous Dynamical Systems - B*, 2001, 1 (2): 233-256. doi: 10.3934/dcdsb.2001.1.233
- [32] Jayant Baliga et.al, "Energy Consumption in Optical IP Networks," *J. Lightwave Technol.*, vol. 27, pp. 2391-2403, 2009.
- [33] Platform Overview. <https://www.opendaylight.org>.
- [34] Mininet Overview. <http://mininet.org/overview/>
- [35] <https://grnet.gr/infrastructure/network-and-topology/>
- [36] Student's t-distribution, https://en.wikipedia.org/wiki/Student%27s_t-distribution

Chapter 4

Dynamic User Plane Function Allocation in 5G Networks

Contents

4.1. Chapter Introduction	50
4.2. Problem Formulation based on theoretical EGT	52
4.2.1. System Model	52
4.2.2. Application in 5G networks	54
4.2.3. Numerical Results and Discussion.....	56
4.3. Problem Formulation based on Lab Measurements	56
4.3.1. System Model	56
4.3.2. Results and Discussion	58
4.4. Summary	60
References	60

4.1. Chapter Introduction

5G communication systems rely on an open and flexible network paradigm to address the requirements of both telecom operators and vertical stakeholders in a cost and energy efficient manner. To this end, the concepts of hardware programmability and network softwarisation are adopted to develop suitable interfaces that can be used to a) interconnect a variety of wired and wireless network technologies forming a common transport network (TN) and, b) decouple Control and User plane functionalities. The former facilitates implementation of a variety of 5G-RAN deployment options, while CUPS allows flexibility in network deployment and operation as well as cost efficient traffic management.

A big part of the user plane functionality in 5G systems is handled by the UPF, which has to be designed to support challenging 5G services with very tight performance requirements. It connects with external IP networks hiding mobility related aspects from the external networks. Moreover, it performs different types of processing of the forwarded data, such as packet inspection, redirection of traffic and application of different data rate

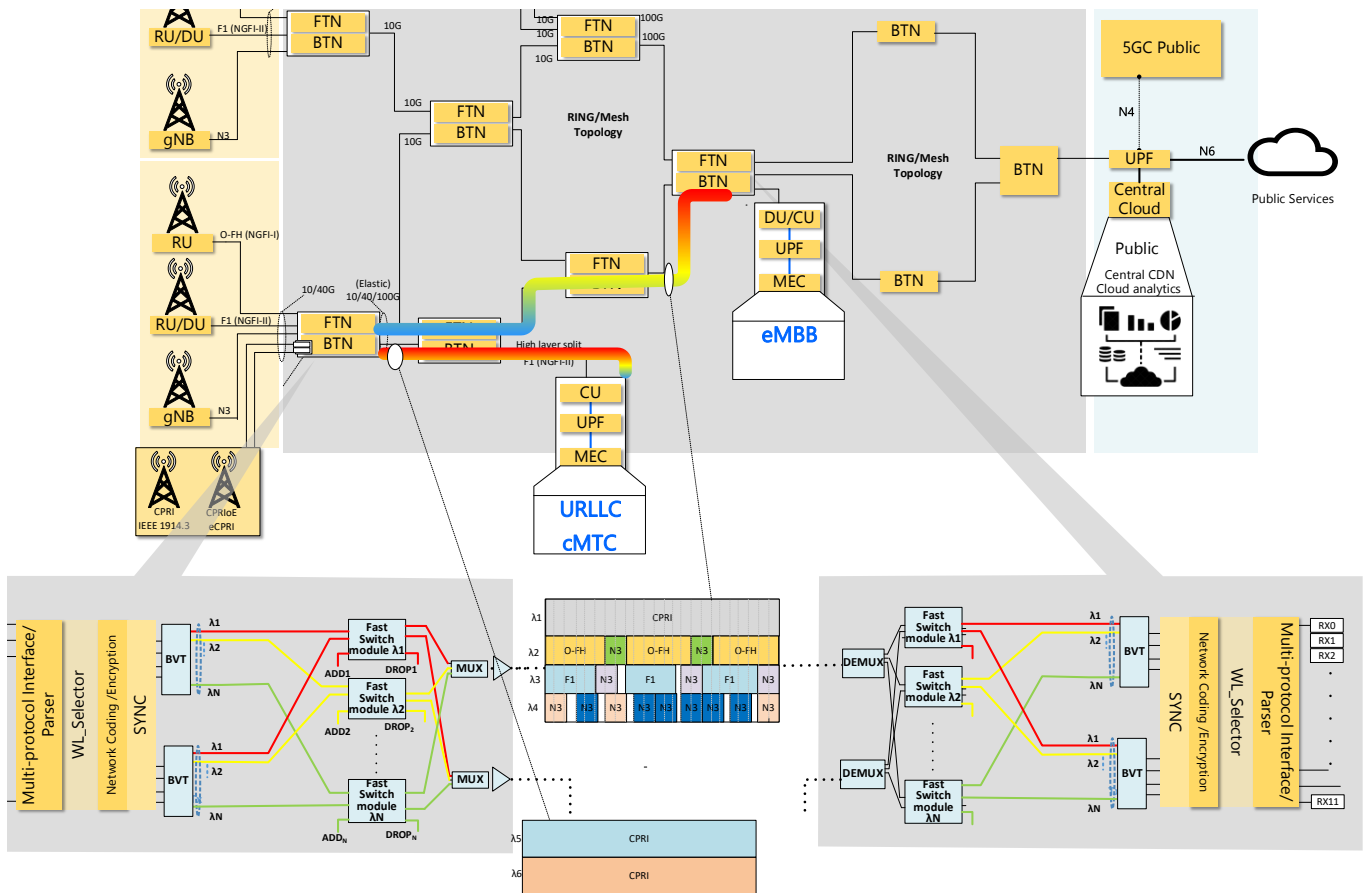


Figure 4. 1: 5G System architecture for access to two (e.g. local and central) data networks [3]

limitations. 5G-CUPS supporting multiple UPFs enables 5G edge capabilities; one of the key 5G advancements compared to 4G. UPF related processing can be dynamically deployed and configured depending on the applications' needs. Overall, UPFs act as termination points for various interfaces and protocols and are also responsible to take several actions (rules) [1] including: (a) Mapping of traffic to the appropriate tunnels based on the QFI information. This requires UPFs to be able to perform Deep Packet Inspection and identify the necessary values in the General Packet Radio Service Tunnelling Protocol (GTP-U) header, associate QFIs with the appropriate Differentiated Services Code Point (DSCP) codes in the external IP network and perform the relevant protocol adaptations at line rate. (b) Steering of packets to the appropriate output port and take the necessary packet forwarding actions. (c) Packet counting for charging and policy control purposes. (d) Deep packet inspection for security and anomaly detection purposes. (e) Buffering and queuing management for traffic service differentiation and assurance of end-to-end delays.

To perform these actions UPFs should support an extensive set of protocols such as, GTP-U, PFCP (Packet Forward Control Packet), IP and also assist in the operation of SDAP and PDCP through mapping of DSCP classified IP traffic coming from the external DN. It should be also capable of handling legacy and new protocols such as enhanced CPRI/Open RAN and Radio over Ethernet (RoE) at high-rates. Towards this direction, *programmable optical networks (e.g. Time Shared Optical Network-TSON [5]) can be effectively used to support transport network requirements as well as classify and*

steer traffic. This is performed adopting specific interfaces for control plane (N1/2, N4), user plane (N3, N6) and UPF handover (N9) communication requirements. For example, the Network Interface Cards (NICs) can steer control plane protocols packets such as PFCP packets into the SMF or the control plane part of UPF and can steer UE sessions based on the PDU session, the flow, the QoS class etc. through N3 and N6 interfaces. Programming can be also used to support extended header (EH) for 5G user plane traffic.

A high-level view of a 5G deployment option is shown in **Figure 4. 1**. The concept of disaggregated 5G-RAN places the RU, DU, and CU entities at different sites (a MEC or a CC) along the transport infrastructure, depending on the split options and network scale. Taking into consideration the need for a common network infrastructure to carry legacy and 5G traffic, a transport network with nodes that can handle both FH and BH traffic is necessary. These FTN/BTN nodes use both wired and wireless technologies including optical, electrical, mmWave, and THz, with capacities ranging from 1Gbps to 100Gbps and key features like synchronization. They can operate in a complementary manner, allowing for a wide range of transport network connectivity options such as point-to-point, point to multi-point, and multi-point to multipoint. Low cost FTN/BTN nodes can be placed close to RUs to connect radio transmission points while aggregating transport traffic from various cell sites to a central location (hub site) for centralized processing. FTN/BTN nodes with higher capacity and more sophisticated connectivity can collect traffic from hub sites and aggregate it into edge sites, and further aggregate it through the core network responsible for collecting traffic from aggregation networks. The mobile core entities can be located closer to the end-user for low-latency services like URLLC or cMTC.

In this environment multiple UPFs are supported. Based on the 5G-RAN deployment option and the type of service that needs to be provided, UPF nodes can be placed closer or further away from the 5G-RAN. Through this approach, UPF elements placed close to the network edge can redirect traffic to MEC servers reducing latency, whereas UPF nodes placed deeper into the network can send traffic to central cloud facilities. In this context, as the network dimension grows, a larger number of rules is required to support policies, whereas network resources (e.g., switch memory) are limited. Therefore, if the available UPF computational resources are occupied with existing services any new upcoming service required by the UE will be blocked. This may result in increased service delay as the number of flows requiring UPF processing increase.

To address this problem, we propose a novel scheme based on EGT that allows dynamic selection of the optimal UPF elements. So far, a very limited set of studies exist addressing the problem of UPF selection [2]. Section 4.2 formulates the optimal UPF selection problem considering a specific optical node implementation [5] and using accurate modeling of the delays introduced when this programmable optical node is adopted to act as UPF element as shown in **Figure 4. 1**. Section 4.3 extends the analysis by considering realistic cost functions in the EGT model developed that have been calculated using lab measurements derived from an open source 5G platform hosted in an optical datacenter cloud environment.

4.2. Problem Formulation based on theoretical EGT

4.2.1. System Model

We consider the uplink transmission of a 5G network shown in **Figure 4. 1**. The UEs initiate the PDU Session Establishment process by transmitting the relevant request to the

AMF. The AMF contacts the SMF, which in turn checks whether the UE requests are compliant with the user subscription. Once subscription information has been verified the SMF selects a UPF to serve the PDU Session. This is a key decision to be taken as a UPF at close proximity to the RAN, may be the optimal choice at first sight, since it should result in reduced latency. However, if all UEs are associated with this UPF congestion may arise resulting in increased latency. To address this challenge, we propose a scheme that allows dynamic selection of the UPFs by the UEs. In this approach users try to optimize their own performance selfishly. The choice adaptation process of the UEs can be formulated as an evolutionary game. To formulate this problem, we consider a set of UEs each requesting a service of class $g \in G$ where G is the total number of available service classes. Let also $S^g = \{UPF_1^g, \dots, UPF_{N_g}^g\}$ be the set of available strategies in users belonging to g -group. For each group, each UE tunnel needs to be terminated at a specific UPF. Assuming that N_g denotes the available UPFs for group g , then the population of the UEs in group g can be described at each time instance by vector $\mathbf{x}^g(t) = [x_1^g(t) \dots x_{N_g}^g(t)]$ where $x_i^g(t)$ is the proportion of UEs in group g that are currently being served by UPF_i . Each UE belonging to a specific group remains associated with a UPF for a time interval, and reviews its choice periodically. When a revision occurs, the UE switches from UPF_i to another UPF_j according to a switching probability $p_{ij}^g(\mathbf{x}) = x_j^g$, that is equal to the population probability distribution of strategies, where $\mathbf{x} = [x^1(t) \dots x^g(t)]$, is the population state of the system. If a switch occurs, the UE receives a payoff $u_j^g(\mathbf{x})$ that quantifies its satisfaction level associated with the selection of UPF_j . The obtained payoff affects the arrival rate of the revision opportunities. Assuming that the number of reviews of a UE that uses strategy i can be described by a Poisson process with arrival rate $r_i^g(\mathbf{x})$, and all UEs' Poisson processes are statistically independent, we can use the law of large numbers to approximate the adaptation process with the following deterministic dynamic model [4]:

$$\dot{x}_i^g(t) = \sum_{j \in S^g} x_j^g(t) r_j^g(\mathbf{x}) p_{ji}^g(\mathbf{x}) - x_i^g(t) r_i^g(\mathbf{x}) \quad (1)$$

The UE updates its review rate, by linearly decreasing it to its current payoff. This means that the average review rate of a UE that uses strategy i is:

$$r_i^g(\mathbf{x}) = a - \beta u_i^g(\mathbf{x}), \quad \beta > 0 \text{ and } \frac{a}{\beta} > u_j^g(\mathbf{x}) \quad (2)$$

This results in forcing UEs with higher payoffs to revise their UPF choice at lower rates than the rest, leading to the replicator dynamics:

$$\dot{x}_i^g(t) = \beta \left(u_i^g(\mathbf{x}) - \bar{u}^g(\mathbf{x}) \right) x_i^g \quad (3)$$

According to this equation, a selected strategy will either survive or be eliminated in the long run depending on whether its payoff is better or worse than the average payoff of all strategies. Since the objective of the UEs is to optimize their performance in terms of latency, greater payoffs correspond to lower delays.

The observed latency can be decomposed into two main components. The first component is the propagation delay between the UE and the UPF and is proportional to the distance between the two entities. Assuming an underlying optical transport network, the propagation delay due to the propagation time in the fiber links corresponds to 5 μ s per kilometer (km) of fiber. The second component is the delay of processing inside the UPF and can be modeled by adding the processing and the transmission delay, that are constant, and the variable queuing delay. Mechanisms for bounding the processing delay

within a network node can be found both in literature and in standardization. In this analysis, we assumed that the UPF, uses the bounded mechanisms described in [3].

Considering these we formulate the payoff on a user of group g that selects action i , when the population state is $\mathbf{x}(t)$, as

$$u_i^g(\mathbf{x}) = 1/t_{prop_i^g} + t_{UPF_i}(\mathbf{x}) \quad (4)$$

where t_{prop} is the propagation delay and $t_{UPF_i}(\mathbf{x})$ the UPF_i delay that can be approximated by an exponential function [5]:

$$t_{UPF_i}(\mathbf{x}) = e^{k_i \rho_{UE} \sum_{g=1}^G M_g x_i^g} \quad (5)$$

where ρ_{UE} is the traffic of one UE, M_g is the UE-population of group g and k_i is a variable related with UPF_i and depends on the characteristics of the UPF node implementation (**Figure 4. 1**) including data rate, number of ports (fibres, wavelengths), buffering capability etc.

4.2.2. Application in 5G networks

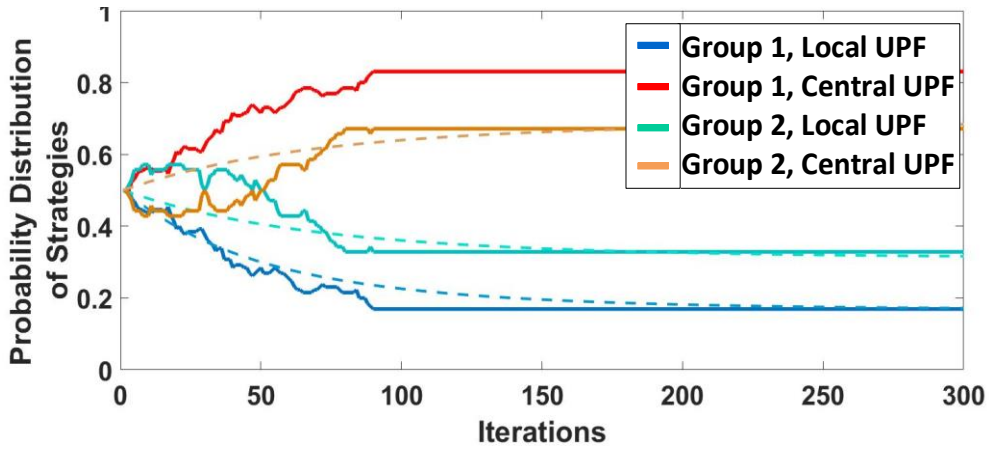
Based on the replicator dynamics of the EGT, we developed a scheme to attain the evolutionary equilibrium. The following steps summarize the algorithm:

(1) *Initialization*: Every UE in each group chooses a strategy at random and observes its payoff u . Then it calculates its review rate λ according to the formula $\lambda = \alpha - \beta u$, where α, β are constants.

(2) *Revision*: A revision opportunity may occur to each UE with probability equal to $p_{revision} = \lambda \cdot dt$, where dt is the time interval between two loops. If the revision occurs, the UE chooses to imitate at random one of the UEs of its group. Then it recalculates λ according to the obtained payoff. The same process is applied until the difference of each strategy's payoff compared with the average payoff of the population is lower than a limit ε .

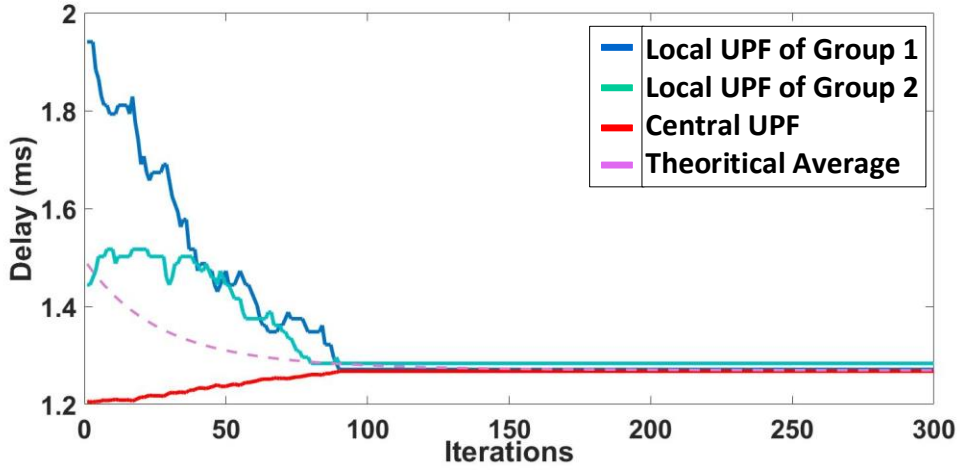
Note that the strategy adaptation process in the proposed EGT-based algorithm does not rely on the knowledge of the strategy selection of the other players. For the evolution a UE requires a random matching with an opponent, a function that can be offered by a central controller (e.g. the SMF). Therefore, the amount of information exchange is reduced. The central controller will randomly match the UEs and stop the evolution process, if all payoffs are equal or differ by a small quantity.

The time interval (dt) between two repetitions must be higher than the communication time between the UE, the AMF, the SMF and the UPF that is going to carry the PDU session. dt is highly affected by the number of UPFs that are under the control of the SMF, since a large number of UPFs may result in increased processing delay for the SMF. Taking into consideration the timing requirements of the network service ($t_{service}$), and the



Full lines: simulation results, Dotted lines: theoretical results

a)



b)

Figure 4. 2: (a) Trajectories of proportions of population and (b) convergence of the algorithm to the equilibrium (for $M_1 = 130, M_2 = 70, \frac{a}{b} = 1, \frac{k_{mec}}{k_{cc}} = 10$). In the equilibrium 16% of group1 UEs and 32% of group 2 UEs are served by their local UPFs, while the remaining are served by the central UPF.

number of iterations of the algorithm (L), the number of UPFs (N) under the SMF's control can be evaluated so that the following relationship is true:

$$N < F^{-1} \left(\frac{t_{service}}{L} \right) \quad (6)$$

Where F^{-1} is the inverse function that relates dt with N .

4.2.3. Numerical Results and Discussion

This section presents simulation results to validate our theoretical findings and evaluate the proposed algorithm performance. In the following we assumed a population of UEs that are organized into two groups. The UEs in each group can decide whether they want to use a local UPF at the edge of the network, that connects to a MEC, or to a UPF four times further away ($t_{prop_{cc}}/t_{prop_{mec}} = 4$), that connects to a central cloud as shown in **Figure 4. 1**. The UPF in the central cloud can process a greater number of requests, compared to the local UPFs, and is shared by all groups in the UE population whereas the local UPF is dedicated to the population inside a group. The traffic generated by each UE is assumed to be $\rho_{UE} = 100 \text{ Mbps}$. The limit ε of the algorithm is set to a payoff difference of 0.01. **Figure 4. 2** illustrates our simulation results (full lines) and the theoretical results derived through the model of the replicator dynamics (dotted lines) demonstrating good agreement between theory and simulation. More specifically, **Figure 4. 2 (a)** plots the evolution of strategy shares among the population of UEs. It can be observed that the system converges after some iterations to the equilibrium. In equilibrium all UEs achieve the same delay (**Figure 4. 2 (b)**) indicating the fairness of the scheme. The number of total iterations of the algorithm is of vital importance for network planning. As it was discussed in IIB, the number of iterations in combination with the time requirements of the service, can give an estimate (Eq. (6)) of the number of UPFs that the SMF can control, without compromising the stability of the system. **Figure 4. 2** shows that less than 100 iterations are needed for the system to converge.

4.3. Problem Formulation based on Lab Measurements

4.3.1. System Model

We consider a 5G system as shown in **Figure 4. 1** where a set of UEs initiate a PDU Session Establishment process in order to get access to computational resources hosted either at a MEC or at a CC facility. The relevant session is established through the interaction of the AMF and the SMF that assigns a UPF to serve the corresponding PDU Session. In this study, this decision is taken in a distributed manner adopting an evolutionary optimization scheme.

Let a set of N UEs each requesting access to a 5G network. Their demands can be fulfilled either through the establishment of a PDU session connecting the UE with the local MEC (denoted as MEC strategy) or through a PDU session between the UE and the CC (denoted as CC strategy). UEs, depending on the selected UPF/MEC node, receive a payoff that quantifies its satisfaction level associated with that choice. The payoff depends on the compute and optical network resources used for the establishment of these connections. Thus, the payoff of a UE is given as the sum of the total compute and network resources used to carry its request. For a UE connected to the server through the UPF node s , the payoff is given by:

$$\pi_s(i) = \frac{h_c + h_n}{h_c \cdot CPU_s(i) + h_n \cdot NET_s(i)} \quad (7)$$

In (7), i is the number of UEs served by UPF s , CPU_s is compute cost for UPF s and NET_s is the optical network cost for the interconnection of the UE with node $s \in$

$\{UPF_{mec}, UPF_{cc}\}$, h_c and h_n are weighting factors which can be determined according on which resource (cpu or network) should be given advantage.

The distributed PDU session establishment problem is solved using the following evolutionary process.

At each time step, a random UE is chosen for updating its strategy from the entire population. The UE measures its payoff and imitates with some probability an opponent, that is a randomly chosen UE from the rest of the population. This probability depends on the payoff difference of the two UEs and is given by the Fermi distribution [7]. According to this, a UE will imitate the strategy of its opponent with probability higher than 50% if its payoff is lower than the payoff of its opponent.

At each time step, i changes by 1 (at most) according to the following probabilities:

$$P_{i,i\pm 1} = \frac{i}{N} \frac{N-i}{N} \frac{1}{1 + e^{\mp w(\pi_{mec}(i) - \pi_{cc}(i))}} \quad (8)$$

$$P_{i,i} = 1 - P_{i,i+1} - P_{i,i-1} \quad (9)$$

where w governs, how much influence has the payoff to the adaptation process.

The overall system can now be considered at each step as a Markov chain with i , $i \in \{0, 1, \dots, N\}$ being the state of the population. Solving this system, the probability to reach the state where all UEs are redirected to the local MEC, given that the initial state is k , can be written as [7][8]:

$$x_k = \frac{1 + \sum_{l=1}^{k-1} \prod_{m=1}^l \gamma_m}{1 + \sum_{l=1}^{N-1} \prod_{m=1}^l \gamma_m}, \text{ with } x_1 = \frac{1}{1 + \sum_{l=1}^{N-1} \prod_{m=1}^l \gamma_m} \quad (10)$$

where $\gamma_i = \frac{P_{i,i-1}}{P_{i,i+1}} = e^{-w(\pi_{mec}(i) - \pi_{cc}(i))}$.

This probability reveals significant information on the evolutionary stability of the system. If we denote with $\rho_{mec} = x_1$, $\rho_{cc} = 1 - x_{N-1}$ the probability of one initial UE assigned to the local MEC and one UE to the CC, respectively, and $h_i = \pi_{mec}(i) - \pi_{cc}(i)$ the payoff difference of the two strategies at state i , then there are four possible outcomes [9][10]:

1. *Domination of MEC-strategy*

$$\rho_{cc} < \frac{1}{N} < \rho_{mec}, h_1 > 0, h_{N-1} > 0 \quad (11)$$

There is no interior equilibrium point.

2. *Domination of CC-strategy*

$$\rho_{cc} > \frac{1}{N} > \rho_{mec}, h_1 < 0, h_{N-1} < 0 \quad (12)$$

There is no interior equilibrium point.

3. *Coexistence*

$$\rho_{mec} > \frac{1}{N}, \rho_{cc} > \frac{1}{N}, h_1 > 0, h_{N-1} < 0 \quad (13)$$

Table 4. 1

CPU UTILIZATION FOR 4 VMs WITH DIFFERENT CHARACTERISTICS IN TERMS OF CPU, MEMORY AND DISK

UEs	CPU Utilization			
	Small	Medium	Large	XLarge
	1vcores, 2GB, 20GB	2vcores, 4GB,40GB	4vcores, 8GB, 80GB	8vcores, 16GB, 160GB
1	16.50%	10.24%	7.65%	4.00%
2	19.50%	11.88%	8.88%	4.79%
3	24.00%	14.20%	10.82%	4.94%
4	29.50%	18.58%	11.28%	5.91%
5	40.10%	22.35%	12.95%	8.20%
6	46.80%	24.48%	16.17%	8.90%
7	55.70%	25.20%	17.40%	9.05%
8	60.60%	26.92%	18.30%	9.14%
9	70.15%	36.10%	20.80%	10.70%
10	75.70%	39.72%	22.66%	11.17%

The interior equilibrium point of replicator dynamics is stable.

4. Bi-stability:

$$\rho_{mec} < \frac{1}{N}, \rho_{cc} < \frac{1}{N}, h_1 < 0, h_{N-1} > 0 \quad (14)$$

The interior equilibrium point of replicator dynamics is unstable.

A comparison of ρ_{mec} and ρ_{cc} is also of interest in favor of which strategy the process spends more time. The system spends more time in whichever strategy's corresponding ρ is greater (because that strategy needs less invasion attempts to fixate).

4.3.2. Results and Discussion

To evaluate the evolutionary UPF selection problem over an optical transport network we deployed Free5GC, an open source 5G Core platform, over a virtualized optical data center network. The relevant measurements have been used to quantify processing and network resources used by the 5G network to support the required UE connectivity. Measurements have been recorded for 10 UEs requesting connectivity whereas the virtualized 5G system is hosted in VMs with different allocated resources (i.e namely small, medium, large and xlarge). The relevant CPU results are shown in **Table 4. 1**.

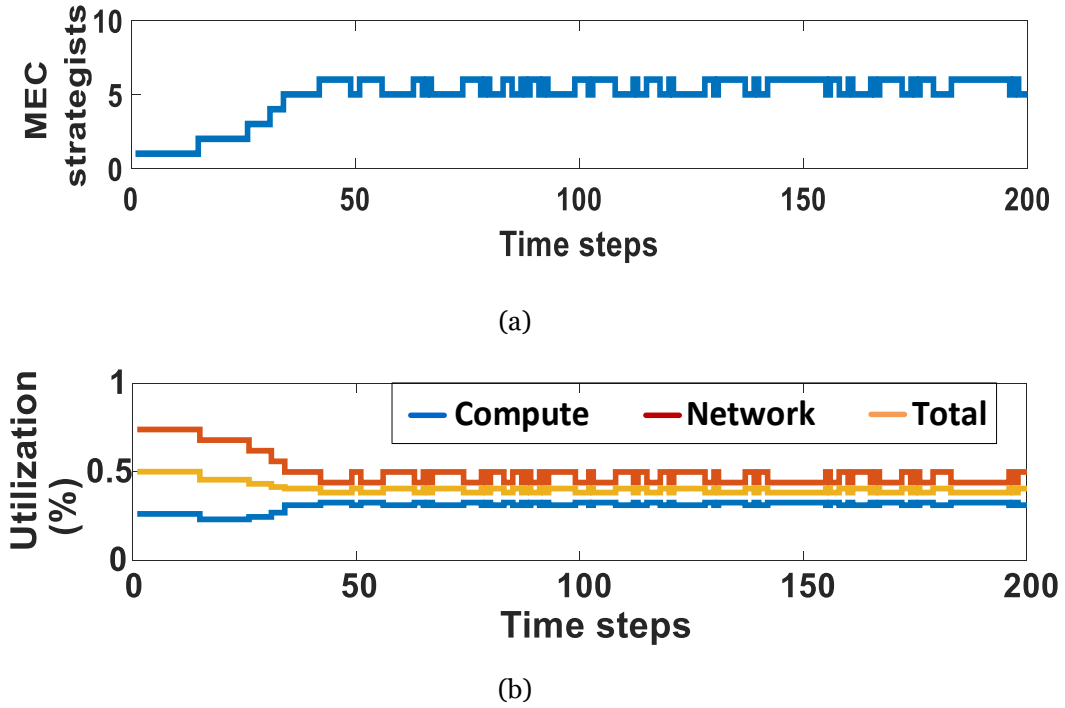


Figure 4. 3: (a) Evolution of UEs that choose MEC (b) Resource Utilization during the evolution of the UE population.

The network utilization of the PDU sessions, depends on the requested services as well as the number of links (distance) between the UEs and the selected UPF. For each link, utilization is defined as the ratio of the PDU traffic over the link capacity. Therefore, the network utilization is calculated as the weighted average of the link utilization.

Figure 4. 3 shows the results of the proposed scheme for the case of a system with a UPF directing traffic to a local MEC, and a UPF connected to the CC. The local MEC is lightweight (hosting a small VM) whereas the CC has higher processing capabilities (hosting a medium VM). The local MEC is 1 hop away from the UE whereas the remote CC is located four hops further from the edge. From Eq (7), (9) we calculate $\rho_{mec} = 0.3385$, $\rho_{cc} = 0.1740$, $h_1 = 2.0418$ and $h_{N-1} = -1.9669$. This means that if the population of UEs was infinite, the system would have a stable interior equilibrium where both strategies would gain the same payoff, so they would co-exist (Eq (13)). In the case of only 10 UEs the system will oscillate around this equilibrium point. The fact that $\rho_{mec} > \rho_{ccc}$ indicates that in the interior equilibrium, the number of MEC-UEs exceeds the one of CC-UEs. Indeed, if we calculate the payoff difference h_i for all the possible states of the system we find that the state that minimizes this difference is when 6 UEs choose MEC and 4 CC. It is worth mentioning that the number of iterations (L) it took the algorithm to converge has a great impact in the network design. Taking into consideration the timing requirements of the network service, the larger the number of iterations the smaller the time interval between two repetitions of the algorithm must be and thus the lower the processing time of the SMF. This can give an estimate of the number of UPFs that can be assigned under the control of a single SMF without compromising the stability of the system. As shown in **Figure 4. 3 (a)** the system converges (oscillates around the equilibrium) relatively fast ($L = 40$ time steps). **Figure 4. 3 (b)** depicts the system trade-off between compute and network costs aiming to minimize the total consumption of resources. **Figure 4. 4** depicts the dependence of the interior equilibrium from the system parameters. Specifically in **Figure 4. 4 (a)** we can see how the processing capabilities of the CC affect the interior

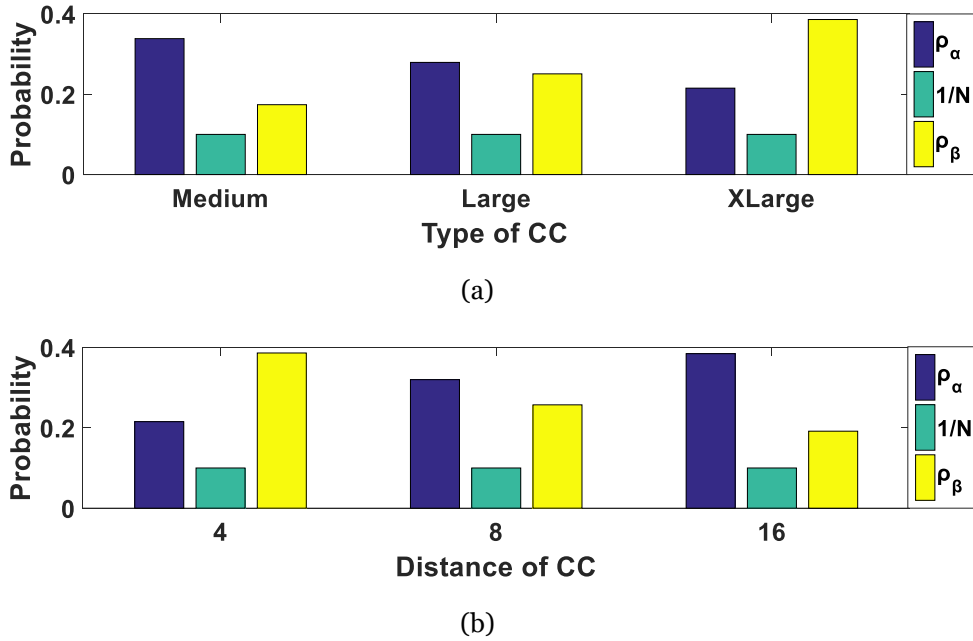


Figure 4. 4: Fixation probabilities for three CC with (a) different processing capabilities (b) different distance from the UEs.

equilibrium. As the remote processing capabilities increase, the interior equilibrium is pushed towards the CC strategy that in this case is more favorable than the MEC strategy (the fixation probability of the MEC strategy decreases, and that of the CC strategy increases). The opposite happens if we increase the network requirements of remote processing by increasing the distance of the CC (**Figure 4. 4 (b)**).

4.4. Summary

This chapter studies optimal UPF placement in 5G networks supported by optical transport networks and a UPF solution that adopts a programmable optical node implementation. A purposely developed EGT model performs dynamic selection of the UPF processing minimizing the overall service delay. The EGT model is capable of selecting the optimal location of the UPF processing unit based on various factors, including network traffic, service requirements, and available resources. To ensure that the model accurately reflects the performance of real-world network environments, the cost functions in the EGT model developed have been calculated using lab measurements derived from an open source 5G platform hosted in an optical datacenter cloud environment.

References

- [1] System architecture for the 5G System (5GS), (3GPP TS 23.501 version 16.6.0 Release 16)
- [2] I. Leyva-Pupo et al., "Dynamic Scheduling and Optimal Reconfiguration of UPF Placement in 5G Networks", In Proc. of MSWiM '20, 2020 NY, USA, 103–111
- [3] H2020 Project 5G-COMPLETE, Deliverable D2.1, "Initial report on 5G-COMPLETE network architecture, interfaces and supported functions". [Online]

- [4] J. W. Weibull, 1997. "Evolutionary Game Theory", MIT Press Books, The MIT Press, edition 1, volume 1, number 0262731215.
- [5] A. Tzanakaki et al., "Wireless-Optical Network Convergence: Enabling the 5G Architecture to Support Operational and End-User Services", *IEEE Comms. Mag.*, vol. 55, no. 10, pp. 184-192, 2017.
- [6] System architecture for the 5G System (5GS), (3GPP TS 23.501 version 16.6.0 Release 16)
- [7] A Traulsen, M. A. Nowak, J. M. Pacheco, "Stochastic dynamics of invasion and fixation", *Phys Rev E Stat Nonlin Soft Matter Phys*, Jul. 2006 ;74(1 Pt 1):011909. doi: 10.1103/PhysRevE.74.011909. Epub 2006 Jul 17. PMID: 16907129; PMCID: PMC2904085.
- [8] L. J. S. Allen, *An introduction to stochastic processes with applications to biology*, Upper Saddle River, N.J., Pearson/Prentice Hall, 2003.
- [9] T. Antal, I. Scheuring, "Fixation of Strategies for an Evolutionary Game in Finite Populations", *Bull. Math. Biol.*, vol. 68, pp. 1923–1944, 2006.
- [10] C. Taylor , D. Fudenberg, A. Sasaki, M. Nowak, "Evolutionary Game Dynamics in Finite Populations", *Bulletin of mathematical biology* vol. 66, pp. 1621-44, 2004.

Chapter 5

Distributed Service Provisioning for 6G Network Infrastructures

Contents

5.1.	Chapter Introduction	63
5.2.	Service Slicing in 5G Systems	65
5.2.1.	QoS Architecture	65
5.2.2.	5G Data Collection and Analytics	68
5.2.3.	Problem Statement.....	69
5.3.	Related Work.....	70
5.4.	Theoretical Background: Evolutionary Game Theory	72
5.5.	Problem Formulation.....	73
5.5.1.	Game Formulation	73
5.5.2.	Payoff Function.....	75
5.5.3.	Dynamics of Adaptation Process	76
5.5.4.	Stability Analysis	77
5.5.5.	From infinite to finite population.....	77
5.6.	5G System Profiling.....	79
5.6.1.	5G Platform Description	80
5.6.2.	Evaluation Process	81
5.6.3.	Cost and Charging Functions.....	83
5.7.	Numerical Results	84
5.7.1.	Simulation with Experimental Values.....	84
5.7.2.	Extended simulation with multiple AFs.....	88
5.8.	Summary.....	89
Appendix		91
A.	Derivation of Fokker-Planck Equation	91
B.	Derivation of Replicator Equation	92
C.	Calculation of π_{HPxt} and π_{LPxt} derivates.....	92
References		93

5.1. Chapter Introduction

An important challenge that needs to be addressed by 6G systems is associated with the need to concurrently support a large variety of very demanding services with very different QoS requirements, such as high bandwidth, low latency, agility, increased resilience and security. All these services should be provisioned over a common infrastructure that is being flexibly and efficiently shared. Sharing of the underlying physical infrastructure is achieved using the concepts of hardware programmability and network softwarisation [1]. The former involves the development of appropriate interfaces that can be used to integrate a variety of wired and wireless networking solutions in support of the required services. The latter enables migration from the notion of network elements to NFs which are interconnected in accordance to the SBA [1]. The softwarisation of 5G systems covers all network domains and planes including the CN, the RAN, the UE and the AF [2]. This approach enables much greater deployment and operational flexibility compared to previous system generations as NFs can now be dynamically instantiated either to a central or a MEC facility [3]. This new concept opens up the possibility to support a wide range of services spanning from services with reduced end-to-end latency and transport network capacity requirements to services requiring to reach to a massive list of recipients [4].

To differentiate service provisioning, 5G systems introduced a new QoS architectural model based on QoS Flows [4]. According to this scheme, applications running at the UEs can generate traffic that can be mapped to flows and treated according to their QoS profile. Differentiated traffic flow management at the application level guarantees that high-priority services maintain session connectivity. In 6G systems, the QoS assignment of end-to-end services having different QoS requirements is facilitated through the concept of network slicing as defined by [5] and [13]. This relies on the partition of traffic to multiple logical networks that are all executed on and share a common physical infrastructure [6]. These logical networks are independent and offer distinct Service Qualities, known as SLAs. Through these the system can provide optimizations specific to different services or operators.

A key enabler of network slicing is the separation of the control and user plane functionality. CUPS allows flexibility in network deployment and operation, as well as cost efficient traffic management. CP functions are handled by the 5G core whereas UP functions are carried out, to a large extent, by the UPF. The UPF connects the mobile system with external IP networks hiding mobility related aspects from the external networks. Overall, UPF acts as a termination point for various interfaces and protocols and is responsible to take several actions (in the form of applied rules) including packet inspection, redirection of traffic and application of different data rate limitations [1]. UPFs can be also dynamically deployed and configured depending on the applications' needs.

CUPS with the support of multiple UPFs, provides a suitable framework to enable “the edge” concept; one of the key 5G enhancements compared to 4G. A high-level view of a

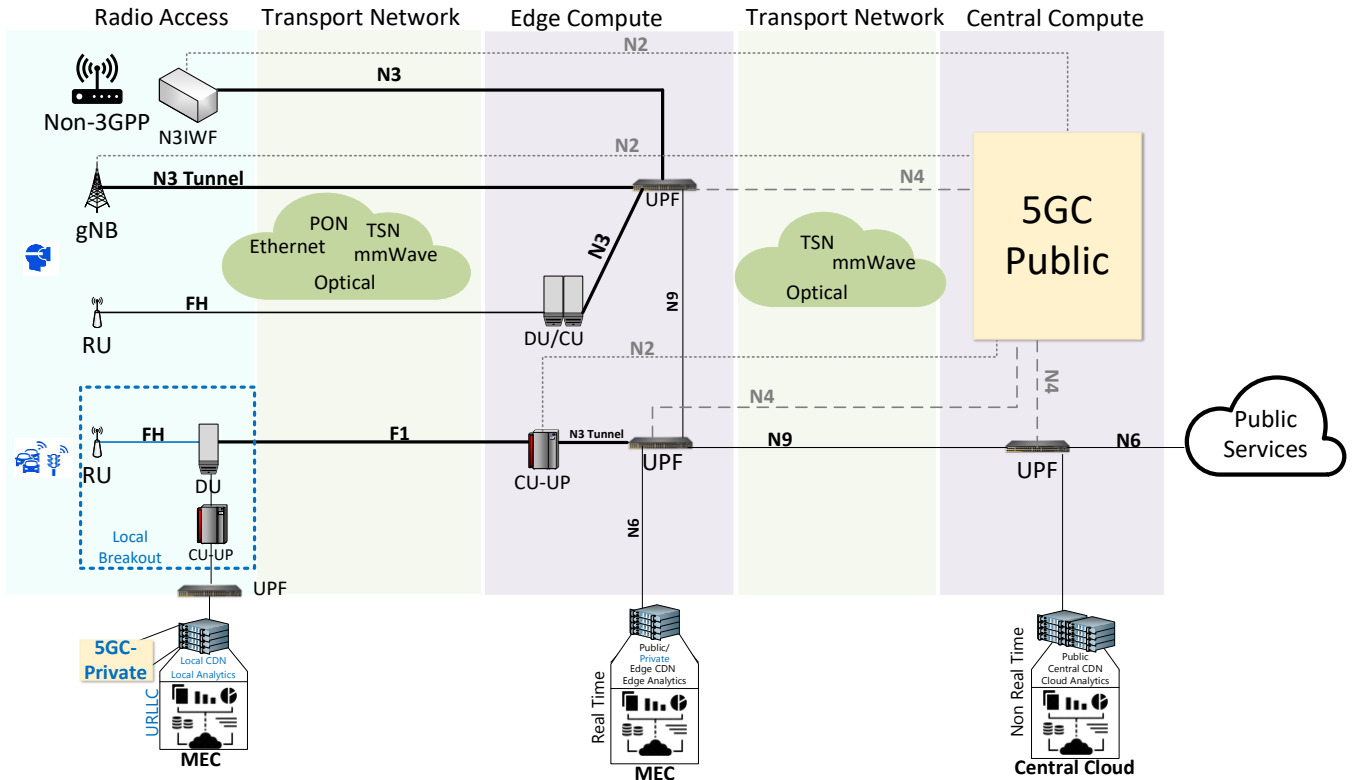


Figure 5. 1: 5G System architecture.

5G deployment option adopting multiple UPFs is shown in Figure 5. 1. This figure, also illustrates the NG-RAN architecture. The 5G-Base Station, that is referred to as gNB, can be split in building blocks, comprising the RU that is responsible for the physical layer functions of the 5G protocol stack, the DU that handles real time layer 1 (L1) and 2 (L2) scheduling functions, and the CU that performs the non-real time processing functions. Those units can be either collocated or instantiated at different compute facilities i.e., at the MEC or the CC nodes. The UP elements can be also placed close to the network edge redirecting traffic to MEC nodes or deeper into the network forwarding traffic to a CC facility. The communication of the UP elements with the 5G core is achieved through their point-to-point interfaces (N2, N3, N4, N6, N9). The 5G Transport Network is a combination of wireless/optical technologies.

6G systems are expected to operate in fully distributed manner allowing applications to directly intervene in the decision-making process. This process has been envisioned by the 3GPP standardization body to be performed by the AF controlling the application running at the UEs and the Application Server (AS) supporting the operation of the service [8]. The role of the AF in the provisioning of services with high QoE requirements is critical as based on the feedback that it receives from the application it can: influence the traffic routing and steering decisions, select the DN where the corresponding AF is hosted, trigger rate adaptation, publish statistics to the analytics function etc. However, allowing AFs to operate in an uncontrolled manner may cause stability issues to the system.

To address this problem, we propose a novel scheme based on EGT that allows AFs to control UE applications in order to dynamically select their service strategy. Thus,

facilitating QoE optimization minimizing at the same time charging costs. The scheme has been designed using realistic constraints and cost functions and has been tested over an operational cloud-based 5G testbed.

The chapter's contributions are summarized as follows:

1. We design, implement and evaluate, both theoretically and experimentally, a fully distributed decision-making framework solving the flow assignment problem in 6G systems. It is proven that under specific conditions the scheme always converge to a stable point providing the optimal tradeoff between QoE and cost efficiency.
2. The cost functions that are used by the EGT model incorporate both network and compute costs. These costs are realistically derived adopting a detailed profiling process over an operational 5G testbed. This profiling process enables modeling of the system performance and requirements under different operational scenarios that can practically assist in the optimized lifecycle management of the provisioned services.
3. We consider practical constraints imposed by the existing protocols used to manage 5G systems.

The rest of the chapter is organized as follows. In Section 5.2 we discuss the 5G network architecture and identify the problem under investigation. Section 5.3 is dedicated to present an overview of published work that falls in the area of focus. The necessary theoretical background is examined in Section 5.4. In section 5.5 the problem is formulated and analyzed in terms of the EGT theory. The experimental setup and the 5G platform profiling results are presented in Section 5.6. Finally, conclusions are drawn in Section 5.7.

5.2. Service Slicing in 5G Systems

5.2.1. QoS Architecture

One of the key internal differences of 5G compared to 4G is a new QoS architecture based on PDU Sessions. This enables the establishment of multiple 'QoS flows' between UEs and the terminating DN connected to the 5G system through the UPF. QoS assurance in 5G and beyond systems is achieved through slicing at transport network, user plane and control plane layers as shown in Figure 5. 2.

At the user plane, slicing between UEs and the DN is performed through two consecutive tunnels: a PDU session between the UE and UPF and an N6 Point-to-Point tunnel between the UPF and the DN. The PDU session consists of a set of SDFs, that, as it was mentioned earlier, are application specific end-to-end packet flows [9]. For example, one PDU session may carry the traffic for a web browsing application together with an email application. Depending on the desired QoS characteristics, an SDF is associated with a 5G QoS Flow. According to the 3GPP specifications [2], [10], the QoS Flow is the finest granularity of QoS differentiation in the PDU Session. The QoS Flows are identified by the QFI. Flows with the same QFI within a PDU Session receive the same traffic forwarding treatment (e.g. scheduling, admission threshold). A typical

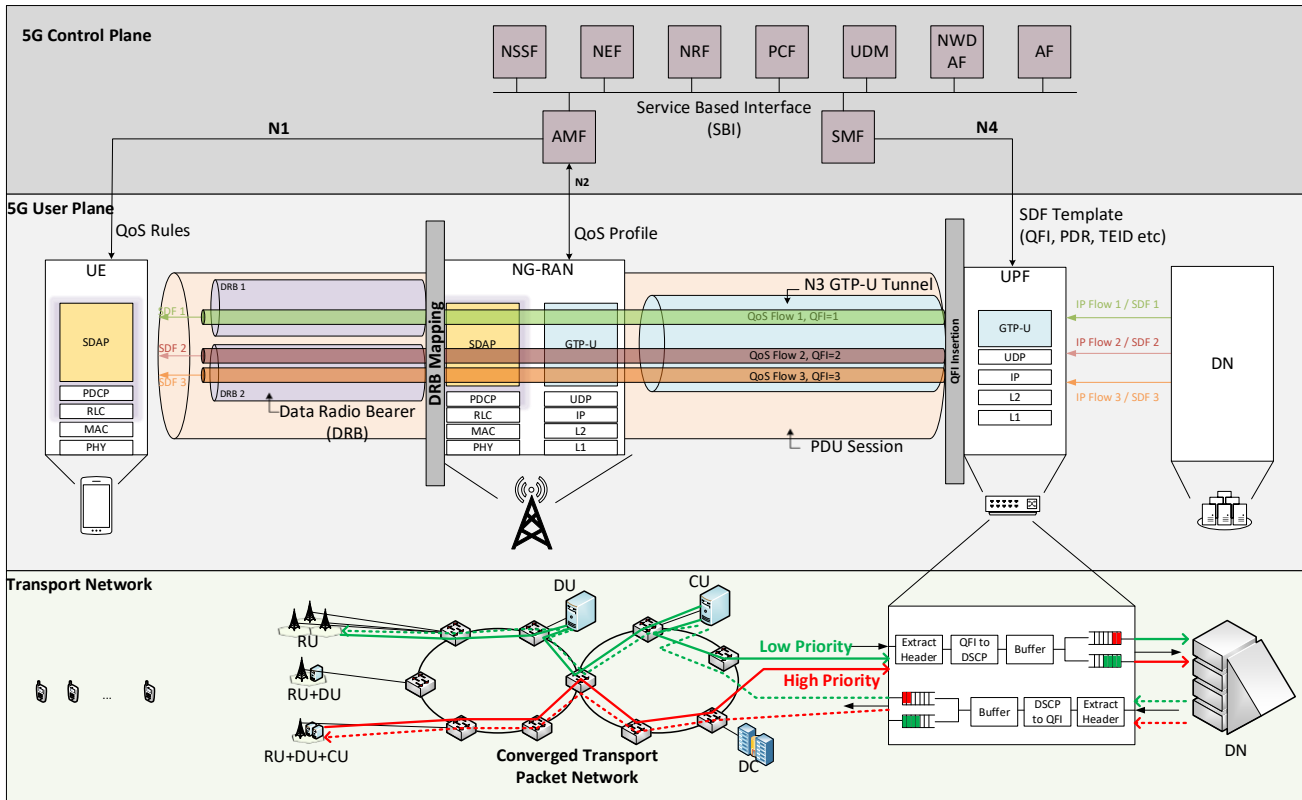


Figure 5. 2: 5G QoS service flows and protocol adaptation.

example of QoS flows defined in the current generation of 3GPP systems is shown in **Table 5. 1** [2].

The first step in the PDU Session Establishment process is the registration of the UE with the 5G core AMF. Then, the AMF contacts another core function the SMF that manages the PDU sessions and controls the 5G QoS Flows. The SMF is responsible for verifying that the UE request matches the user subscription. For this task, the SMF reaches out two other 5G core network functions, namely the UDM and the PCF [11] and [12]. The UDM provides information about the user subscription and authentication. The PCF provides policy rules for QoS Flows and flow-based charging to the SMF.

After the validation of the user subscription information, the SMF decides the appropriate user plane path that will carry the requested service and communicates the necessary information to the UPFs. In the final step, the AMF notifies the RAN and the UE with the QoS information retrieved from the SMF [2]. At the end, when all QoS information is shared among the involved entities, the PDU session is established, first for the UL and then for the DL data. **Figure 5. 2** also shows the tunnels used to implement the PDU session including the wireless DRB, between the UE and the RAN and wired GTP-U tunnel (N3) between the NG-RAN and the UPF.

For the UL transmission, a 5G service starts from the UE’s application layer. The packets from the applications are associated with specific QoS Flows, based on the QoS Rules that are preconfigured at the UE, or are received from the SMF. According to that association, the data packets are loaded to specific DRBs and transferred to the RAN. A

Table 5. 1

STANDARDIZED 5QI TO QoS CHARACTERISTICS MAPPING [2]

5QI Value	Resource Type	Default Priority Level	Packet Delay	Example Service
1	GBR	20	100ms	Conversational Voice
2	GBR	40	150ms	Conversational Video (Live Streaming)
3	GBR	30	50ms	Real Time Gaming, V2X messages Process
5	Non-GBR	10	100ms	IMS Signalling
80	Non-GBR	68	10ms	Low Latency eMBB applications Augmented Reality
82	Delay-critical GBR	19	10ms	Discrete Automation
84	Delay-critical GBR	24	30ms	Intelligent transport systems

DRB can be associated with one or more QFIs. The RAN marks the packets with the suitable QFI, based on the priority level of the QoS Flow that is associated with this QFI. The marking is performed by adding the QFI value in the encapsulation at the GTP-U header. When the tunnel ends at the UPF, the UPF translates the QFIs to IP flows, according to the information received from the SMF, the Packet Detection Rules (PDRs). The prioritization of the QFIs can be performed through appropriate mapping of DSCP to QFI values in the IP header (lower level of Figure 5. 2).

The DL process starts from the IP packets that arrive at the UPF. The UPF relies on the PDRs, that are received from the SMF, to perform the classification of packets to QoS Flows. A header with the value of the associated QFI is added to the packets. Then, through the N3 GTP-U tunnel, packets are transmitted to the RAN, where they are mapped to specific DRBs according to their QFI value and the QoS profile that is provided to the RAN. The DRBs are forwarded to the UEs, where the data packets reach their destination.

Figure 5. 2 provides a graphical representation of the entire process, as well as a high-level overview of how the various components of the 5G systems interact to set up the end-to-end service flow with QoS guarantees. In the protocol stack the SDAP protocol is introduced. The SDAP is responsible for the QoS Flow management between the UE and the RAN. In the downlink it maps a specific QoS Flow inside a PDU session to a suitable DRB. In the uplink it marks the packets with the appropriate QFI, in order to be treated appropriately by the core network.

To provide the required service level guarantees, the underlying transport network infrastructure should be also appropriately managed providing the necessary levels of isolation and capacity to co-host flows with very different service requirements. In the IEEE Standard for Packet-based Fronthaul Transport Networks [13], this is performed through transport network slicing that based on the technologies used can have either *deterministic* or *statistical behavior*. Adopting transport technologies based on e.g. Time Division Multiplexing (TDM) or optical networks supporting WDM can offer *Deterministic* slicing. In this case, a network slice handling a user plane tunnel can be associated to a specific timeslot(s) or wavelength, and then can be multiplexed in the frequency or time domain over a line [13]. Statistical multiplexing can be applied in cases where the underlying infrastructure comprises some form of packet switching e.g. through Ethernet switches. In this case, slicing is implemented through appropriate queuing management and scheduling policies at the ingress/egress ports of each switch providing the required QoS assurance. The relevant policies are applied to these switches through either local or SDN controllers that continuously interact with the CN.

5.2.2. 5G Data Collection and Analytics

AF is expected to play key role in 6G systems as it can influence traffic routing and steering decisions taken by the SMF. With a newly introduced interface (N5) [12], applications can also interact with the PCF of the CN and provide dynamic session information that is essential for the PCF to operate. In case where an AF is not authorized to access NF functions of the CN directly, the NEF can be used to act as a mediator between external AFs and the CN [8]. Through the NEF AFs can interact with the CN subscribing/publishing information that is essential for the applications to operate [2].

This also includes monitoring, profiling and analytics data provided by the Network Data Analytics Function (NWDAF). NWDAF is a newly introduced NF and has been added to the 5G Core in an effort to facilitate the introduction of AI and ML techniques in the 5G system. It collects data from NFs and AFs regarding user mobility, application load, QoS etc, along with data from the 5G management plane, related to fault management and radio analytics. Those data are fed to the ML/AI algorithms inside the NWDAF, for NF optimization through statistical analysis and prediction information, QoE optimization, efficient resource allocation, dynamic policy provisioning, service profiling etc. NWDAF consists of two services, namely the events subscription and the analytics information. The events subscription grants access to NFs to subscribe/unsubscribe and get notified about analytics events. The analytics information service is responsible for the application of the AI/ML algorithms [14].

In order to collect data related to a specific UE or group of UEs, NWDAF must first discover the NFs that are serving this UE/group of UEs. The UDM entity, that is responsible for the access control and session establishment of network subscribers, can provide the necessary information for the discovery of the SMF and AMF. For all other NFs discovery, including AFs, NWDAF should communicate with the NRF. NRF stores

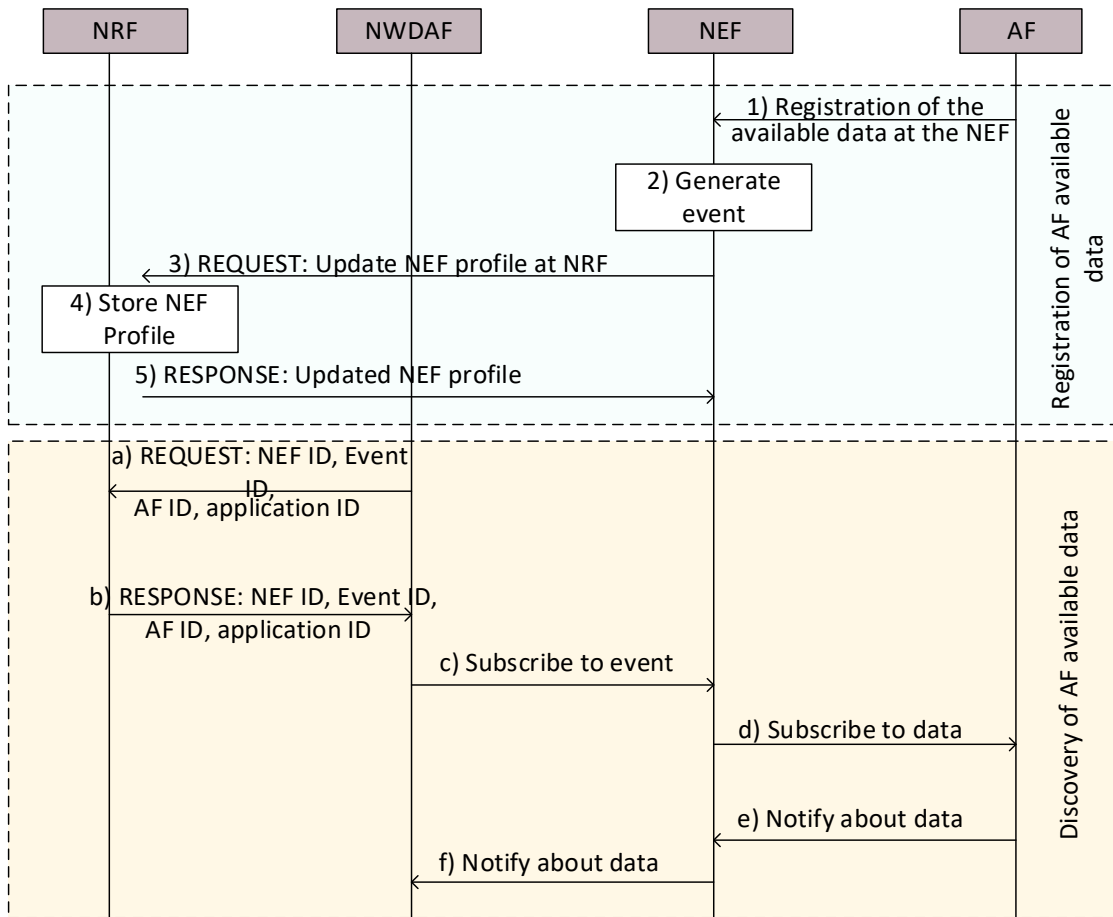


Figure 5. 3: NWDAF Data Collection from AF via NEF.

a record for all 5G existing elements in the network, as well as their supported services. Using NRF, NFs can discover, subscribe and get notified about each other’s events [15].

As an example, Figure 5. 3 depicts the whole 5G Data Collection for a non-trusted AF. The AF at first registers with the NEF. The NEF creates a new event for the AF, and informs the NRF. When the NWDAF required metrics are collected from the AF, it consults the NRF, for identifying the NEF that the AF is registered to. A subscription to the AF event is sent to the NEF, which in its turn sends a data subscription message to the AF. The AF informs the NEF about the data, and the NEF notifies the NWDAF.

5.2.3. Problem Statement

To control network traffic, UPF nodes have been designed to perform a set of tasks including [2]: mapping of traffic to appropriate tunnels, steering of packets to appropriate output ports, packet counting for charging and policy control purposes, buffering and queuing management for traffic service differentiation and assurance of end-to-end delays, deep packet inspection etc. To perform these actions UPFs also support an extensive set of protocols such as, GTP-U, PFCP, IP, PDCP, QoS mapping etc. through appropriate mapping of DSCP classified IP traffic coming from the external DN etc.

However, as discussed in [7] and [5] UPF functionalities come at the cost of increased computational requirements that increase exponentially with the number of flows they need to process. If the computational resources allocated for the operation of a UPF are not sufficient, the quality of the services using this UPF will be degraded. This problem is further exaggerated under scenarios involving a large number of flows with low latency requirements requesting termination at a MEC facility placed close to the edge. At first sight, the selection of a UPF node that is at close proximity with the RAN can be the optimal choice, since this may result in reduced latency and transport network resource requirements. However, if the core network associates all UEs with the same UPF, congestion may arise resulting in poor performance and increased latency.

In current 5G systems, this decision is taken centrally by the SMF. For every traffic flow and depending on its service QoS requirements, the SMF is responsible to decide how and by which NFs this flow will be handled. It is clear that without appropriate optimization, users seeking the highest throughput, would uncontrollably exploit both network and compute resources, eventually leading to poor QoS for all users. On the other hand, it is also obvious that, centralized decision making (taken by a centralized entity e.g. MANO) does not scale with the number of connected devices.

In this work, we consider a 6G system with multiple UEs requiring the establishment of end-to-end connections with ASs hosted either at a local MEC or at a remote CC facility. Each AF that is responsible for the management of its associated application communicates with the SMF and following the process described above decides the path and compute facility where the relevant session will be terminated. However, in case where a large number of connections traverse the same resources (i.e., use the same UPF node), the system will be congested resulting in QoE degradation. This is especially true for services with strict latency constraints as the majority of the AFs will have the incentive to terminate their connections to a closely located MEC server.

To address this issue, we propose an optimization framework based on EGT. EGT allows AFs to periodically reconsider the routing and traffic steering policies having a primary objective to optimize the QoE. At the same time, when an AF selects a specific policy, it is charged by the Core Charging Function (CHF) for the compute and network resources it has utilized. The relevant monitoring and profiling data are provided to the CHF through the NWDAF while the price paid by UEs for each application increases with the resources used. Therefore, a secondary objective for the AFs is to minimize the relevant costs for the UEs through the identification of the optimal network/compute resources used by the corresponding service slices that takes place during policy updates.

5.3. Related Work

Currently QoS in 5G systems is an area of intense scientific research, due to the high dependence of QoS provisioning schemes on the system level implementation. From an architectural perspective, end-to-end QoS management in 5G and beyond networks is addressed through slicing and mapping of flows to the appropriate tunnel at the different segments of the network [2]. For example, in [16] the authors proposed a network slicing architecture allowing multiple virtualized 5G functions (based on OpenAirInterface) to run over a common physical infrastructure. However, topics related to service isolation and dynamic QoS Flow control have not been addressed. In [17] the authors addressed the problem of isolation at network level using Active Queue Management (AQM) algorithms. AQM policies are adopted to address limited queues and guarantee the delay limits of QoS Flows, while achieving maximum throughput.

The relevant analysis is limited to the UPF function while application related issues are not considered. Communication between the 5G network and the applications, in support of QoS requirements is investigated in [18]. Based on this communication, the number of Resource Blocks that are necessary for the operation of the application are optimally assigned through a scheduling mechanism. [19] focuses on the 5G system operation and the potential of resource allocation using QoS Flows and network slicing. The objective of the study is to identify how QoS Flows and network slicing can optimally be exploited to ensure a high customer QoE, while efficiently utilizing available network resources.

Our study is differentiated from the previous approaches by a) introducing mechanisms to ensure isolation using the concept of 5G network slicing and assigning flows to the appropriate queues at the transport network, b) creating dynamic QoS Flows by terminating sessions to the appropriate servers where the application server is hosted, c) considering the role of the application function that tries to maximize QoE of the end users in the decision making process, d) periodically reconsidering the optimal QoS Flow problem using statistics collected by the NWDAF through an operational 5G network.

The introduction of data collection and analytics through the NWDAF entity, plays a key role in supporting 5G services with different QoS requirements. NDWAF facilitates processing of collected data corresponding to specific network metrics through AI and ML algorithms. In [14] the authors employ ML algorithms for network load prediction, and anomaly detection. In [20], [21] and [22] the QoE level is determined through the use of NWDAF. More specifically, in [20] the QoE level is derived by correlating network statistics and the application data within NWDAF, while [21] performs QoE evaluation exploiting AI techniques. Finally, [22] investigates the trade-off between cost and accuracy of ML-based methods when applied for QoE estimation through the NWDAF. In our analysis, monitoring data from the 5G testbed are exposed to the NWDAF in order to extract metrics related to network and compute resource utilization. The relevant metrics are also exploited by the CHF to calculate the cost and charging models.

Trying to identify suitable models that can be used to analyze and optimally operate complex systems such as 5G networks EGT appears to be a promising candidate. EGT is widely used to model and analyze the performance of various types of networks, spanning from social networks to computer and communication networks. What makes EGT a great candidate as an optimization tool for 5G systems is its distributed nature along with its resilience in environments that lack global knowledge. There are several examples in the literature that rely on EGT to solve relevant problems e.g. EGT has been used to identify the optimal Radio Access Technology (RAT) for a set of users in a heterogeneous wireless cellular network in [23]. The RAT selection problem is formulated as an incomplete information game where users choose among multiple access technologies selfishly. In our previous work [24], [25] we used EGT to address the concept of flexible functional splits of BBU functions supported by a suitable FH network. In this environment, the RUs can act as non-cooperative/self-optimizing players in an evolutionary game, with the objective to unilaterally minimize the infrastructure operational costs in terms of power consumption. Based on the stability of the system, the controller placement problem is also discussed. In [26] we presented a distributed model to dynamically place the UPF in a 5G network targeting to minimize the overall service delay. In [27] we extended our previous work, by using real lab measurements to evaluate our proposed scheme in a real 5G environment. In the present study, the EGT model is used by the AFs to select the optimal policies for the UEs in an end-to-end fashion by a) calculating the relevant payoff values based on the

QoE and cost models exposed by the 5G system and b) recommending optimal policies to the 5G core and transport network controllers.

5.4. Theoretical Background: Evolutionary Game Theory

EGT studies strategic games of non-cooperative individuals that interact repeatedly in time [28]. It differs with regards to classic Game Theory as it considers that players are irrational and can only learn through experience. The focus is directed on how the strategies can "survive" through evolution and how they can help players who choose these to "strengthen" and better meet their needs.

Mainly, an evolutionary process involves two major procedures. The first is a mutation procedure that produces varieties. The role of mutation is highlighted by the notion of Evolutionary Stable Strategies (ESS) – which is a refinement of the Nash Equilibrium (NE) – a strategy is defined to be an ESS if its fitness is greater than the fitness of any rare mutant. Mutants are agents with a different strategy than the rest of the population. The second procedure refers to selection and is responsible for favoring certain varieties over others. The most popular selection model is the replicator dynamics model, which assumes that the percentage of individuals inside a population grows or declines, depending on the performance of the adopted strategy. The replicator equation assumes that the populations are infinite, homogeneous, and well mixed (each individual is equally likely to interact with any other individual) and with no mutations. The mathematical formula is given below:

$$\dot{x}_j(t) = x_j(t) \left(\pi_j(\mathbf{x}(t)) - \bar{\pi}(\mathbf{x}(t)) \right), \quad (1.a)$$

$$j \in S$$

where S is the set of strategies that are available to the population, $\mathbf{x}(t) = [x_1(t) \ x_2(t) \ \dots \ x_j(t) \ \dots]^T$ is the population state at time t with $x_j(t)$ being the proportion of the population that uses strategy j at time t . $\pi_j(\mathbf{x}(t))$, $\bar{\pi}(\mathbf{x}(t))$ are used to describe the expected payoff of strategy j and the mean payoff, respectively [28]. The payoffs provide a means to measure the level of satisfaction of choosing a specific strategy. Thus, according to (1), the growth rate \dot{x}_j/x_j of the strategies that are currently used is equal to the excess of the current payoff $\pi_j(\mathbf{x}(t))$ over the average population's payoff, $\bar{\pi}(\mathbf{x}(t))$.

The replicator equation can be extended to multiple populations to cover scenarios where the interacting agents are drawn from different populations. In this case, each population is characterized by a state vector \mathbf{x}_i that depicts the distribution of the available strategies inside the population. At each interaction a player is drawn from each population to play the game. The multi-population dynamics equation can be described as follows:

$$\dot{x}_{ih}(t) = x_{ih}(t) \left(\pi_i(e_i^h, \mathbf{x}_{-i}(t)) - \bar{\pi}(\mathbf{x}(t)) \right), \quad i \in N, h \in S \quad (1.b)$$

where N is the number of populations, S is the set of strategies that are available to the populations, $\mathbf{x}_i(t) = [x_{i1}(t) \ x_{i2}(t) \ \dots \ x_{ih}(t) \ \dots]^T$ is population's i state at time t with $x_{ih}(t)$ representing the proportion of the population i that uses strategy h at time t , and

$x_{-i}(t)$ the state of all populations except i . Finally, $\pi_i(e_i^h, x_{-i}(t))$, $\bar{\pi}(x(t))$ is the expected payoff of strategy of an agent in population i that uses strategy h and the mean payoff of all strategies, respectively.

The replicator equation provides a deterministic analysis of a stochastic process, only if the population involved is considered infinite. However, real world scenarios comprise populations of finite size, exhibiting stochastic effects that may play a sizable role to the accuracy of the replicator equation. For a long time, areas in physics try to analyze the impact of stochastic effects in finite-size systems. However, the majority of the studies on stochastic evolutionary dynamics in finite populations, have been performed mostly numerically. For example, individual-based computer simulations frequently include finite-size populations, which naturally integrate such stochastic effects. In these, several levels of selection were implemented and tested, ranging from a strong selection framework captured by the finite-population equivalent to the replicator dynamics to an extreme selection pressure reflected by imitation dynamics, used as a metaphor for cultural evolution [29], [30].

A process that takes into consideration the random drifts in the selection mechanism of a finite population is the Moran process. The corresponding microscopic dynamics is characterized by three simple steps [31]:

1. Selection: An individual from the entire population is chosen to reproduce, with probability proportional to its fitness.
2. The individual produces one identical offspring.
3. The offspring replaces an existing individual randomly.

One may think of the Moran process as a random walk that moves towards the direction of the strategy that is optimal given the population state [30]. The fitness of an individual in the original Moran process is genetically determined, independent of the frequency of the individual's strategy, and is not affected by interactions with others inside the population. The fitness is generalized in order to study evolutionary dynamics in finite populations, by making it dependent on the frequency of the strategy inside the population. This leads to a new definition of the ESS concept, when dealing with finite populations: (a) rare mutants cannot invade the strategy and (b) the probability that a rare mutant strategy overtakes the resident one is smaller than it is for random drift.

5.5. Problem Formulation

5.5.1. Game Formulation

In the scenario studied in this paper we consider the uplink transmission of a resource constrained 5G network shown in **Figure 5. 2**. Without loss of generality the analysis has been limited to the uplink direction. However, given that the operations for both directions are similar, the same procedure can be applied to the downlink direction. UEs with strict priority service requirements (e.g., Low Latency eMBB applications such as augmented reality with 5QI value equal to 80 as shown in **Table 5. 1**) initiate the PDU Session Establishment process by transmitting the relevant request to the AMF. The AMF contacts the SMF, which in turn checks whether the UE requests are compliant with the user subscription. Once subscription information has been verified the SMF selects a UPF to serve the PDU Session. The related QoS information (QoS profile, QoS

rules, PDRs) is then provided to the involved entities (RAN, UE, UPF). The selection of the UPF and MEC facility where the sessions will be terminated is a key decision determining the network and compute resources required to support the corresponding services. In contrast to conventional systems where these decisions are calculated by a central authority, the proposed scheme assumes that these decisions are taken in a distributed manner.

Specifically, in the present study we consider a set of N UEs each one requesting the establishment of a service with strict QoS characteristic (i.e., 5QI value 80). For this type of service, all UEs will have the incentive to be interconnected with an AS hosted in a local MEC facility. As the network is resource constrained, the demands of all applications cannot be addressed in this location and, therefore, a portion of user requests is forwarded to a more centralized cloud facility. To keep the analysis tractable, we consider two compute node classes. However, the analysis, can be extended to address scenarios with multiple ASs hosted at different locations. Therefore, two possible ASs are available for selection by the UEs:

A closely located AS (CL-AS) hosted at a MEC facility. This AS is connected to the 6G system using a UPF node that is also placed close to the edge. This type of implementation has been primarily envisioned by the MEC forum as an obvious approach to reduce service latencies.

A remotely located AS (RL-AS) hosted at a central cloud facility. In this case, the UPF node is placed deeply into the network. To support services with strict latency constraints the underlying network resources are also adjusted. In this case, network switches take appropriate actions (i.e., assign these flows to a higher priority queuing port than the standard, as shown in the lower part of **Figure 5. 2**) to compensate the increased path length interconnecting UEs with the RL-ASs.

We will formulate the problem using the concepts of the Moran process described in Section 5.4(A), for the following reasons:

- It assumes a finite population, thus it maps to the number of UEs in our case
- The changes in the population occur at each timestep. At each timestep the network metrics can be used to update the decisions of the UEs.
- At each timestep only one UE reconsiders its strategy, thus the exchanged information for the policy update does not overwhelm the network.

Each UE selects one of the two available ASs for termination and depending on their choice, they receive a payoff that quantifies the satisfaction level associated with that choice. The evolution of the strategies is governed by the following evolutionary process. At each time step, a random UE is chosen from the entire population to update its strategy. The UE measures its payoff and imitates with some probability an opponent, that is a randomly chosen UE from the rest of the population. This probability depends on the payoff difference of the two UEs and is given by the Fermi distribution [29]:

$$p = \frac{1}{1 + e^{-w(\pi_A - \pi_B)}} \quad (1)$$

In (1), π_A, π_B are the payoffs of the UE of interest and its opponent respectively. According to this equation a UE will imitate the strategy of its opponent with probability higher than $1/2$ if its payoff is lower than the payoff of its opponent.

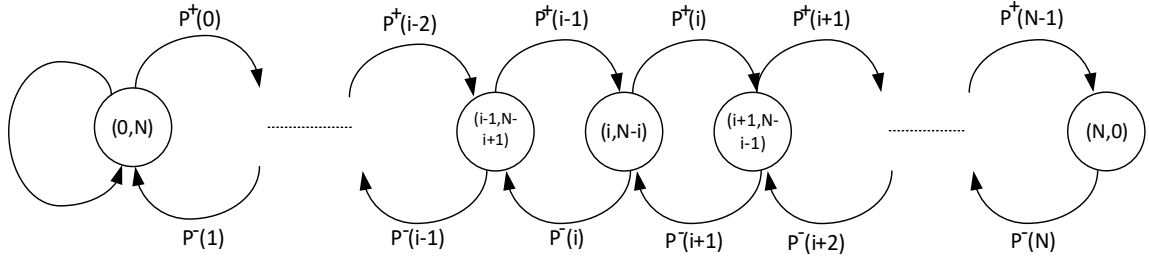


Figure 5. 4: UE population state transition probability.

Now let i be the number of UEs that are currently choosing the CL-AS. The remaining $N - i$ users will be allocated to the RL-AS. According to the proposed evolutionary dynamics model, at each time step, the population i can change at most by one. Hence, assuming that the system is in state $(i, N - i)$, the transition probability to state $(i + 1, N - i - 1)$ will be given by

$$P_{i,i+1} = \frac{i}{N} \frac{N-i}{N} \frac{1}{1 + e^{-w(\pi_{HP}(i) - \pi_{LP}(i))}} \quad (2)$$

Similarly, the probability that the system will move to state $(i - 1, N - i + 1)$ is:

$$P_{i,i-1} = \frac{i}{N} \frac{N-i}{N} \frac{1}{1 + e^{+w(\pi_{CL}(i) - \pi_{RL}(i))}} \quad (3)$$

Finally, the probability that the UEs will remain at the same service class will be given through:

$$P_{i,i} = 1 - P_{i,i+1} - P_{i,i-1} \quad (4)$$

In (2)-(4), w is the intensity of selection which governs, how much influence has the payoff to the adaptation process. According to these probabilities a UE will imitate the strategy of its opponent with probability higher than $1/2$ if its payoff is lower than the payoff of its opponent.

5.5.2. Payoff Function

A key parameter affecting the performance of the proposed model is the selection of the appropriate cost function. Since we are interested in identifying the optimal tradeoff between compute and network utilization, we formulate the payoff function for each UE as the sum of the total compute and network resources needed to carry out its request. The payoff that each UE receives when it selects an CL or RL AS is given by:

$$\pi_{CL}(i) = \frac{h_c + h_n}{h_c \cdot CPU_{CL}(i) + h_n \cdot NET_{CL}(i)} \quad (5)$$

$$\pi_{RL}(i) = \frac{h_c + h_n}{h_c \cdot CPU_{RL}(N - i) + h_n \cdot NET_{RL}(N - i)} \quad (6)$$

In equations (5) and (6), CPU_s and NET_s correspond to the compute and network impact, respectively, of the two available choices s ($s \in \{CL, RL\}$), h_c and h_n are biasing factors that sum to 1, which can be determined according to the resource (CPU or

network) will be given advantage to. CPU_s and NET_s are evaluated using a profiling process carried out over an operational 5G testbed described in Section 5.6.

5.5.3. Dynamics of Adaptation Process

In order to find the dynamics of the stochastic process described above, first, we formulate the master equation of the Markovian chain described in **Figure 5. 4**:

$$\begin{aligned}
 k^{\tau+1}(i) - k^{\tau+1}(i) & \quad (7) \\
 & = k^{\tau}(i-1)P^+(i-1) \\
 & + k^{\tau}(i+1)P^-(i+1) \\
 & - k^{\tau}(i)P^+(i) - k^{\tau}(i)P^-(i)
 \end{aligned}$$

Where $k^{\tau}(i)$ is the probability that the system is in state i at time τ , and $P^+(i), P^-(i)$ are the transition probabilities $P_{i,i+1}$ and $P_{i,i-1}$ respectively. We shift the variables by introducing the notations: $x = i/N, t = \tau/N, d(x, t) = Nk^{\tau}(i)$ (probability density function). Eq. (8) yields to the form:

$$\begin{aligned}
 d\left(x, t + \frac{1}{N}\right) - d(x, t) & \\
 & = d\left(x - \frac{1}{N}, t\right)P^+\left(x - \frac{1}{N}\right) \\
 & + d\left(x + \frac{1}{N}, t\right)P^-\left(x + \frac{1}{N}\right) \\
 & - d(x, t)P^+(x) \\
 & - d(x, t)P^-(x)
 \end{aligned} \quad (8)$$

For large populations ($N \gg 1$), the expansion of the probability density and the transition probabilities leads to the Fokker-Planck Equation [32] (See Appendix A):

$$\begin{aligned}
 \frac{\partial}{\partial t}d(x, t) & \\
 & = -\frac{\partial}{\partial x}[(P^+(x) - P^-(x))d(x, t)] \\
 & + \frac{1}{2}\frac{\partial^2}{\partial x^2}\left[\frac{P^+(x) + P^-(x)}{N}d(x, t)\right]
 \end{aligned} \quad (9)$$

With drift coefficient $\mu(x) = (P^+(x) - P^-(x))$ and diffusion coefficient $\sigma^2(x) = \frac{P^+(x) + P^-(x)}{N}$. The drift describes the average motion of the system, while the diffusion “widens” the probability distribution in the course of time. From this, we can extract a Langevin equation, after implementing the Ito calculus (noise is uncorrelated with time) [33]:

$$\dot{x} = \mu(x) + \sigma^2(x)\xi \quad (10)$$

where ξ corresponds to uncorrelated Gaussian noise. For large populations, the diffusion term vanishes and under weak selection ($w \ll 1$) this AS selection process can be approximated by the replicator dynamics equation (1..a) written in the following form [29] (see Appendix B):

$$\dot{x}(t) = x(t)(1 - x(t))\left(\pi_{CL}(x(t)) - \pi_{RL}(x(t))\right) \quad (11)$$

with x being the percentage of population adopting the CL strategy that is equal to $x = i/N$. When multiple AFs interact, the above analysis leads to the multi-population model of the replicator equation, described in Section 5.4.

5.5.4. Stability Analysis

According to the replicator equation the population evolves until it reaches a situation where all individuals gain reach the same payoff, and thus have no incentive to deviate from their current strategy. When this happens, the population has reached an equilibrium. Evolutionary equilibrium is defined as the fixed point of the replicator dynamics. In the single-population replicator dynamics defined in (12) two types of evolutionary equilibrium, namely, boundary evolutionary equilibrium and interior evolutionary equilibrium exist. The boundary evolutionary equilibrium corresponds to the case where there exists a population share $x_i = 1$, while $x_j = 0$ for all $i \neq j \in S$. The interior equilibrium x^* corresponds to the case where we have $x_i \in (0, 1), \forall i \in S$. The boundary equilibria are not stable in the sense that any small perturbation will make the system deviate from the equilibrium state.

To evaluate the stability of the interior evolutionary equilibrium, the eigenvalues of the Jacobian matrix corresponding to the replicator dynamics need to be evaluated. The system is stable if all eigenvalues have a negative real part. With the analytical expressions obtained from the stochastic geometry-based analysis, the stability of the equilibrium is analytically tractable in some cases.

As we have one system equation, an equivalence to this condition is that the Jacobian should be negative definite. To this end, we first denote by f the right hand side of (12). Accordingly, we have:

$$\begin{aligned} & \frac{df}{dx} \\ &= x(1-x) \left(\frac{d\pi_{CL}(\mathbf{x}(t))}{dx} - \frac{d\pi_{RL}(\mathbf{x}(t))}{dx} \right) \quad (12) \\ &+ (1-2x) \underbrace{\left(\pi_{CL}(\mathbf{x}(t)) - \pi_{RL}(\mathbf{x}(t)) \right)}_{0 \text{ at equilibrium point}} \end{aligned}$$

Consequently, the stability of the equilibrium point is highly affected by the payoff function. In section VI we extract the relations for the CPU and network consumption, and substitute the results to Eq. (6) and (7). Then, we calculate the derivatives of $\pi_{CL}(\mathbf{x}(t))$ and $\pi_{RL}(\mathbf{x}(t))$ (see Appendix) and substitute the results to Eq. (13). We obtain that $\frac{df}{dx}$ is strictly negative for any non-zero value of x . This means that all the interior equilibrium points are stable.

The stability analysis for the case of multiple AFs, differs from the single-population case. According to [28], no interior population state can be stable in the multi-population replicator equation. This means that for each AF population of UES, the only one strategy will eventually be present in the equilibrium state.

5.5.5. From infinite to finite population

If we take the number $i, i \in \{0, 1, \dots, N\}$, as the state of the population, then the state evolves as a Markov chain. Boundary states $i = 0$ and $i = N$ are absorbing, in the sense

that if the system reaches one of these, it will remain there forever. It can be seen that all other states $i = 1, 2, \dots, N - 1$ are transient. Fixation, i.e., the chain eventually entering one of the absorbing states, will happen with probability 1 [28]. However, a sufficiently large expected fixation times may indicate co-existence [30]. In general, the fixation probability to reach the absorbing state with 100% UEs terminating their connection at the CL-AS given that the initial state is k can be written as in [29] and [34]:

$$x_k = \frac{1 + \sum_{l=1}^{k-1} \prod_{m=1}^l \gamma_m}{1 + \sum_{l=1}^{N-1} \prod_{m=1}^l \gamma_m} \quad (13)$$

with

$$x_1 = \frac{1}{1 + \sum_{l=1}^{N-1} \prod_{m=1}^l \gamma_m} \quad (14)$$

In (15), γ_i is defined through the following equation:

$$\gamma_i = \frac{P_{i,i-1}}{P_{i,i+1}} = e^{-w(\pi_{HP}(i) - \pi_{LP}(i))} \quad (15)$$

Due to the stochastic character of the strategy update, a single mutant with a more beneficial strategy may be eliminated after a few time steps in a uniform population. Thus, the fixation probability of a single initial mutant (ρ) in combination with the payoff difference of the two strategies in the initial state can give us significant information on the evolutionary stability of the system. We denote with $\rho_{CL} = x_1$, $\rho_{RL} = 1 - x_{N-1}$ the fixation probability of one initial UE in the population that terminates its connection at the CL-AS and RL-AS, respectively, and $\rho_{ND} = 1/N$ is the fixation probability of a random drift, that is the neutral case where at each timestep, a UE updates its strategy by imitating another random UE from the population. According to [35], where the authors analyze the possible evolutionary scenarios of the Moran process, if $g_i = \pi_{CL}(i) - \pi_{RL}(i)$ is the payoff difference of the two strategies at state i , there are four possible outcomes:

1. *Domination of CL-strategy:* In this case, the following relation holds

$$\rho_{RL} < \rho_{ND} < \rho_{CL}, g_1 > 0, g_{N-1} > 0 \quad (16)$$

indicating that there is no interior equilibrium point.

2. *Domination of RL-strategy:* As in case 1, an interior equilibrium point does not exist:

$$\rho_{RL} > \rho_{ND} > \rho_{CL}, g < 0, g_{N-1} < 0 \quad (17)$$

3. *Coexistence*

$$\rho_{RL} > \rho_{ND}, \rho_{CL} > \rho_{ND}, g_1 > 0, g_{N-1} < 0 \quad (18)$$

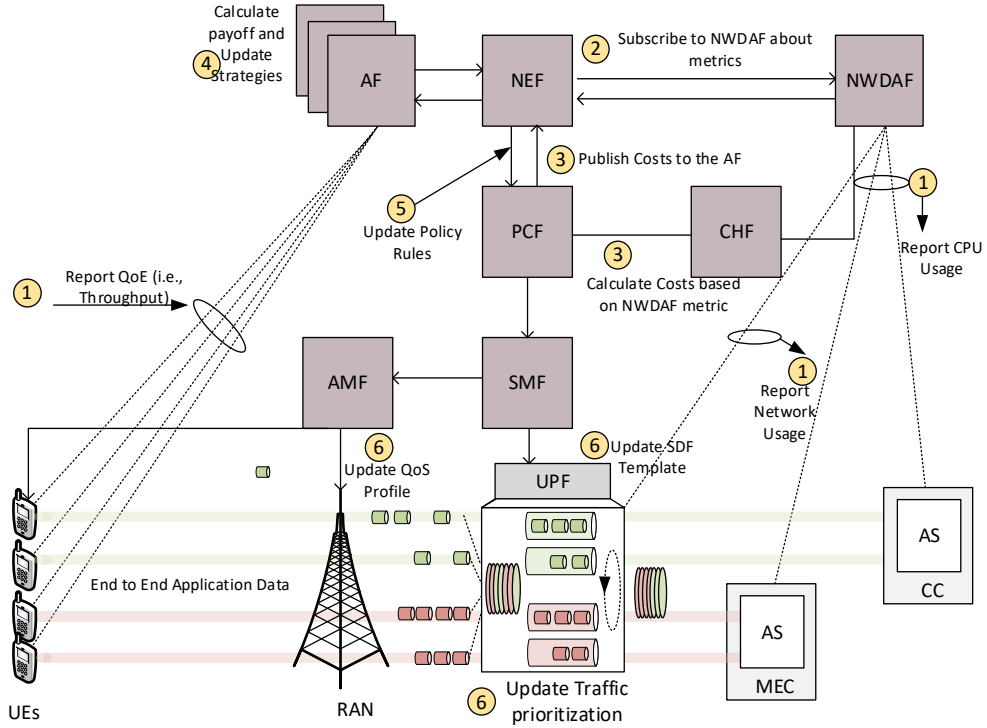


Figure 5.5: Strategy Adaptation Process

The interior equilibrium point of replicator dynamics is stable.

4. Bi-stability

$$\rho_{CL} < \rho_{ND}, \rho_{RL} < \rho_{ND}, g_1 < 0, g_{N-1} > 0 \quad (19)$$

The interior equilibrium point of replicator dynamics is unstable.

A comparison of ρ_{CL} and ρ_{RL} is also of interest in favor of which strategy the process spends more time. The system spends longer in the strategy for which the corresponding ρ is greater (because this strategy needs less invasion attempts to fixate).

5.6. 5G System Profiling

In the present study, the policy adaptation process relies on EGT. The EGT model requires at each timestep a random matching between two UEs. **Figure 5.5** presents the overall proposed procedure, that can be summarized in the following steps:

1. AF is informed regarding the throughput of each UE, and the CPU-usage of the servers that fulfilled the requests. NWDAF collects utilization statistics of the network and compute resources.
2. The NWDAF subscribes to NEF to publish the collected metrics to other NFs.
3. CHF calculates the network/compute costs based on the statistics published by the NWDAF. The relevant results are exposed through NEF to the AF.
4. Each AF calculates its payoff function and updates its strategy adopting the EGT model.

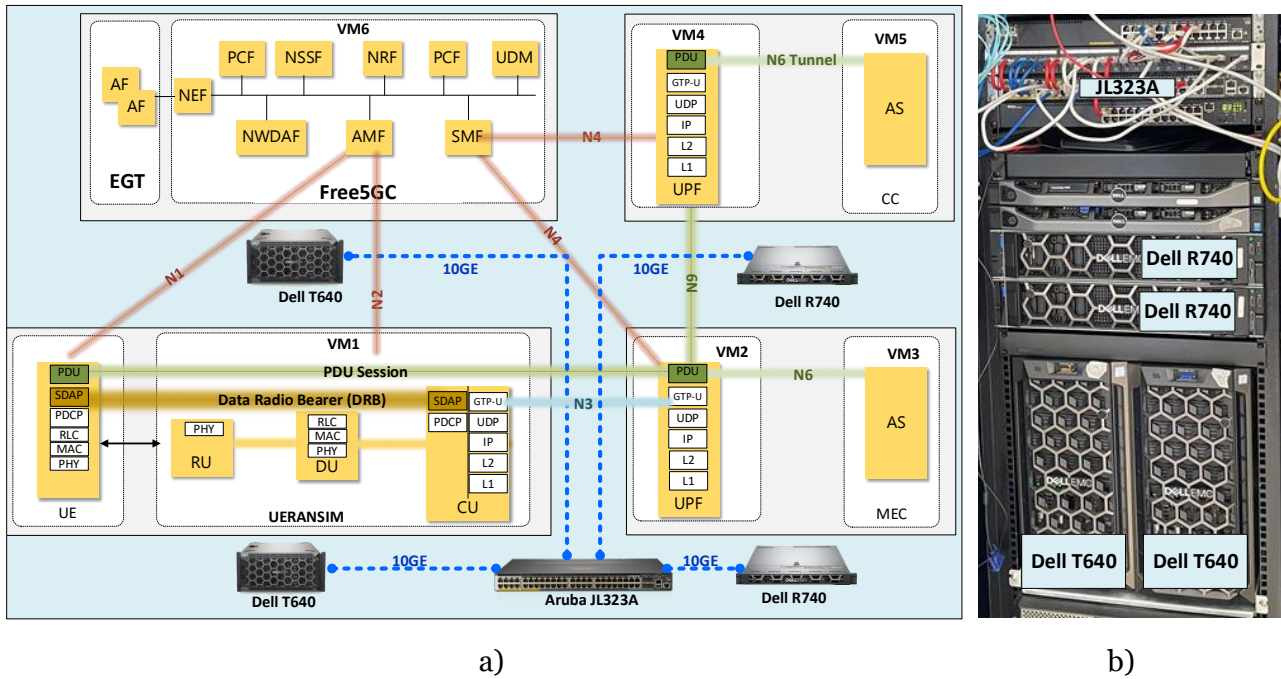


Figure 5. 6: a) 5G network deployment and infrastructure connectivity b) Servers/switch used in the experimentation

5. AF informs the PCF to update its policy rules.
6. PCF informs SMF, AMF and subsequently UPF to update QoS profiles and apply traffic steering and prioritization rules according to the new profiles. The new profile indicates whether a UE in the population changed the AS that it was connected to.

This procedure is repetitive. The timestep between two repetitions is related to the time elapsed for the six above steps to be completed.

In the following we describe the 5G Platform and the cost function that was used for the evaluation of the EGT scheme.

5.6.1. 5G Platform Description

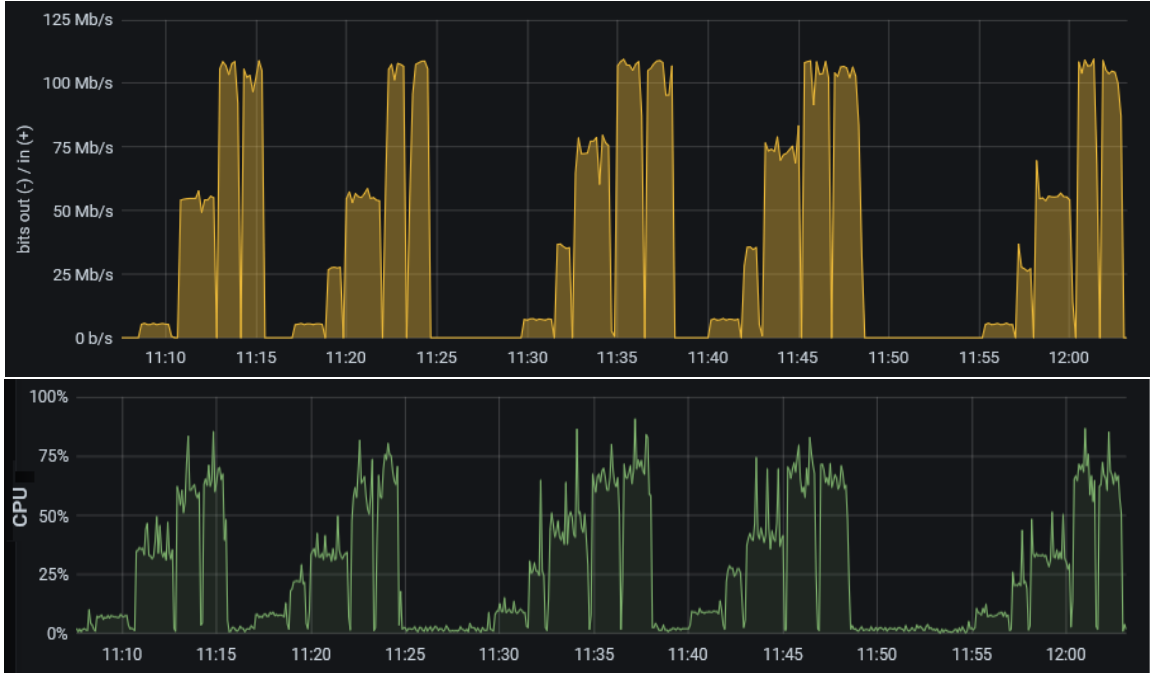
To evaluate the evolutionary service connection establishment problem, an open source 5G platform has been deployed over a virtualized private data center network managed by Openstack. The RAN segment of the 5G testbed platform relies on UERANSIM [36] while the core functionality is provided through Free5GC [37]. Servers hosting virtualized NFs are interconnected with an SDN controlled electro-optical switch (Aruba JL323A) using 10GBASE-T copper and 10G SFP+ ports through an external module (JL085A). This switch provides the ability to perform real-time traffic classification into eight priority levels mapped to eight different queues.

The 5G network deployment and the associated infrastructure connectivity used for the experimentation is shown in **Figure 5. 6**. Virtual Machines (VMs) hosting UPF and AS functionalities have been also placed at the same physical machine (Dell R740). The AS for both the local and the remote deployment is supported by the same resources whereas for the UPFs different sizing options are available. These options along with the key configuration parameters of the VMs are shown in **Table 5. 2**.

Table 5. 2

VM CONFIGURATIONS FOR THE VM HOSTING UPF

VM	No of CPU Cores	RAM [GB]	Storage [GB]
Small	1	2	20
Medium	2	4	40
Large	4	8	80
X Large	8	16	160

**Figure 5. 7:** UE traffic and CPU utilization for the VM hosting a UPF NF

To perform the experimentation a set of UEs is created requesting the establishment of end-to-end connections through the UPF to the AS. For the AS a simple FTP server has been deployed and users request the transfer of large files. The compute/network resources of the system are monitored through Prometheus [38] and visualized using Grafana [39]. An example of a generated traffic pattern per UE and the associated CPU utilization for the UPF is shown in **Figure 5. 7**.

5.6.2. Evaluation Process

To evaluate the performance of the problem an extended experimentation campaign has been performed. The relevant measurements have been used to quantify compute and network resources used by the 5G network to support the required UE connectivity. Measurements have been recorded for a range of 1 up to 10 UEs per gNB requesting the establishment of services with different data rates. All UEs (applications running on the UEs) are treated equally by the 5G system as the same QoS configuration has been applied to all slices.

The impact of the number of UEs (all having the same service requirement) on the CPU utilization of VMs hosting the UPF under different sizing options is shown in **Figure 5. 8**. We observe that the CPU consumption scales almost linearly with the number of the connected UEs. It is also obvious that the sizing option of the VM that hosts the UPF

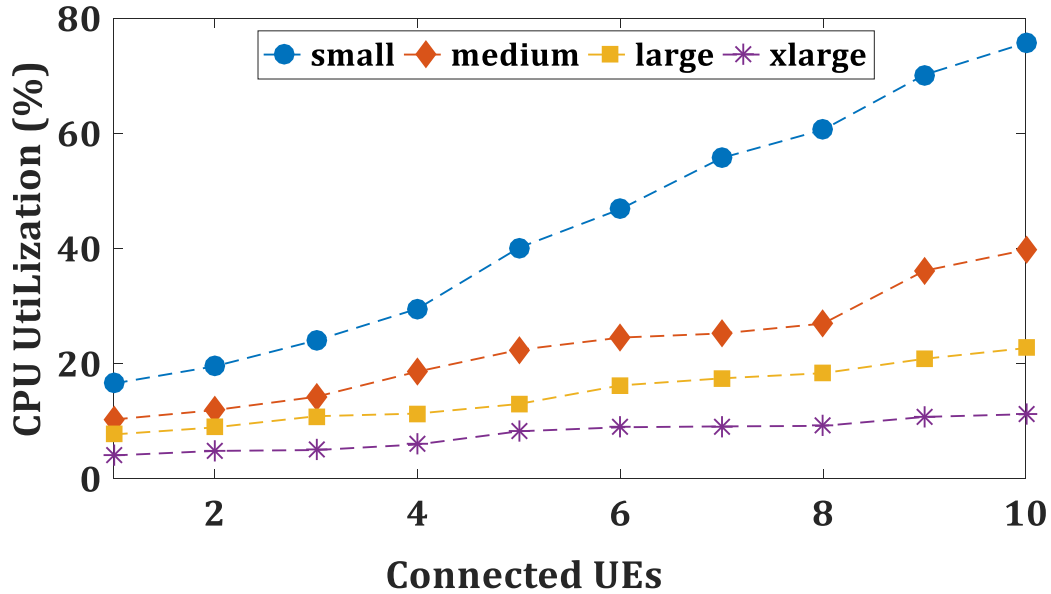


Figure 5. 8: CPU utilization for 4 VMs with different characteristics in terms of cpu, memory and disk.

affects the CPU consumption. The higher the processing capabilities of the VM, the lower is the CPU consumption to carry the UE request.

We then evaluate the impact of network slicing achieved through queuing prioritization on the requested service connections. As most of the contemporary switches support only statistical slicing, a critical parameter that needs to be evaluated is the level of isolation achieved when multiple 5G network slices are established over the same network switches. At this point it should be mentioned that the analysis is limited only to backhaul traffic (gNB-UPF through N3 connections, UPF-UPF through N9 interfaces) as we have assumed that all building blocks (RU, DU, CU) of the gNB are hosted in the same machine. Therefore, the fronthaul/midhaul traffic is implemented through virtual network interfaces. The same holds for UPF-DN interfaces as the associated VMs run also on the same physical server.

To evaluate the impact of network slicing policies on the requested services, we have considered multiple UEs – AS connection scenarios terminated either at the MEC or the central cloud. To implement network slicing we assign connections to different queues of the switch. For simplicity reasons we have considered only two queues (a high priority (HP) namely Q8 and low priority (LP) namely Q3) but the analysis can be extended to multiple queues. The Q8 queue is served in a strict method, meaning that it gets all the bandwidth it needs, and any remaining bandwidth is shared among the non-strict queues based on their need and their configured bandwidth profiles. The Q3 queue, is configured with 30% guaranteed minimum bandwidth (GMB) of port throughput.

The SDN controller of the switch based on the recommendations received from the AF can either assign traffic connections to HP queue Q8 in order to carry time sensitive requests or to the LP queue (Q3) in case latency constraints are relaxed. To evaluate the performance of these two queuing policies we classify UE connections in two types: HP UEs assigned to the Q8 and LP UEs assigned to Q3. Other configurations can be also supported based on the scenario and traffic type. In total, 10 UEs were emulated by the

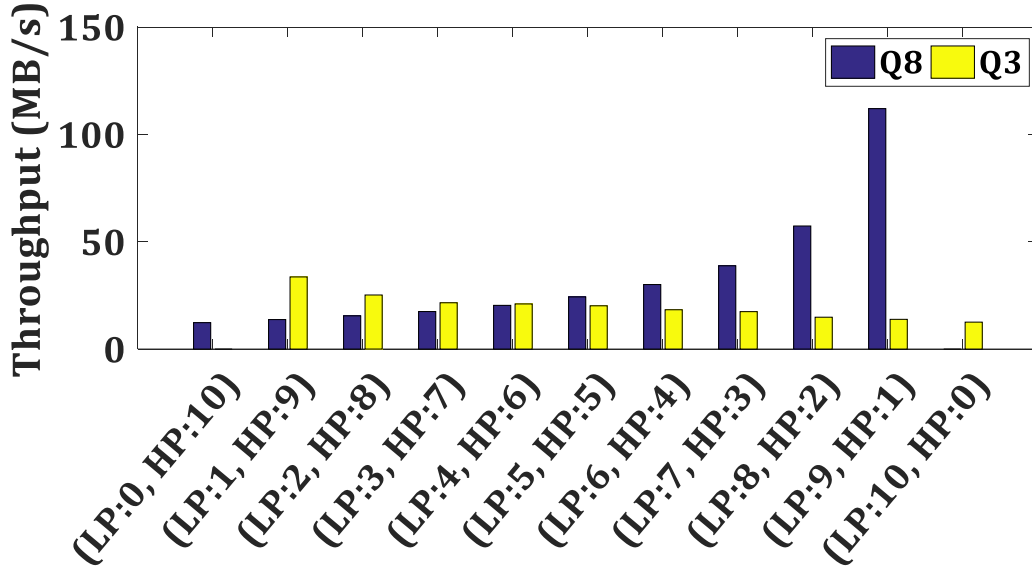


Figure 5. 9: Average throughput of the high and low priority queues.

UERANSIM, with all possible Q8 and Q3 combinations, and their achieved throughput was recorded. The experiment was repeated 10 times to increase statistical validity.

Figure 5. 9 shows the average throughput of the two UE types. The throughput of the Q8 queue connections follows a power trendline curve with negative exponent, that is well aligned with the relevant theoretical models [40]. However, the throughput of the Q3 queue connections, follow a more linear response. This is due to that, when there are Q8 queue connections in the system, only a small percentage of the switch port available bandwidth (30%) is allocated to the Q3 queue, that is shared between the UEs inside the queue. It should be noted that network dimensioning and optimized capacity allocation per queue is out of the scope of the present study.

5.6.3. Cost and Charging Functions

Taking into consideration the values obtained from **Figure 5. 8.**, we can approximate the relation of a node's CPU consumption (CPU_s) with the traffic q that traverses the node using a linear function:

$$CPU_s(q) = a_s q + \beta_s \quad (20)$$

Table 5. 3

CPU COST FUNCTION APPROXIMATION RESULTS OF FIGURE 5. 8

VM	a_s	β_s	RMSE
Small	0.07	0.05	0.9893
Medium	0.03	0.06	0.9573
Large	0.02	0.05	0.9873
X Large	0.008	0.03	0.9545

Where a_s, β_s depend on the type of the node (depending on whether it connects to a small, medium, large or xlarge VM). **Table 5. 3** depicts the parameters a_s, β_s of the 4 types of VMs, with the Root-Mean-Square Error (RMSE) of the fitting curve.

From **Figure 5. 9** we can extract the relationship of the throughput and the traffic (q) that traverses the HP (Q8) and LP (Q3) queues. As before, we can have the relations:

Table 5. 4:

NETWORK COST FUNCTION APPROXIMATION RESULTS OF FIGURE 5. 9

Queues	a_ℓ	β_ℓ	RMSE
HP (Q8)	111.75	0.954	0.9999
LP (Q3)	34.133	0.392	0.9749

$$NET_\ell(q) = \frac{D}{Throughput_\ell} = \frac{D}{a_\ell q^{-\beta_\ell}}, \ell = HP, LP \quad (21)$$

Where D is the size of the requested file and a_ℓ, β_ℓ are positive constants that depend on the type of queuing policy adopted. **Table 5. 4** shows the approximation of these parameters for the 2 types of queues together with their approximation error (RMSE).

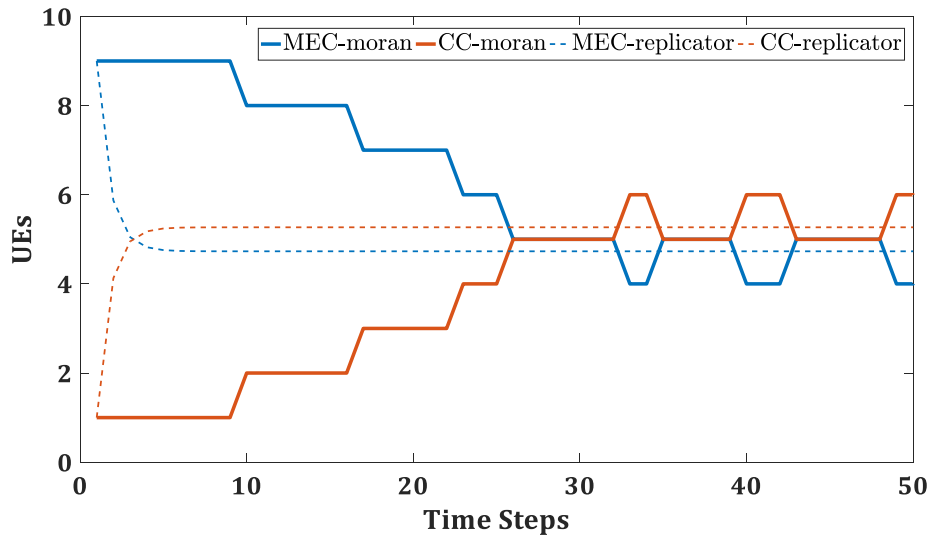
5.7. Numerical Results

5.7.1. Simulation with Experimental Values

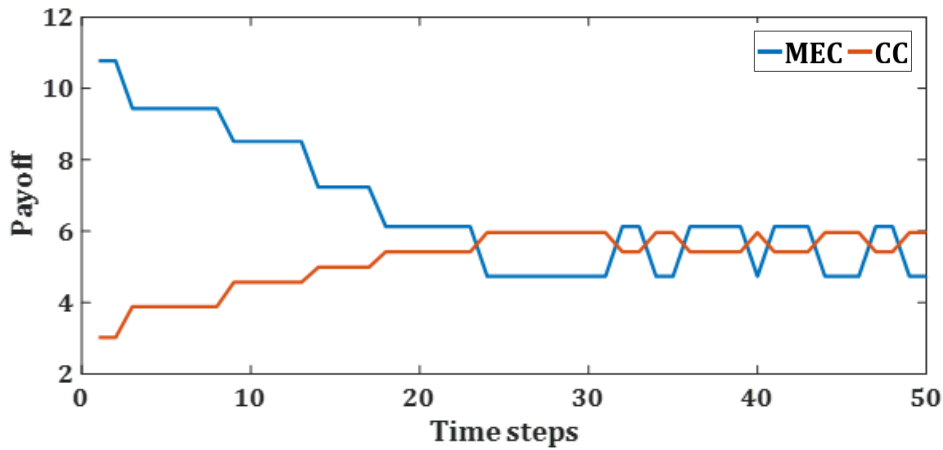
This section presents simulation results to validate our theoretical findings and evaluate the performance of the proposed model. To achieve this, we assumed a system with a UPF directing traffic to a local MEC and the CC. The local MEC is lightweight (hosting a small VM) whereas the CC has higher processing capabilities (hosting a medium VM). The UEs that choose the CL-AS strategy are served by the local MEC, whereas the others that choose the RL-AS are served by the CC. Inside the UPF, the transport level marking of QoS Flows to IP-Flows is equivalent to the JL323A Aruba 2930M switch that we measured before. The QoS Flows marked with CL-AS were served by the Q8 queue of the UPFs, whereas the RL-AS Flows were served by the Q3 Queue.

From Eq. (6), (7), (14) and (15) we calculate $\rho_{HP} = 0.2016$, $\rho_{LP} = 0.1509$, $g_1 = 1,2777$ and $g_{N-1} = -0.7359$. This means that if the population of UEs was infinite, the system would have a stable interior equilibrium where both strategies would gain the same payoff, so they would co-exist (Eq. (19)). In case of only 10 UEs the system will oscillate around this equilibrium point. According to the replicator equation the interior equilibrium point for an infinite population is $x^* = 0.4731$. It is worth noting that this is in agreement with the numerical implementation: if we calculate the payoff difference h_i for all the possible states of the system of 10 UEs we find that the state that minimizes this difference is when 5 UEs choose MEC and 5 CC.

Figure 5. 10 shows the results of our proposed scheme. As we observe in **Figure 5. 10 (a)** the system oscillates around the theoretical state given by replicator equation (5UEs in MEC, 5 in CC), and succeeds in finding the stable equilibrium, where the minimization of the payoff difference happens. **Figure 5. 10 (b)** depicts how the payoffs of the two strategies evolve in time, with respect to the initial state. We see that in the



(a)



(b)

Figure 5. 10: (a) Evolution of the number of UEs inside the population (b) Evolution of the payoff of the two strategies in relation with the initial state.

equilibrium the two strategies obtain the same payoff, indicating the fairness of the algorithm. The impact of the proposed scheme in the compute and network resources is shown in **Figure 5. 11**. It is evident that the algorithm reduces effectively the total compute resources by approximately 25%, while maximizing the average UE throughput.

As it was expected, the results indicate, that some UEs will eventually be served by the CC, in order to minimize the CPU consumption and thus the power consumption of the servers used to carry the request. **Figure 5. 12** depicts the dependence of the interior equilibrium from the compute parameters. Specifically, we can see how the processing capabilities of the system affect the interior equilibrium. As the remote processing capabilities increase, the interior equilibrium is pushed towards the CC strategy that in this case is more favorable than the MEC strategy (the fixation probability of the MEC strategy decreases, and that of the CC strategy increases).

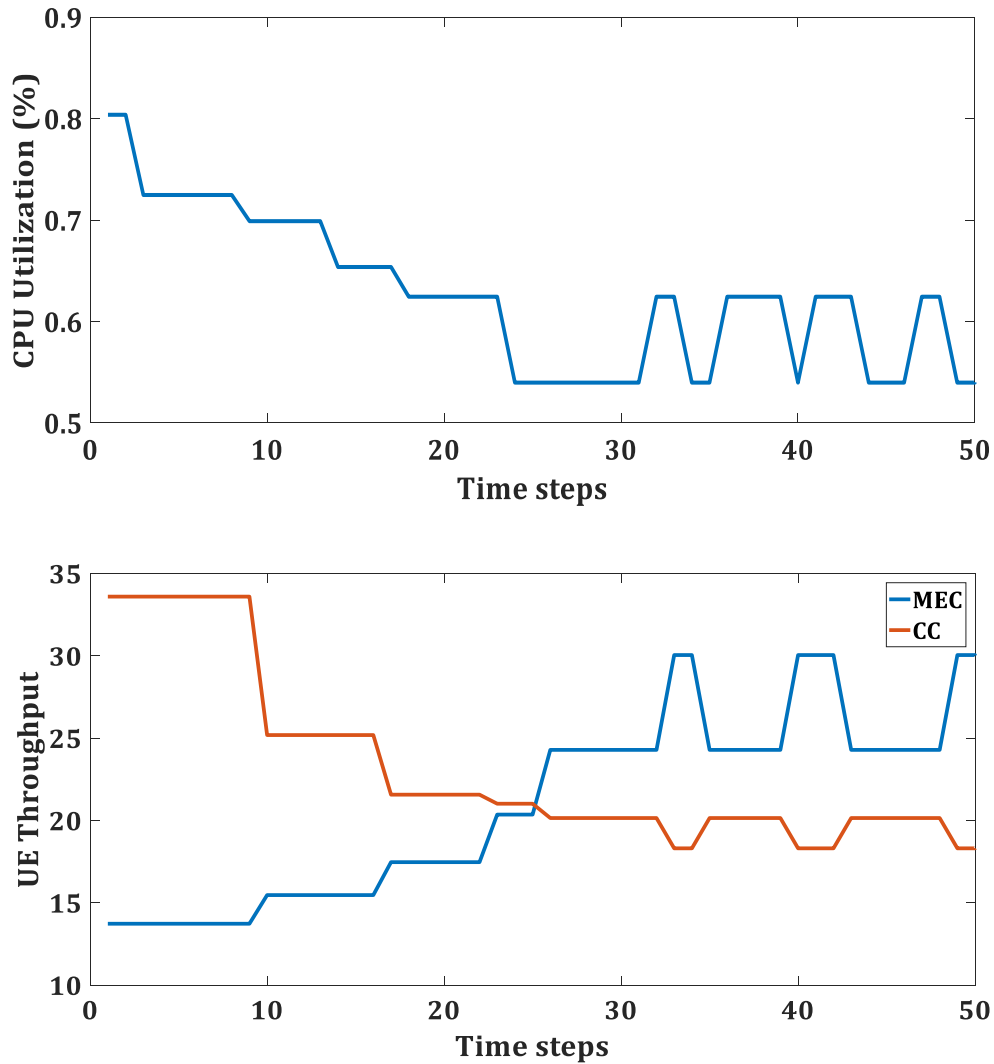


Figure 5. 11: Impact of the proposed scheme on the compute (MEC and CC) and network resources.

To validate the performance of our proposal, we compared our results with a baseline centralized algorithm having global knowledge of the system. In this baseline approach, the payoff function is known in advanced. The algorithm iterates through all possible states of the system and chooses the one that minimizes the selected cost function. The cost function is identical to that used by the EGT model in Section 5.5(B). The objective of the algorithm is to maximize the average payoff of the two strategies, while guaranteeing throughput fairness for the UEs. In order to evaluate how compute and network resources impact the allocation of the UEs to the different ASs, we run the algorithm for three different scenarios. First, we set the weighted parameter h_n of the payoff function to zero, so that the algorithm minimizes the total CPU consumption. For the second scenario, the algorithm tries to maximize the throughput of the UEs (h_c parameter is set to zero), while taking into consideration the fairness between them. The results are shown in **Figure 5. 13**. The CPU-minimization algorithm converges to the state where all UEs are served at the CC, in this state the CPU consumption is minimized (40% of the total CPU used) and all UEs gain the same throughput. However, the queue of the switch that serves the CC UEs is overloaded, resulting in poor average

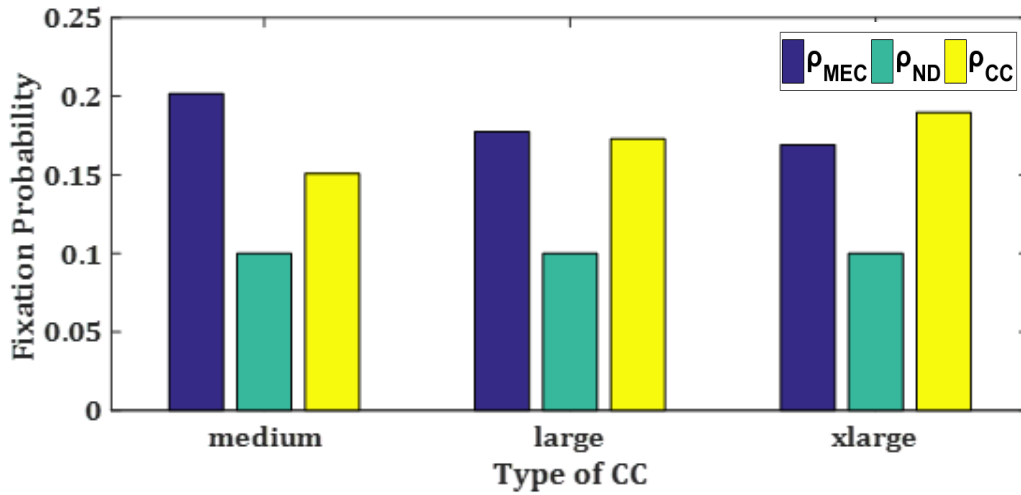


Figure 5. 12: Dependence of the interior equilibrium from the processing capabilities of the CC.

throughput for the UEs (12,5 MB/s). The second algorithm (Throughput Fairness) iterates through all possible states of the system and picks the one that maximizes the throughput of the UEs while securing fairness among them. In our example, this happens when 6 UEs are served by the MEC and 4 by the CC. In this case, the UE's throughput increases, achieving almost equal throughput for each strategy, but so is the system's CPU consumption. Both costs were equally considered for the third case. The centralized approach provides the lower combined network and compute cost. The fairness of the method is provided through an upper limit for the accepted difference between the costs of the two strategies. The lower the limit, the more fair the algorithm becomes. Finally, the proposed scheme that is based on EGT, is also depicted in **Figure 5. 13**. The performance of the EGT scheme is similar to that achieved through the centralised optimization approach, and can provide even better fairness among the UEs, without the requirement to provide explicit limits. This is due to the nature of the replicator equation that terminates the evolution, when every agent in the population has no incentive to deviate from its strategy. Note that the most important advantage of the proposed scheme is its distributed nature that eliminates the need of global system knowledge. Centralized optimization requires that the cost function is already known for all the possible states of the system, something that may not be feasible in dynamic environments. On the other hand, the payoffs are provided to the AF at each repetition of the algorithm. The strategy adaptation process in the proposed EGT-based algorithm does not rely on the knowledge of the strategy selection of the other players. For the evolution a UE requires a random matching with an opponent. Therefore, the amount of information exchange is reduced. The Afs strategies are randomly matches and stop the evolution process, if all payoffs are equal or differ by a small quantity.

Finally, it is important to test the validity of our proposed scheme in different traffic scenarios. In order to do this, we implemented the algorithm for different UE requested data rates (R), as shown in **Figure 5. 14**. From the results we see that for UE requested data rates above 30MB/s, 5 UE mitigate to the CC strategy. The system in this case is over-loaded and cannot offer the requested rate to the UEs. In this case, our scheme

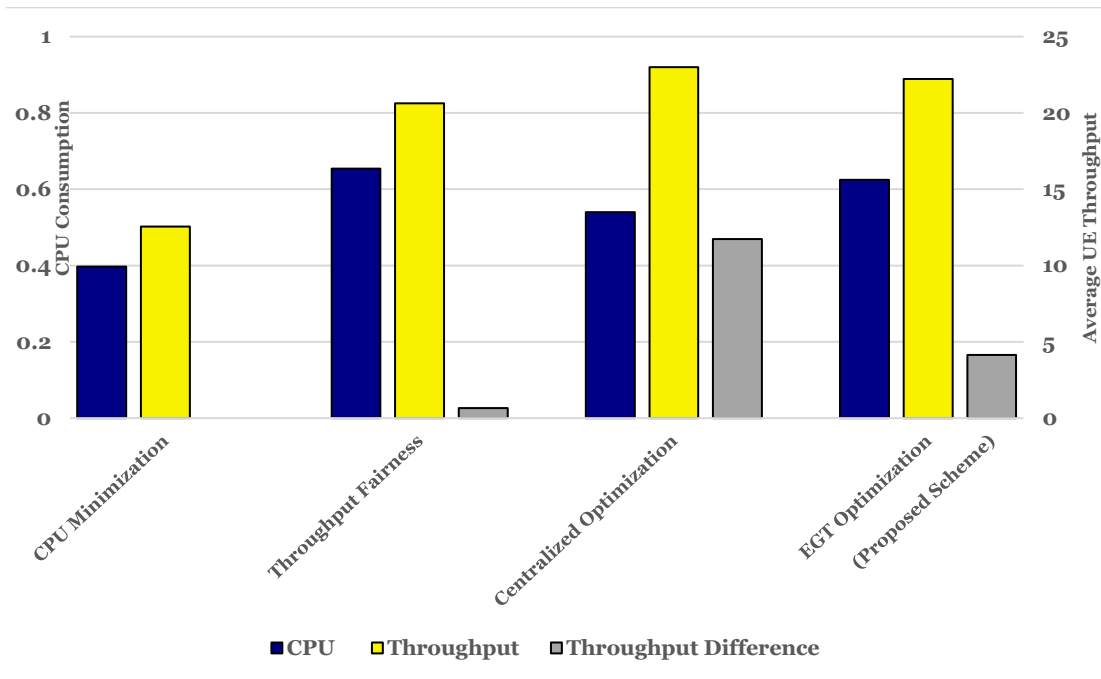


Figure 5. 13: Performance of the proposed scheme versus central baseline algorithm.

tries to offer fairness among the UEs, while minimizing the CPU consumption. As the requested data rate decreases, more UEs can choose the CC strategy, in order to minimize further the CPU consumption.

5.7.2. Extended simulation with multiple AFs

It is important to validate the behavior of the proposed scheme under the existence of multiple AF instances. For this purpose, we divided the population of the UEs in two groups, with each group allocated to one AF. All UEs are initially served at the MEC, when one random UE in each group switches to the CC. The game begins, and at each timestep one UE updates its strategy by observing another UE in the population as in the previous case. The difference here occurs due to that the population of the UEs is not well-mixed, as in the case of classical Moran process. The UEs that belong to different AFs cannot be matched. The spatial relationship of the UEs (whether they belong to the same AF or not) is depicted with the help of a weight matrix $W^{N \times N}$, where w_{ij} is the probability that j UE will be selected as opponent of i UE. As it is obvious $w_{ij} = 0$ if i, j belong to different AFs. **Figure 5. 15** compares the evolutionary process for the cases of one and two AFs. The payoff function is constructed with the experimental values as before. The scenario with one AF converges quickly to the interior equilibrium that was predicted from the analysis of the single-population replicator equation. The important point to note for the case of two AFs is that although it takes more time to converge, it converges to the exact equilibrium (it doesn't oscillate around it). The explanation behind this comes from the stability analysis of the multi-population replicator equation. As it was mentioned in section 5.5(D) each group of UEs will ultimately be redirected to one of the two strategies. When this happens, there will be no mutant, i.e. a UE with different strategy, in each population for the update process to continue, so the evolutionary process will end, and no oscillations will emerge. Finally, the results confirm that the total compute and network resource consumption

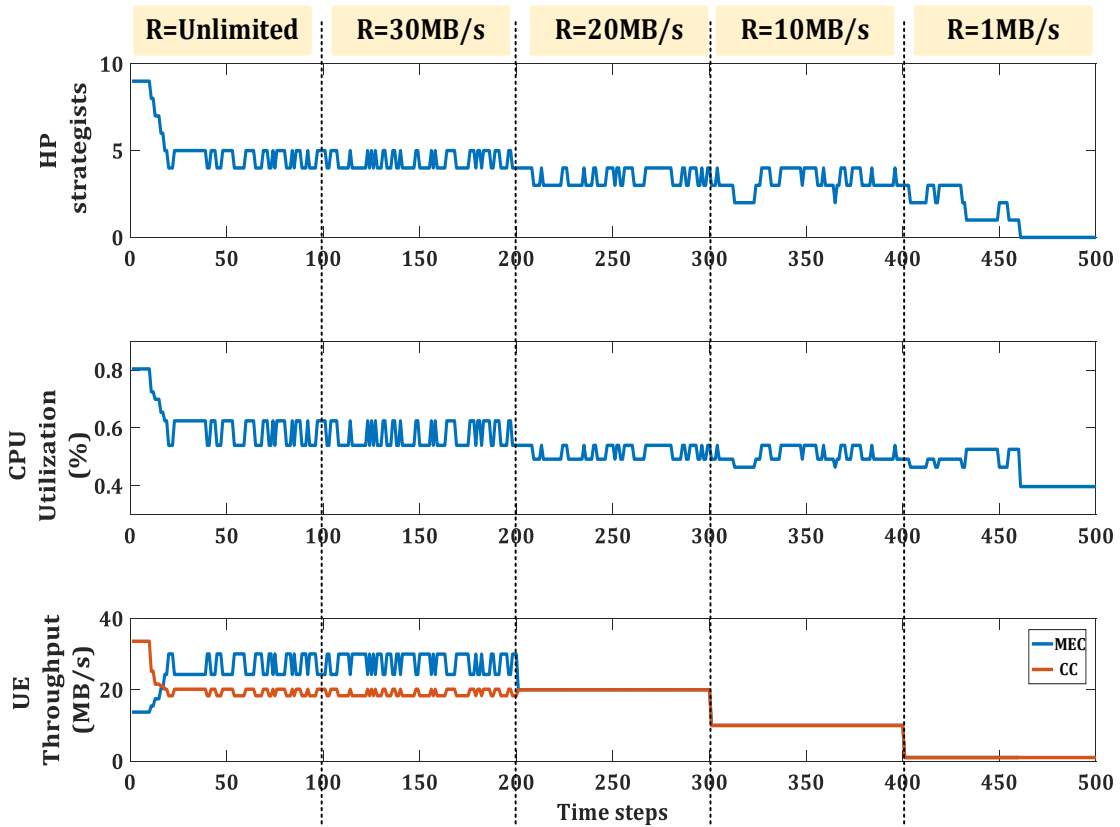


Figure 5.14: Performance of the proposed scheme in different traffic scenarios.

is reduced (**Figure 5.15 (c) and (d)**), as in the case of a well-mixed population of UEs. The same procedure applies to more AFs, by altering the weight matrix, and increasing the number of the UE population.

5.8. Summary

6G systems are expected to operate in a fully distributed manner allowing applications to directly intervene in the system decision-making processes. To address this need, this paper proposes a novel scheme based on EGT that allows AFs to control UE applications. This will enable them to dynamically select their service strategy, in order to optimize the QoE delivered to the end users minimizing at the same time charging related costs. The scheme has been designed using realistic constraints and cost functions and was evaluated over a real cloud-based 5G testbed. The overall scheme can minimize the total compute resources by approximately 25%, while offering the highest available throughput to each UE in a fair matter.

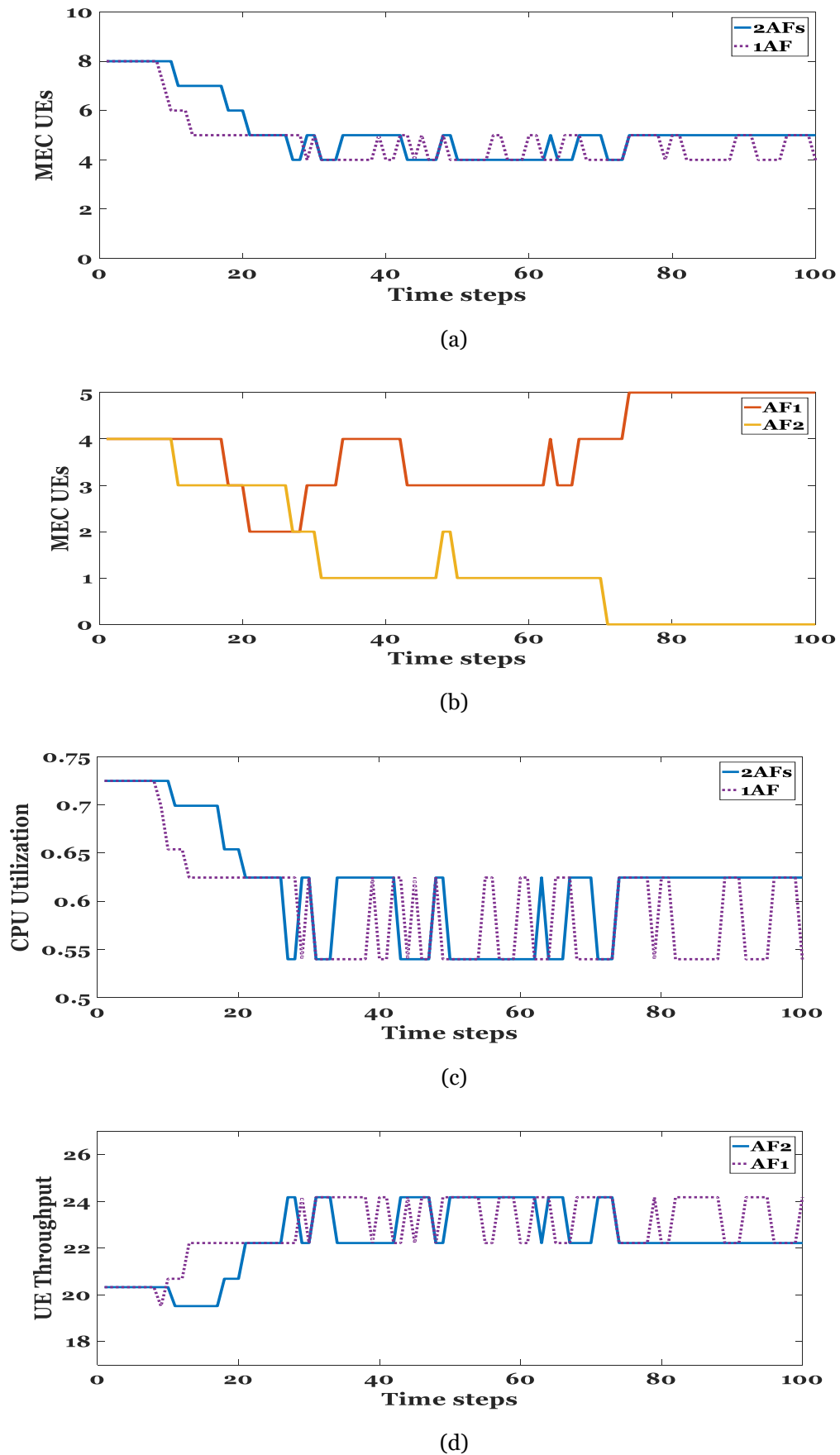


Figure 5.15: (a) Evolution of the UEs population under the control of one or two AFs. (b) Evolution of each AF's population of UEs that are served at MEC facility. Impact of the proposed scheme on (c) the compute (MEC and CC) and (d) network resources.

Appendix

A. Derivation of Fokker-Planck Equation

The Taylor expansion of the probability density around t and x is respectively:

$$d\left(x, t + \frac{1}{N}\right) \approx d(x, t) + \frac{1}{N} \frac{\partial}{\partial t} d(x, t) \quad (22)$$

$$d\left(x + \frac{1}{N}, t\right) \approx d(x, t) + \frac{1}{N} \frac{\partial}{\partial x} d(x, t) + \frac{1}{2N^2} \frac{\partial^2}{\partial x^2} d(x, t) \quad (23)$$

The Taylor expansion of the transition probabilities around x is:

$$P^\pm\left(x + \frac{1}{N}\right) \approx P^\pm(x) \pm \frac{1}{N} \frac{\partial}{\partial x} P^\pm(x) + \frac{1}{2N^2} \frac{\partial^2}{\partial x^2} P^\pm(x) \quad (24)$$

Substituting the above expressions in Eq (9), we extract the following equation:

$$\begin{aligned} \frac{\partial}{\partial t} d(x, t) = & \left(-d(x, t) \frac{\partial}{\partial x} P^+(x) - P^+(x) \frac{\partial}{\partial x} d(x, t) + d(x, t) \frac{\partial}{\partial x} P^-(x) + P^-(x) \frac{\partial}{\partial x} d(x, t) \right) \\ & + \frac{1}{2N} \left(d(x, t) \frac{\partial^2}{\partial x^2} P^+(x) + 2 \frac{\partial}{\partial x} P^+(x) \frac{\partial}{\partial x} d(x, t) + P^+(x) \frac{\partial^2}{\partial x^2} d(x, t) \right. \\ & \left. + d(x, t) \frac{\partial^2}{\partial x^2} P^-(x) + 2 \frac{\partial}{\partial x} P^-(x) \frac{\partial}{\partial x} d(x, t) + P^-(x) \frac{\partial^2}{\partial x^2} d(x, t) \right) \\ & + O\left(\frac{1}{N^3}\right) + \dots \end{aligned}$$

That can be written in the following form by neglecting the higher order terms of N :

$$\begin{aligned} & \frac{\partial}{\partial t} d(x, t) \\ & = -\frac{\partial}{\partial x} \left[-d(x, t) (P^+(x) - P^-(x)) \right] \\ & + \frac{1}{2} \frac{\partial^2}{\partial x^2} \left(d(x, t) \frac{P^+(x) + P^-(x)}{N} \right) + O\left(\frac{1}{N^3}\right) \\ & + \dots \end{aligned} \quad (25)$$

B. Derivation of Replicator Equation

From Eq. (11), for large populations ($N \gg 1$):

$$\dot{x} = \mu(x) + \sigma^2(x)\xi \quad (26)$$

Substituting $\mu(x) = (P^+(x) - P^-(x))$ and $\sigma^2(x) = \frac{P^+(x)+P^-(x)}{N}$ into (26) results the following equation:

$$\begin{aligned} \dot{x} &= (P^+(x) - P^-(x)) + \lim_{N \rightarrow \infty} \frac{P^+(x) + P^-(x)}{N} \\ &= (P^+(x) - P^-(x)) \end{aligned} \quad (27)$$

(27) after some calculations can be written in the following form:

$$\begin{aligned} \dot{x} &= x(1-x) \left(\frac{1}{1 + e^{-w(\pi_{CL}(x) - \pi_{RL}(x))}} - \frac{1}{1 + e^{+w(\pi_{CL}(x) - \pi_{RL}(x))}} \right) \\ &= x(1-x) \left(\frac{e^{+w(\pi_{CL}(x) - \pi_{RL}(x))} - e^{-w(\pi_{CL}(x) - \pi_{RL}(x))}}{2 + e^{w(\pi_{CL}(x) - \pi_{RL}(x))} + e^{-w(\pi_{CL}(x) - \pi_{RL}(x))}} \right) \\ &= x(1-x) \left(\frac{\sinh w(\pi_{CL}(x) - \pi_{RL}(x))}{\cosh w(\pi_{CL}(x) - \pi_{RL}(x)) + 1} \right) \\ &= x(1-x) \tanh 2w(\pi_{CL}(x) - \pi_{RL}(x)) \end{aligned}$$

Under weak selection ($w \ll 1$), we can substitute $\tanh 2w(\pi_{CL}(x) - \pi_{RL}(x))$ with its Taylor expansion and thus extracting the replicator equation (neglecting the higher order terms of w):

$$\begin{aligned} \dot{x} &= x(1-x)(2w(\pi_{CL}(x) - \pi_{RL}(x)) \\ &\quad + O(w^3)) \end{aligned} \quad (28)$$

C. Calculation of $\pi_{HP}(\mathbf{x}(t))$ and $\pi_{LP}(\mathbf{x}(t))$ derivates

In (21) and (22) equations, we express traffic (g) in terms of x :

$$q_{CL} = Nx \quad (29)$$

$$q_{RL} = N(1-x) \quad (30)$$

Now the derivatives of the payoffs can be calculated:

$$\begin{aligned}
 \frac{d\pi_{CL}(\mathbf{x}(t))}{dx} &= \frac{d}{dx} \left(\frac{h_c + h_n}{h_c \cdot (a_{CL}Nx + \beta_{CL}) + h_n \cdot \frac{D}{a_{CL}(Nx)^{-\beta_{CL}}}} \right) \\
 &= -(h_c + h_n) \cdot \frac{h_c a_{CL}N + h_n \cdot \frac{D \cdot \beta_{CL} \cdot (Nx)^{\beta_{CL}-1}}{a_{CL}}}{\left[h_c \cdot (a_{CL}Nx + \beta_{CL}) + h_n \cdot \frac{D}{a_{CL}(Nx)^{-\beta_{CL}}} \right]^2} < 0 \\
 \frac{d\pi_{RL}(\mathbf{x}(t))}{dx} &= \frac{d}{dx} \left(\frac{h_c + h_n}{h_c \cdot (a_{RL}N(1-x) + \beta_{RL}) + h_n \cdot \frac{D}{a_{RL}(N(1-x))^{-\beta_{RL}}}} \right) > 0 \\
 &= \frac{(h_c + h_n) \left[h_c a_{RL}N + h_n \cdot \frac{D(N(1-x))^{\beta_{RL}-1}}{a_{RL}} \right]}{\left[h_c \cdot (a_{RL}N(1-x) + \beta_{RL}) + h_n \cdot \frac{D}{a_{RL}(N(1-x))^{-\beta_{RL}}} \right]^2}
 \end{aligned}$$

Thus,

$$\frac{d\pi_{CL}(\mathbf{x}(t))}{dx} - \frac{d\pi_{RL}(\mathbf{x}(t))}{dx} < 0$$

References

- [1] Tzanakaki et al., "Wireless-Optical Network Convergence: Enabling the 5G Architecture to Support Operational and End-User Services", *IEEE Comms. Mag*, vol. 55, no. 10, pp. 184-192.2017.
- [2] ETSI TS 123 501 V16.11.0 (2021-12). 5G;. System architecture for the 5G System (5GS).
- [3] H. Liu, F. Eldarrat, H. Alqahtani, A. Reznik, X. de Foy and Y. Zhang, "Mobile Edge Cloud System: Architectures, Challenges, and Approaches", *IEEE Systems Journal*, vol. 12, no. 3, pp. 2495-2508, 2018. Available: 10.1109/jsyst.2017.2654119.
- [4] "Cloud RAN and MEC: A Perfect Pairing", Etsi.org, 2019. [Online]. Available: https://www.etsi.org/images/files/ETSIWhitePapers/etsi_wp23_MEC_and_CRAN_e d1_FINAL.pdf.
- [5] Mataj, Enida, "Network slicing and QoS in 5G systems and their impact on the MAC layer." Master Thesis, Department of Electronics and Telecommunications, Politecnico di Torino, 2020. Accessed 2021-04-05.
- [6] Stefan Rommer, Peter Hedman, Magnus Olsson, Lars Frid, Shabnam Sultana, Catherine Mulligan - "5G Core Networks_ Powering Digitalization", *Academic Press*, 2019.
- [7] Tzanakaki, A. Manolopoulos, M. Anastasopoulos, and D. Simenidou, "Optical Networking in Support of User Plane Functions in 5G Systems and Beyond," in *Photonics in Switching and Computing 2021*, W2B.3.
- [8] ETSI TS 129 517 V16.5.0 (01/2021). 5G; 5G System; Application Function Event Exposure Service; Stage 3. Accessed 2022-03-01

- [9] Devopedia, "5G Quality of Service." Version 3, April 8 2021. Accessed 2022-01-18. <https://devopedia.org/5g-quality-of-service>
- [10] ETSI TS 138 300 V16.4.0 (2021c). 5G; NR; NR and NG-RAN Overall description; Stage-2.". Accessed 2021-04-05.
- [11] H2020 Project 5G-COMplete, Deliverable D2.1: Initial report on 5G-COMplete network architecture, interfaces and supported functions.
- [12] ETSI TS 129 513 V16.6.0 (01/2021d). 5G; 5G System; Policy and Charging Control signalling flows and QoS parameter mapping; Stage 3.
- [13] "IEEE Standard for Packet-based Fronthaul Transport Networks," in *IEEE Std 1914.1-2019*, vol., no., pp.1-94, 21 April 2020, doi: 10.1109/IEEESTD.2020.9079731.
- [14] S. Sevgican, M. Turan, K. Gökarslan, H. B. Yilmaz and T. Tugcu, "Intelligent network data analytics function in 5G cellular networks using machine learning," in *Journal of Communications and Networks*, vol. 22, no. 3, pp. 269-280, June 2020, doi: 10.1109/JCN.2020.000019.
- [15] ETSI TS 123 288 V16.10.0. 5G (01/2022); Architecture enhancements for 5G System (5GS) to support network and data analytics services. Accessed 2022-03-01.
- [16] K. Katsalis, N. Nikaein, E. Schiller, A. Ksentini and T. Braun, "Network Slices toward 5G Communications: Slicing the LTE Network," in *IEEE Communications Magazine*, vol. 55, no. 8, pp. 146-154, Aug. 2017, doi: 10.1109/MCOM.2017.1600936.
- [17] M. Irazabal, E. Lopez-Aguilera and I. Demirkol, "Active Queue Management as Quality of Service Enabler for 5G Networks," in *proc. of EuCNC*, 2019, pp. 421-426, doi: 10.1109/EuCNC.2019.8802027.
- [18] F. Akyildiz et al., "xStream: A New Platform Enabling Communication Between Applications and the 5G Network," in *Proc. of IEEE Globecom Workshops*, 2018, pp. 1-6, doi: 10.1109/GLOCOMW.2018.8644183.
- [19] M. Bosk et al., "Using 5G QoS Mechanisms to Achieve QoE-Aware Resource Allocation", *In Proc. of 17th International Conference on Network and Service Management 2021*, pp. 283-291, doi: 10.23919/CNSM52442.2021.9615557.
- [20] Schwarzmann, Susanna & Marquezan, Clarissa & Bosk, Marcin & Liu, Huiran & Trivisonno, Riccardo & Zinner, Thomas, "Estimating Video Streaming QoE in the 5G Architecture Using Machine Learning. ", 2019. 7-12. 10.1145/3349611.3355547.
- [21] R. Vidhya, P. Karthik and S. Jamadagni, "Anticipatory QoE Mechanisms for 5G Data Analytics," *In Proc. of International Conference on COMMunication Systems & NETWORKS (COMSNETS)*, 2020, pp. 523-526, doi: 10.1109/COMSNETS48256.2020.9027358.
- [22] S. Schwarzmann, C. C. Marquezan, R. Trivisonno, S. Nakajima and T. Zinner, "Accuracy vs. Cost Trade-off for Machine Learning Based QoE Estimation in 5G Networks," in *IEEE International Conference on Communications (ICC)*, 2020, pp. 1-6, doi: 10.1109/ICC40277.2020.9148685.
- [23] P. Naghavi, S. Hamed Rastegar, V. Shah-Mansouri and H. Kebriaei, "Learning RAT Selection Game in 5G Heterogeneous Networks," in *IEEE Wireless Communications Letters*, vol. 5, no. 1, pp. 52-55, Feb. 2016, doi: 10.1109/LWC.2015.2495123.
- [24] Alevizaki, Victoria-Maria & Anastasopoulos, Markos & Tzanakaki, Anna & Simeonidou, Dimitra. "Joint Fronthaul Optimization and SDN Controller Placement in Dynamic 5G Networks", in *proc. of Optical Network Design and Modelling 2019*. 10.1007/978-3-030-38085-4_16.

- [25] Alevizaki, VM., Anastasopoulos, M., Tzanakaki, A. et al., "Adaptive FH optimization in MEC-assisted 5G environments.", *Photon Netw Commun* 40, 209–220 (2020). <https://doi.org/10.1007/s11107-020-00906-8>
- [26] V. M. Alevizaki, M. Anastasopoulos, A. Tzanakaki and D. Simeonidou, "Dynamic Selection of User Plane Function in 5G Environments" *In Proc. of Optical Network Design and Modelling*, 2021
- [27] V. M. Alevizaki, A. I. Manolopoulos, M. Anastasopoulos and A. Tzanakaki, "Dynamic User Plane Function Allocation in 5G Networks enabled by Optical Network Nodes," *2021 European Conference on Optical Communication (ECOC)*, 2021, pp. 1-4, doi: 10.1109/ECOC52684.2021.9606154.
- [28] J. Weibull, *Evolutionary game theory*. Cambridge, Mass.: MIT Press, 2004.
- [29] A Traulsen, M. A. Nowak, J. M. Pacheco, "Stochastic dynamics of invasion and fixation", *Phys Rev E Stat Nonlin Soft Matter Phys*, Jul. 2006 ;74(1 Pt 1):011909. doi: 10.1103/PhysRevE.74.011909. Epub 2006 Jul 17. PMID: 16907129; PMCID: PMC2904085.
- [30] William Norman, "Evolutionary Game Dynamics and the Moran Model. ", Uppsala Universitet, U.U.D.M. Project Report 2020:39
- [31] T. Antal, I. Scheuring, "Fixation of Strategies for an Evolutionary Game in Finite Populations", *Bull. Math. Biol*, vol. 68, pp. 1923–1944, 2006.
- [32] Risken, H. "Fokker-Planck Equation. ", in *The Fokker-Planck Equation. Springer Series in Synergetics*, vol 18. Springer, Berlin, Heidelberg, 1996. https://doi.org/10.1007/978-3-642-61544-3_4
- [33] Dean Foster, Peyton Young, "Stochastic evolutionary game dynamics", in *Theoretical Population Biology*, Volume 38, Issue 2, 1990, Pages 219-232, ISSN 0040-5809, [https://doi.org/10.1016/0040-5809\(90\)90011-J](https://doi.org/10.1016/0040-5809(90)90011-J).
- [34] L. J. S. Allen, "An introduction to stochastic processes with applications to biology", Upper Saddle River, N.J., Pearson/Prentice Hall, 2003.
- [35] Taylor, D. Fudenberg, A. Sasaki, M. Nowak, "Evolutionary Game Dynamics in Finite Populations", *Bulletin of mathematical biology* vol. 66, pp. 1621-44, 2004.
- [36] UERANSIM, [Online]. Available: <https://github.com/aligungr/ UERANSIM>
- [37] free5GC, [Online]. Available: <https://www.free5gc.org/>
- [38] Prometheus, [Online], Available: <https://prometheus.io/>
- [39] Grafana, [Online], Available: <https://grafana.com/>
- [40] Bertsekas and R. Gallager, "Data Networks", 2nd Edition, Prentice Hall, Englewood Cliffs, N.J., 1992

Chapter 6

Network Slicing and Orchestration

Contents

6.1.	Introduction	96
6.2.	NFV-MANO Fundamentals.....	98
6.2.1.	Network Function Virtualization	98
6.2.2.	Network slicing	99
6.3.	Problem Statement	101
6.4.	Experimental Setup	102
6.4.1.	5G Platform Description	102
6.4.2.	Orchestration Platform Overview.....	102
6.5.	Implementation.....	105
6.5.1.	Creation of Network Descriptors	105
6.5.2.	Slice Deployment and Results.....	108
6.6.	Summary	109
	References.....	109

6.1. Introduction

Future mobile communication networks aim to offer services and applications in the most flexible, adaptable and cost-effective manner as possible. The general interest is shifting from QoS-oriented to QoE oriented strategies, that focus on delivering a sufficient end-user experience. QoE looks further than the QoS key network requirements, at the impact of the network behavior on the end user [1][2][3]. In order to provide effective QoE management, B5G networks aim at a fully softwarized network architecture, where hardware and software programming is used for the design, implementation, deployment, management, monitoring and maintenance of network equipment/components/services [4][6].

To this end networks are transforming to open flexible infrastructures that are efficiently shared. Hardware programmability and NFV are two principal methods that can be utilized for sharing physical infrastructure among various network services. Hardware programmability refers to creating a common hardware platform that can accommodate a variety of different network configurations, making it easier to deploy new services without requiring significant hardware changes. NFV takes this concept one step further by

enabling the migration from traditional network elements (such as routers, switches, and firewalls) to network functions that can be hosted on general-purpose servers as VMs. This allows multiple network functions to be consolidated onto a single physical server, reducing the need for specialized hardware and improving service delivery times. Both concepts enable softwarisation of different parts of the network, such as radio access, core, transport, mobile-edge and central clouds, based on each segment's needs and technical characteristics [8]. Softwarizing the 5G transport network and the RAN will enable efficient coordination between them in support of functionalities such as mobility and load balancing [9]. Most 5G core networks and service plane functions are envisaged to be softwarized and implemented as VNFs on Fog/Cloud Computing environments [5gslicing-64,65], located at any network site. As a result, deploying the essential core functions may be performed through a single click, allowing simple remote upgrades as opposed to installing conventional hardware equipment, which requires cumbersome in advance planning. Finally, softwarization of MEC in 5G promises to decrease the amount of data carried to the 5G core network for processing, provide real-time and application flow information, and efficient usage of the available resources.

Within NFV, a network service is essentially a collection of different VNFs. Each VNF performs a different NF in isolation, with the appropriate virtual resources provided by the virtualization layer. The collection and appropriate combination of multiple VNFs gives an E2E network service. The number of VNFs and their arrangement within the virtual infrastructure are typical for each NS that needs to be deployed. VNFs can be located anywhere in the network, as long as the location of each VNF is known so that they can communicate properly. This enables 5G networks to offer network slicing as-a-service [7], where network operators and developers can swiftly construct isolated application-aware networks and network-aware apps, driven by their business needs. To manage these network slices, a comprehensive MANO system is required. MANO is essential for the successful deployment and operation of 5G slicing, as it provides a unified management platform for the creation, modification, and monitoring of network slices. MANO enables network operators to manage resources and services, such as NFs, computing resources, and storage, in an automated and scalable manner. It plays a crucial role in coordinating the various components of the network slice, such as the core network, the RAN and the transport network. The MANO system also ensures that the network slice adheres to the SLA agreed upon with the customer, ensuring that the network slice meets the required performance and QoS levels.

In this chapter, we present our contributions towards automating the deployment of a sophisticated 5G multi-operator network and E2E provisioning of multiple slices, using an NFV-MANO compliant orchestrator. All of our contributions are experimentally driven, applied to a lab testbed. In this context, a softwarized multi-operator 5G platform hosted in a containerized data centre infrastructure was implemented. The servers within the platform were interconnected with optoelectronic/SDN switches. The main contributions of our work are the following:

- Dynamic instantiation of multi-operator 5G network systems (public/private) with shared 5G core functions
- Provisioning of services with different characteristics
- End-to-end connections with the appropriate QoS Levels
- Private/public network slice terminated at the local MEC/remote facility accordingly.

All the activities carried out in the lab environment were related to projects under Horizon 2020, such as 5G-COMPLETE[20] and 5G-VICTORI[21].

The chapter is structured as follows. First the fundamental components and concepts of NFV and network slicing are presented. The problem statement follows, along with the experimental setup for the realization of a 5G environment testbed are described. Finally, we present the procedure that was followed for the automation and creation of 5G network slices and show the experimental results.

6.2. NFV-MANO Fundamentals

6.2.1. Network Function Virtualization

A significant innovation associated with 5G cellular systems is the adoption of the notion of NFV, a concept that enables implementation of network functions using virtualized software, rather than traditional hardware-based appliances. Traditionally, each network function needed its own proprietary device to be deployed. This rendered the networks very rigid and expensive, due to the high cost of the required specialized hardware along with their low reconfigurable features. NFV enables virtualization of NFs, decoupling them from the physical hardware, and enabling their dynamic deployment and allocation on a shared infrastructure based on general-purpose servers. This way, multiple NFs can be hosted in the existing equipment without having to add hardware infrastructure. Changes in the network, like removal or update of a function for all or subset of customers are automated and significantly simplified to operations performed in a single click fashion, making the infrastructure for the service providers (SPs) much more scalable, flexible, manageable and consequently low cost with short deployment time [12].

Future 5G networks will use NFV to guarantee the performance of VNFs, including low latency and failure rates, optimize resource allocation to end users with high QoS and ensure their compatibility with non-VNFs [13]. In order to accomplish the aforementioned advantages, NFV introduces changes to the way network services are provisioned. Since the hardware and software platforms are fully decoupled, they can progress and evolve separately from one another, carrying out diverse tasks at different times. Thus, operators can now introduce new, cutting-edge services using the same hardware platform. Finally, the dynamic scaling of the VNFs allows the network operators to implement services that are specifically designed to meet client needs. [1]

Figure 6. 1 shows the NFV architectural overview. In NFV, there are three primary elements: VNFs, NFV Infrastructure (NFVI) and NFV MANO [12][14][15][16]. Each VNF performs a different NF in isolation. A VNF may consist of several internal components, e.g VMs or containers. For example, a VNF may consists of two VMs, one for the data storage and a second for the main functions. Appropriate virtual resources for the VNFs are provided by the NFVI. The NFVI consists of general-purpose computing hardware and software and could be distributed over several discrete geographical locations. It offers the virtualization layer that is responsible for abstracting and delivering physical resources (compute, storage, network) to support the VNFs. Finally, the MANO of the VNFs is crucial for the rapid and reliable deployment of NSs. [17] NFV-MANO interacts with the VNF and NFVI blocks to manage and coordinate all virtualization-specific duties within the NFV framework. In more detail, it offers virtual machines for the VNFs, sets them up and controls the management of physical resources for the VMs and the lifecycle of the VNFs [12]. NFV MANO is composed of virtualized infrastructure managers (VIMs), VNF

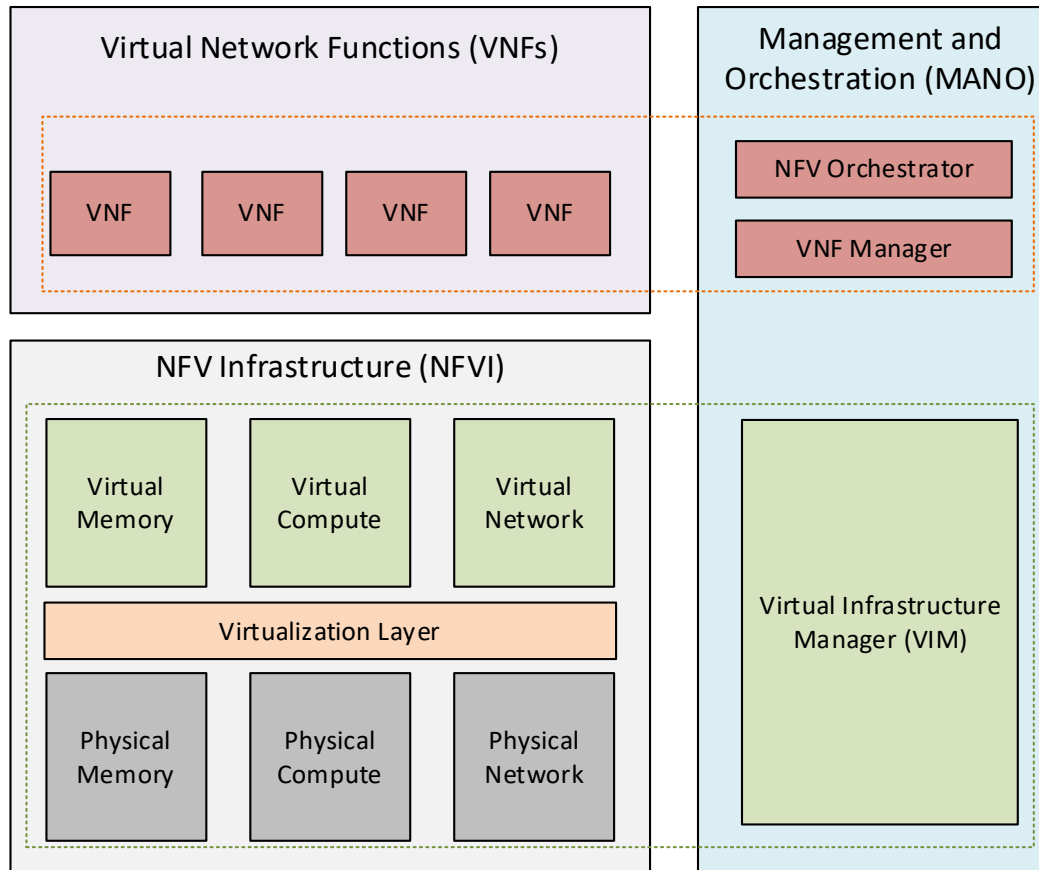


Figure 6. 1: NFV Architectural Framework.

managers and NFV orchestrators. VIMs control and manage the interactions of a VNF with its computing, storage and network resources, and ensures their virtualization. The VNF manager is responsible for managing the entire VNF lifecycle. It is responsible to initialize, query, update and terminate VNF instances. Finally, NFV orchestrators are responsible for orchestrating and managing new network services into a virtual framework, which includes instantiation, policy management, performance measurement and monitoring. Together, these blocks are responsible for deploying and connecting functions and services when they are needed throughout the network.

6.2.2. Network slicing

The network slicing concept was first introduced by the Next Generation Mobile Network (NGMN) [18]. According to NGMN's definition:

“A network slice is an E2E logical network/cloud running on a common underlying (physical or virtual) infrastructure, mutually isolated, with independent control and management that can be created on demand.”

The NGNM divides a slice concept in three basic layers, namely the 5G Service Instance Layer (SIL), the 5G Network Slice Instance (NSI) and the 5G Resource Layer (RL). The SIL provides the services instances (SIs) that are going to be supported. The NSI describes the network requirements for these services, and can be shared across multiple SIs. The NSI may comprise network slice subnet instances (NSSIs), which can be dedicated or shared among several NSIs. Finally, the RL consists of physical and logical resources that are allocated to the slice [28].

Network slicing on 5G softwarized networks is guided by several key principles [5] [21], including automation of network operations for dynamic life-cycle management [20], rapid fault detection mechanisms [22] for high reliability, scalability, and isolation, and programmability for flexible resource allocation and customized service delivery [23]. Network slicing adds another layer of abstraction, with hierarchical abstraction and slice customization possible at all layers of the abstracted network topology. The use of SDN and network resources elasticity also enables value-added AI assisted services [24] and guarantees the desired SLAs for users, regardless of their physical location [25].

A network slice may include components relevant to the access, transport, core and edge networks as well as cross-domain components from various domains under the same or distinct administrations [1]. Depending on the implementation, a network slice can be characterized as hard, if it is completely isolated, or soft if it shares network resources (for example SMF) with other slices. To application providers or distinct vertical sectors without a physical network infrastructure, network slicing can provide radio, cloud, and networking services. By tailoring network operation to customers' needs based on the type of service, it enables service differentiation [19]. For example, a slice could be created for a specific industry, such as healthcare, or for a specific application, such as autonomous vehicles.

The specific characteristics and requirements of a network slice are depicted through two attributes, the Service Slice Type (SST) and the Service Descriptor (SD) feature. The SST is a high-level identifier that represents a specific type of service that the network slice is intended to provide. It provides an abstract representation of the service, such as eMBB, URLLC and mMTC. The SD, on the other hand, is a detailed specification that defines the functional and non-functional requirements of the service, including network functions, network topology, performance metrics, and security policies. It provides a complete description of the slice, which can be used to configure and provision the network slice. The SD is typically composed of several parts, including a slice subnet, network function requirements, performance requirements, security requirements, and charging and billing requirements. The slice subnet refers to the IP address subnet used by the slice, while the functional requirements dictate the types of network functions required. Performance requirements are the non-functional requirements like maximum latency and throughput and security requirements refer to policies like authentication and encryption.

For a unified view of the 5G-slicing concept, the 3GPP has united multiple standard developing organizations (SDOs) to define the concept, use cases, requirements and solutions for MANO of network slices [32]. Based on their efforts, a 3-step MANO lifecycle of 5G slices is established, that comprises of:

1. instantiation, configuration and activation. During this step, all the resources required for the NSI are created and configured. Any other actions that are need to be performed in order for the NSI to be operative and active are also performed in this step.
2. run-time. During this step, the NSI is operational, and can be supervised and monitored. Furthermore, run-time actions, e.g. scaling, can be also performed.

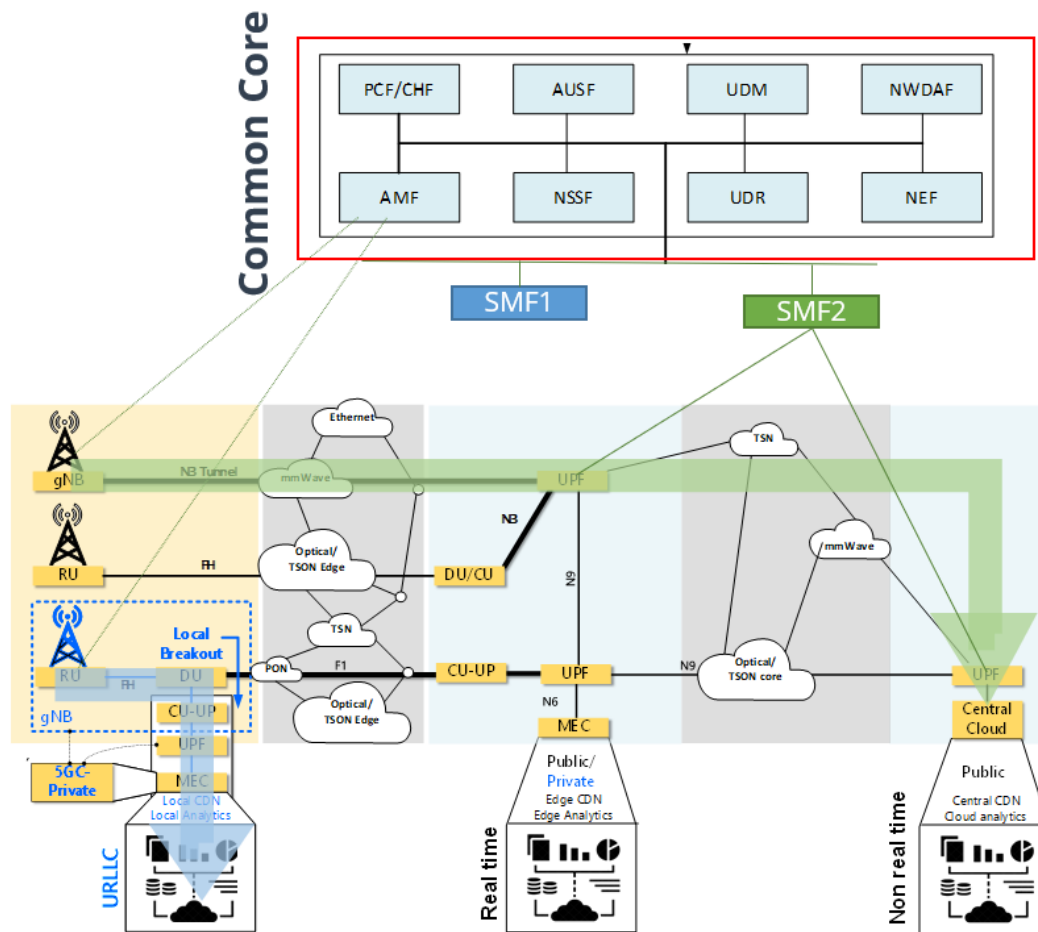


Figure 6. 2: Network topology under consideration.

3. decommissioning. During the decommissioning phase the NSI is deactivated and the resources that were allocated to the NSI are released.

6.3. Problem Statement

A high-level view of a 5G deployment option that was considered is shown in **Figure 6. 2**, where the NG-RAN architecture of 5G is also illustrated. The gNB can be divided into three distinct components: the RU, which manages the physical layer functions of the 5G protocol stack; the DU, which oversees real-time layer 1 (L1) and layer 2 (L2) scheduling tasks; and the CU, which handles non-real-time processing tasks. These components can either be collocated or instantiated in different computing facilities, such as MEC or CC nodes. Additionally, UP elements can be placed close to the network edge to redirect traffic to MEC nodes or deeper within the network to forward traffic to a CC facility. The UP elements communicate with the 5G core through point-to-point interfaces (N2, N3, N4, N6, N9). The 5G Transport Network is a blend of wireless and optical technologies.

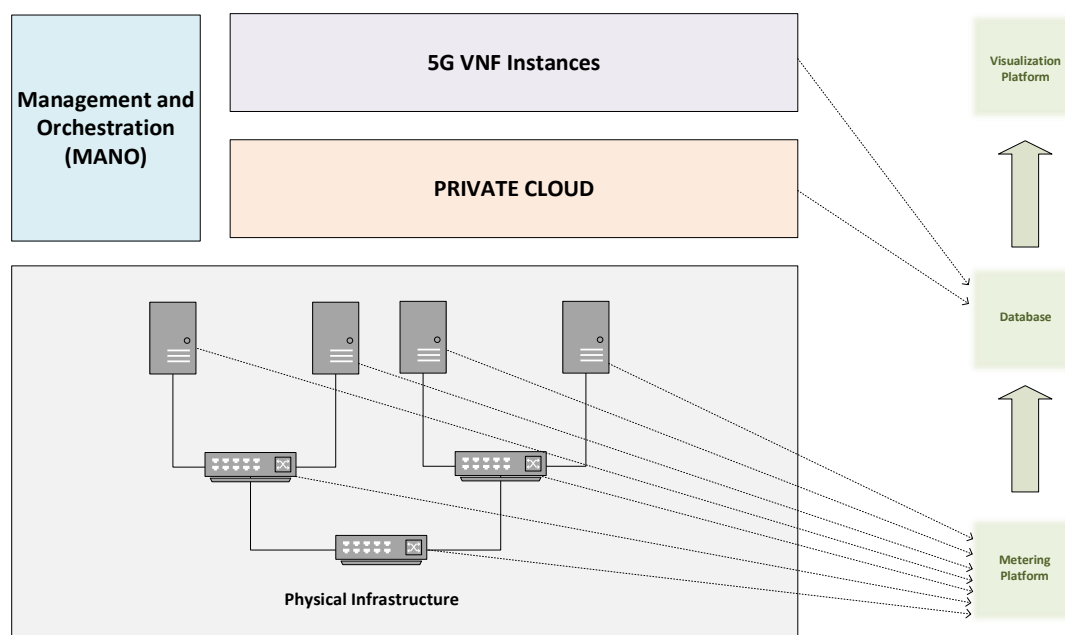


Figure 6. 3: Environment Description

In this environment our goal is to dynamically instantiate new slices belonging to different operators through a MANO platform. The slices may belong to the same or different SSTs, with different SDs. In this example, the slices have different SMF and UPF nodes for the two slices, but share all other CN functions (like AMF, UDR etc..). Furthermore, the end-to-end path of each slice is different and depends on the Service Type. For simplification purposes, we consider the creation of two distinguished slices, one that is served at a local MEC (private network slice), and one served by a remote facility (public network slice).

6.4. Experimental Setup

6.4.1. 5G Platform Description

We begin by describing the Environment of our lab testbed which was used for this scenario, as it shown in **Figure 6. 3**. Going bottom-up we have the physical infrastructure which comprises a router, servers, switches and the links connecting them. Energy metering devices are attached to each piece of our equipment. The physical resources are clustered into an (openstack) cloud platform, which enables the MANO framework to create 5G related VNF instances. Metrics regarding the performance of either a physical or a virtualized compute node are extracted and stored to a Monitoring database, along with the energy metering metrics. All the metrics stored in the Monitoring Platform are then plotted through the Visualization Platform, as shown in **Figure 6. 3**. The entire lab environment described above was used for activities related to Horizon 2020 projects, like 5G-COMLETE [20] and 5G-VICTORI [21].

6.4.2. Orchestration Platform Overview

An VNF orchestrator framework, developed by a division of ETSI, is the OSM [34]. OSM, utilizing well-established open-source tools and development techniques, performs the

role of orchestrator in an E2E Network Service that is capable of modeling and automating network services. It is a community-driven effort that is fully aligned with ETSI-NFV layer/module-based architecture. Each layer has its own degrees of abstraction, in order to provide the necessary independence that will allow networks to scale as required. Furthermore, the modularized nature of OSM layers, enables effortless module replacements when necessary. Its overall architecture significantly advances NFV technologies and standards and gives VNF vendors a way to test and validate how their services operate and interact at the production level and at the commercial level.

OSM hides the technical complexities and offers simplicity in the way services are provided. This is achieved through an ETSI NFV-aligned Information Model (IM) that manages, automates, and monitors the entire lifecycle of network functions, services, and slices. The IM is independent of the underlying infrastructure and can be used across various VIM types and transport technologies. Additionally, OSM provides a unified northbound interface (NBI), based on NFV SOL005 [35], for complete control of system and service operations. OSM extends the concept of a "Network Service" to include network functions across different domains, treating all components in an undistinguishable manner. An NS in OSM may consist of a combination of virtual, physical or hybrid NFs (VNFs-PNFs-HNFs) with on demand transport links amongst different sites.

OSM covers the tasks of NFV Manager and VNF Orchestrator in the classic NFV-MANO architecture. It aims at highlighting the capabilities of NFV, delivering Networks as a Service-NaaS. There are two types of NaaS service objects that OSM is able to provide on demand: the NS that is a composition of VNFs at specific arrangement, and the NSI that is a composition of several NSs that can be treated as a single entity. As an orchestrator, it is responsible for the whole lifecycle of the VNFs. The lifecycle management is performed through the LifeCycle Management (LCM) module and the VNF Configuration and Abstraction (VCA) layer. LCM module in OSM is responsible of managing the workflows associated to life cycle events of VNF and NS such as instantiation, termination, scaling, healing and upgrading. The VCA layer provides a uniform interface to the compute resources available to a VNF, regardless of the underlying virtualization technology. The LCM module and VCA layer work together to ensure that VNFs are managed efficiently and effectively throughout their lifecycle. When a VNF is instantiated, the LCM module will communicate with the VCA layer to ensure that the required compute resources are available, and will allocate those resources as needed. If the VNF needs to be scaled up or down, the LCM module will communicate with the VCA layer to adjust the compute resources accordingly.

The VCA layer is handling VNF modeling and configuration through charms, with primitives and attributes [36]. Charms are a generic set of scripts and metadata for deploying and operating software which can be adapted to any use case. Essentially, a charm captures the DevOps knowledge of experts in a particular product. With the help of charms, applications can be quickly, easily and repeatedly deployed and scaled when needed. They can be written in any programming language that can be executed from the command line and are made up of several YAML configuration files and a number of "hooks." A "hook" refers to a code that installs software, starts and stops services, configures and updates charms, and manages connections with other charms. Charms can be deployed in two ways: as native charms, that are allocated inside the VMs that are part of the VNFs, or as proxy charms that are not allocated in the workload, but in LXC containers inside machine that hosts OSM. Proxy Charms use ssh or other methods to get into the VNF instances and configure them. They are widely used for cases where fixed images are used for the VMs, which cannot be modified. However, when a proxy charm is

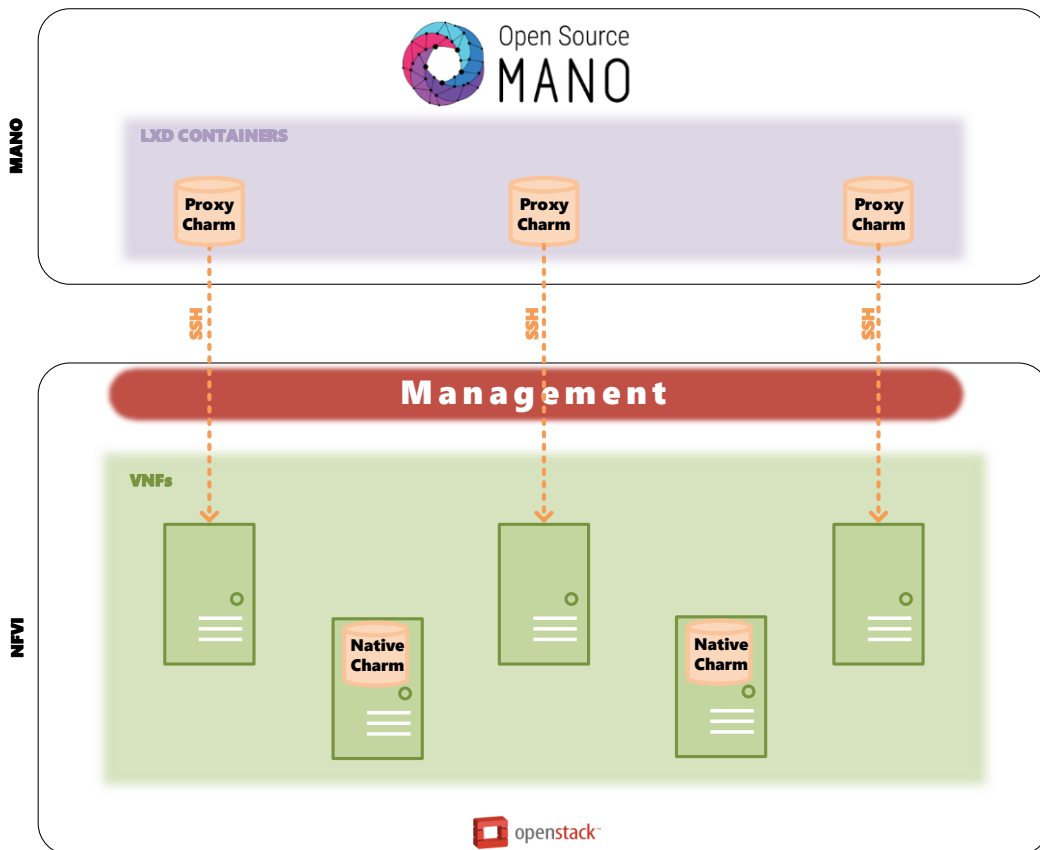


Figure 6. 4: Network service instantiation and configuration with proxy charms [34]

deployed, there are several time-consuming steps that are executed by default. The LXD container that hosts the charm must be launched and configured with juju, and the charm needs to be installed. The deployment can be accelerated by using a distributed VCA located in the same datacenter as the VIM. **Figure 6. 4** shows a graphical illustration of a network service instantiation and configuration with proxy charms.

The VNFs that comprise the service are provided to OSM as packages that include the descriptor of the VNF (VNFD), that describes the specifications of the VMs that comprise the VNF along with their connections, and the charms that need to be deployed for the NS initialization and runtime operations. The arrangement of the VNFs, and the connections between them is given to OSM through a script called network service descriptor (NSD). In general, descriptors, or configuration templates, are what OSM's IM uses to describe the key characteristics of managed objects (e.g VNFs or NSs) in a network. For each component, descriptors specify how it will be deployed and used, as well as how it will interact with other components. Descriptors are written in YAML, a markup language designed for data that is easy to read and understand.

Figure 6. 5 shows the steps that are performed when a network service is instantiating from OSM. First, OSM instructs VIM to create the necessary VMs that are needed for the VNFs, and the network connections between them. Then the life-cycle management of the VNFs begins. That includes three stages: basic instantiation of the VNF (“Day-0 Configuration”), service initialization (“Day-1 Configuration”) and runtime operations

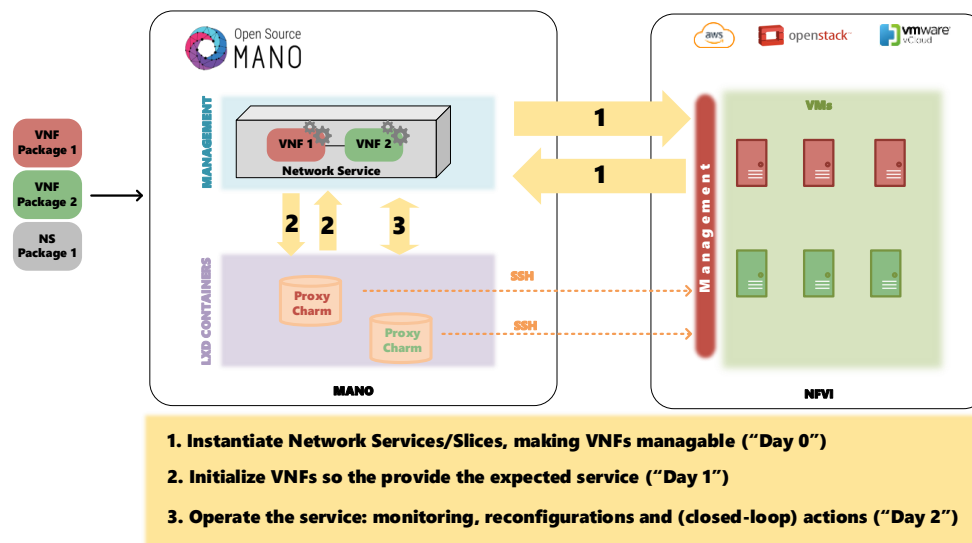


Figure 6. 5: Network Service Instantiation through OSM [34]

("Day-2 Configuration"). Day 0 is configured through cloud-init, a service used for modifying the generic OS configuration of virtual machines on boot. It mainly consists of basic VM configuration like import ssh-keys, create users/pass, network configuration etc. For Day 1 and Day 2 configurations, OSM uses the aforementioned charms to configure and manage the VNFs for the specific service.

6.5. Implementation

6.5.1. Creation of Network Descriptors

The scenario described above assumes a sophisticated 5GC topology with the following characteristics:

- CP and UP separation,
- multiple UPF elements with different roles. The number of UPF elements in the User Plane path may vary. The UPF element is categorized into three main roles, according to the function it performs on the User Plane path.
 - PDU Session Anchor (PSA) UPF: The UPF where the PDU Session is terminated.
 - Intermediate (I) UPF: This UPF is inside the path between the RAN and the PSA-UPF, and is responsible for forwarding data between them.
 - Branching (B) UPF: As it is implied by its name, the B-UPF redirects the uplink traffic to the appropriate UPF that ends the PDU Session, and merges the downlink traffic from different PSA-UPFs, to the UE.

It is important to highlight the fact that a UPF can perform more than one role at the same time. For instance, a UPF may act as a PSA-UPF for some traffic, and as an I-UPF for other traffic.

- configuration of different Network Slices

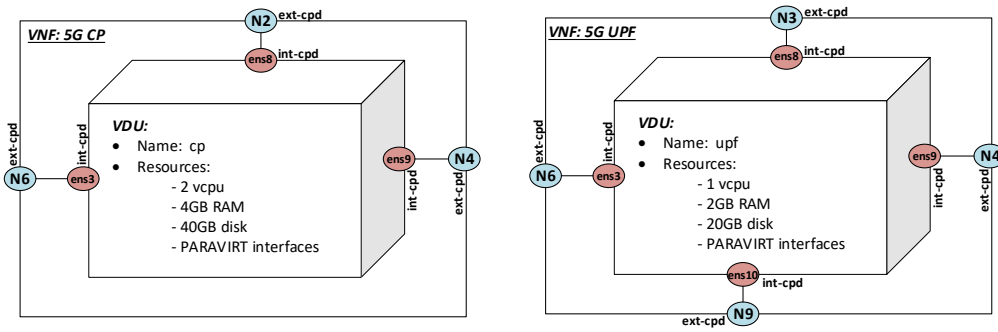


Figure 6. 6: Graphical illustration of the VNFs for 5G CP and UP respectively.

- making use of a Local Breakout so that a user can either be served from a MEC or central Cloud.

After the successful manual deployment of a 5G system dictated by the above-mentioned characteristics a mapping of all the above to appropriate VNFs is necessary, in order to be able to automate the above procedure through the OSM MANO platform. To do this, we followed these steps:

1. Firstly, we defined the required VNFs for the aforementioned 5G topology, namely a VNF for the 5G CP and a VNF for the 5G UP functionalities. We created two VNFs, one for the 5G-Control Plane and one for the UPF. The VNFs describe the specifications of the VMs that will host the CP/UPF of 5G, and the charms that dictate the actions that need to be performed for the proper functionality of those two. **Figure 6. 6** shows a graphical illustration of the CP and UP VNFs.
2. Then we continued by creating two NSDs, one for the CP and one for the UPF. Each NSD contains the corresponding VNF, and describes the network connectivity of the VNF inside the private cloud. **Figure 6. 7** shows a graphical illustration of the CP and UPF NSDs.
3. Finally, we defined the network slices. The NSSIs that comprise each slice were investigated, along with their network connections and sharing capabilities. In

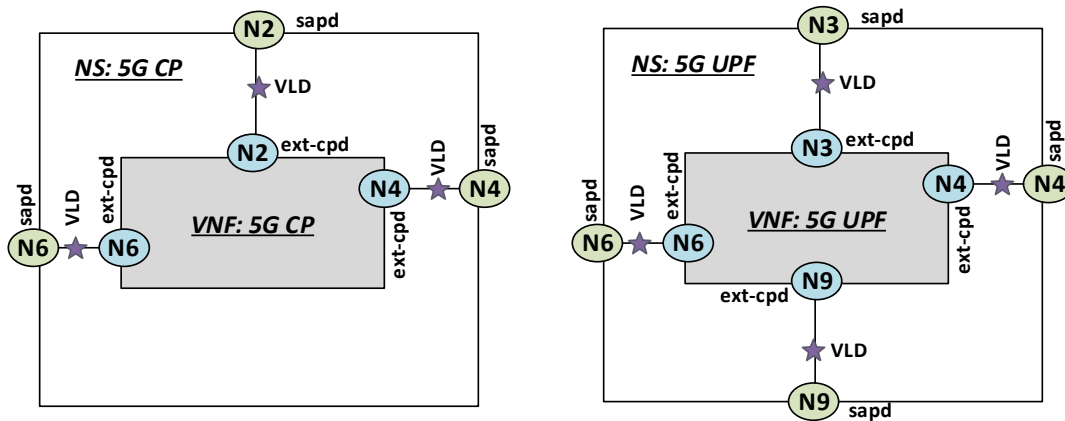


Figure 6. 7: Graphical illustration of the NSDs for 5G CP and UP respectively.

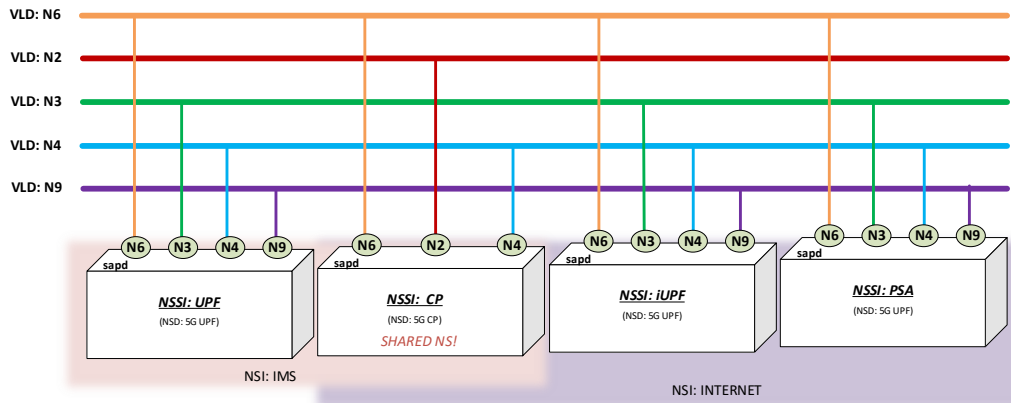


Figure 6. 8: Graphical illustration of the NSSIs of each NSI

OSM all this necessary information about the slice is stored in a descriptor, called Network Slice Templates (NST). In our scenario, we created two NSTs, by combining the CP and UPF NSs. The two slices described in the NSTs share the 5G CP (same CP NS), but have different UP paths (different UPF NS). The first slice has only one UPF, while the second has two, with the ability to serve its clients in a MEC facility. **Figure 6. 8** shows a graphical illustration of the two slices that are going to be supported in the aforementioned 5g topology.

After the creation of the NSTs, appropriate configurations need to be made for the proper deployment of the slices. These are given through juju charms to OSM. Those actions are distinguished to

- actions that are automatically performed at the instantiation of the slices (Day 1 actions or initial-config-primitive in the VNFD descriptor), and
- actions that can be dynamically performed during the deployment of the slice (Day 2 actions or config-primitive in the VNFD descriptor).

Specifically, for Day 1 the following actions are performed

1. Configure ssh access and IPs of the VMs
2. Manipulation of the configuration files for the proper operation of 5G core and UPF
3. Role of UPF (i-upf or psa with N3/N9 interface)
4. Load GTPU tunnel
5. Start UPFs

While as Day 2 the Runtime Operations configured are:

- Start/Stop 5G CORE
- Start/Stop SMFs
- Start/Stop UPFs

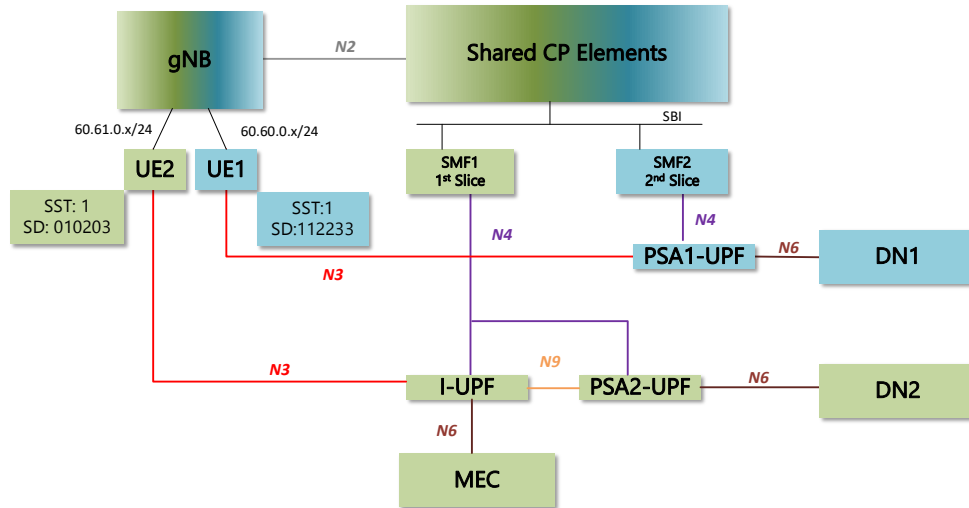


Figure 6. 9: The two slices implemented on the lab

6.5.2. Slice Deployment and Results

Figure 6. 9 shows the logical connection of the network elements of the 5GC topology. The scenario was demonstrated in the 5G-COMPLETE [20] project and involves the instantiation of two slices with the same SST and different SD. The first slice carries the traffic of File Transfer, while the second accommodates a video streaming service. The two slices share the VNF that contains the 5G core functions, but each one has its own UPF VNFs. The slice for the video streaming services can be served either at a MEC server inside the lab, or at the central cloud (local breakout).

For the correct instantiation of slices, initial parameters for the VM IP configuration and the specific UPF roles must be given. The video streaming slice is initially instantiated. A UE (UE1) connects to the slice and begins to receive a video stream generated from a local server inside the lab. With the first slice being operational, the second slice is instantiated. A second UE (UE2) connects to that slice and starts downloading a file from the web. The whole procedure was captured in wireshark. In Figure 6. 10 we can see the wireshark traces, where the differentiation of the two slices is evident.

Through the data visualization platform that is used along with the lab cloud implementation, in Figure 6. 11 we can check and confirm that the traffic of each slice is

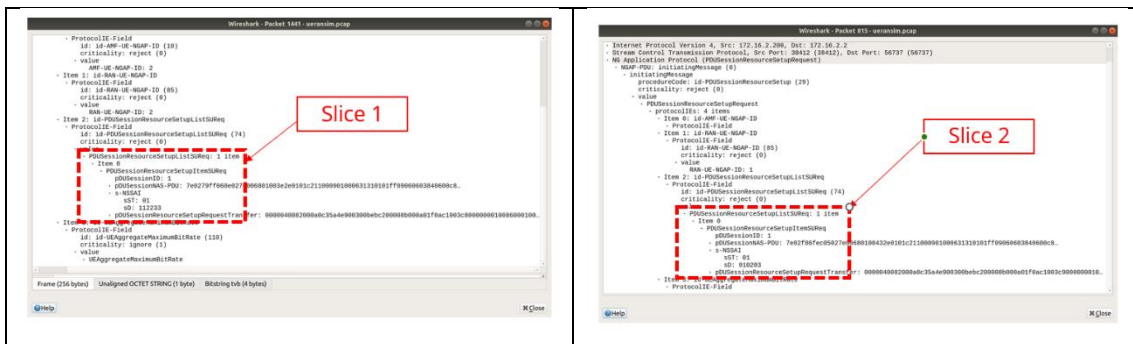


Figure 6. 10: wireshark traces after the instantiation of the two slices

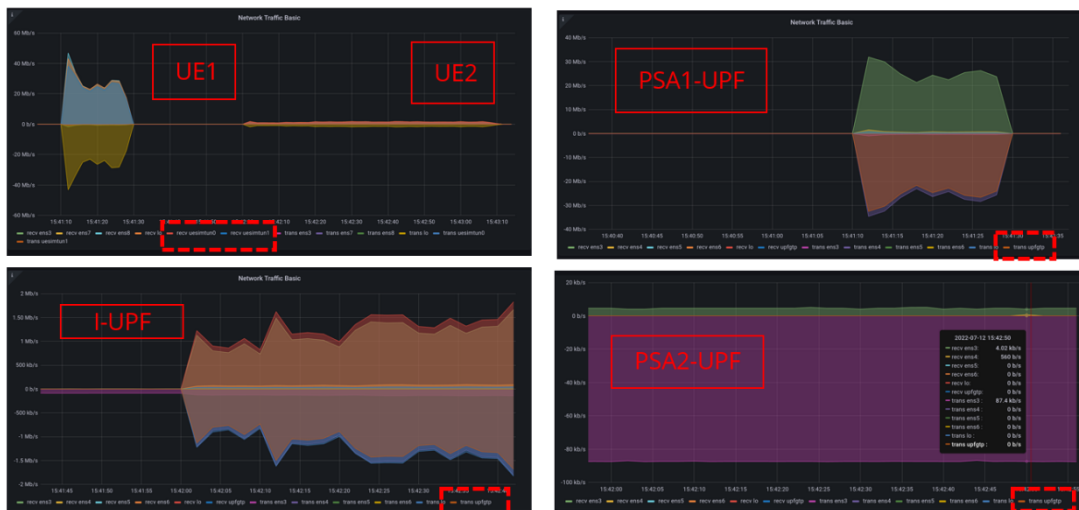


Figure 6. 11: Visualization Platform screenshots for UE traffic of each slice

transmitted through the right path that was configured for the slice. In the first screenshot (up-left) we can see the traffic that is generated by the two UEs. UE1 refers to the UE that connects to the data transfer slice, and UE2 is the UE that connects to the video streaming slice. The other images show the traffic that passes through each one of the three UPFs. Since the UE that receives the video streaming is served by a local MEC server inside the lab, traffic terminates at I-UPF and does not proceed to PSA2-UPF.

6.6. Summary

In this chapter, we presented our contributions towards automating the deployment of a sophisticated 5G multi-operator network and E2E provisioning of multiple slices. For this purpose, we used OSM, an NFV-MANO compliant orchestrator. We created network descriptors for the core and the user plane of a 5G network. Combining those descriptors, we successfully deployed dynamic 5g network slices on top of a softwarized multi-operator 5G platform hosted in a containerized data centre infrastructure.

References

- [1] Barakabitze, Alcardo & Ahmad, Arslan & Hines, Andrew & Mijumbi, Rashid, “5G Network Slicing using SDN and NFV: A Survey of Taxonomy, Architectures and Future Challenges.”, in *Computer Networks*, Vol. 167, 106984, 2020. Available: <https://doi.org/10.1016/j.comnet.2019.106984>.
- [2] A. Barakabitze, N. Barman, A. Ahmad, S. Zadtootaghaj, L. Sun, M. Martini, L. Atzori, “QoE Management of multimedia services in future networks: A Tutorial and survey.”, in *IEEE Communications Surveys & Tutorials*, vol. 22, no. 1, pp. 526-565, Firstquarter 2020. Available: <https://doi.org/10.1109/COMST.2019.2958784>.
- [3] A.A. Barakabitze, I.-H. Mkwawa, L. Sun, E. Ifeachor, “Qualitysdn: improving video quality using MPTCP and segment routing in SDN/NFV.”, *2018 4th IEEE Conference on Network Softwarization and Workshops (NetSoft)*, Montreal, QC, Canada, 2018, pp. 182-186. Available: <https://doi.org/10.1109/NETSOFT.2018.8459917>.

- [4] A. Ibrahim, T. Tarik, S. Konstantinos, A. Ksentini, H. Flinck, "Network slicing & softwarization: A Survey on principles, enabling technologies & solutions.", in *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2429-2453, thirdquarter 2018. Available: <https://doi.org/10.1109/COMST.2018.2815638>.
- [5] I. Afolabi, M. Baga, T. Taleb, H. Flinck, "End-to-End network slicing enabled through network function virtualization.", *2017 IEEE Conference on Standards for Communications and Networking (CSCN)*, Helsinki, Finland, 2017, pp. 30-35. Available: <https://doi.org/10.1109/CSCN.2017.8088594>.
- [6] Condoluci, M., Sardis, F., Mahmoodi, T., "Softwarization and Virtualization in 5G Networks for Smart Cities.", In: *et al. Internet of Things, IoT Infrastructures, IoT360 2015, Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering*, vol 169. Springer, Cham., 2016. Available: https://doi.org/10.1007/978-3-319-47063-4_16.
- [7] X. Zhou, R. Li, T. Chen and H. Zhang, "Network slicing as a service: enabling enterprises' own software-defined cellular networks," in *IEEE Communications Magazine*, vol. 54, no. 7, pp. 146-153, July 2016. Available: <https://doi.org/10.1109/MCOM.2016.7509393>.
- [8] J. Sánchez, I. G. Ben Yahia, N. Crespi, T. Rasheed and D. Siracusa, "Softwarized 5G networks resiliency with self-healing," *1st International Conference on 5G for Ubiquitous Connectivity*, Akaslompolo, Finland, 2014, pp. 229-233. Available: <https://doi.org/10.4108/icst.5gu.2014.258123>.
- [9] *View on 5G Architecture*, 5G PPP Architecture Working Group, Version 3.0, 06-2019. [Online]. Available: https://5g-ppp.eu/wp-content/uploads/2019/07/5G-PPP-5G-Architecture-White-Paper_v3.0_PublicConsultation.pdf
- [10] Y. Liu, J. E. Fieldsend and G. Min, "A Framework of Fog Computing: Architecture, Challenges, and Optimization," in *IEEE Access*, vol. 5, pp. 25445-25454, 2017. Available: <https://doi.org/10.1109/ACCESS.2017.2766923>.
- [11] S. Yi, Z. Hao, Z. Qin and Q. Li, "Fog Computing: Platform and Applications," *2015 Third IEEE Workshop on Hot Topics in Web Systems and Technologies (HotWeb)*, Washington, DC, USA, 2015, pp. 73-78. Available: <https://doi.org/10.1109/HotWeb.2015.22>.
- [12] Akyildiz, Ian & Nie, Shuai & Lin, Shih-Chun & Chandrasekaran, Manoj, "5G Roadmap: 10 Key Enabling Technologies," in *Computer Networks*, Vol 106, 2016. Available: <https://doi.org/10.1016/j.comnet.2016.06.010>.
- [13] R. Menon, R. M. Buehrer and J. H. Reed, "On the Impact of Dynamic Spectrum Sharing Techniques on Legacy Radio Systems," in *IEEE Transactions on Wireless Communications*, vol. 7, no. 11, pp. 4198-4207, November 2008. Available: <https://doi.org/10.1109/T-WC.2008.070155>.
- [14] *Network Functions Virtualisation (NFV); Terminology for Main Concepts in NFV*, ETSI GS NFV 003 V1.2.1 (2014). [Online].
- [15] *Network Functions Virtualisation (NFV); Architectural Framework*, ETSI GS NFV 002 V1.2.1 (2014a). [Online].
- [16] *Network Functions Virtualisation (NFV); Management and Orchestration*, ETSI GS NFV-MAN 001 V1.1.1 (2014b). [Online].
- [17] X. Foukas, G. Patounas, A. Elmokashfi and M. K. Marina, "Network Slicing in 5G: Survey and Challenges," in *IEEE Communications Magazine*, vol. 55, no. 5, pp. 94-100, May 2017. Available: <https://doi.org/10.1109/MCOM.2017.1600951>.

- [18] NGMN 5G Initiative Team, "A Deliverable by the NGMN Alliance: NGMN 5G White Paper". [Online]. Available: https://www.ngmn.org/wp-content/uploads/NGMN_5G_White_Paper_V1_0.pdf.
- [19] M. Jiang, M. Condoluci and T. Mahmoodi, "Network slicing management & prioritization in 5G mobile systems," *European Wireless 2016; 22th European Wireless Conference*, Oulu, Finland, 2016, pp. 1-6.
- [20] H2020 Project 5G-COMLETE. <https://5gcomplete.eu/>. [Online]
- [21] H2020 Project 5G-VICTORI. <https://www.5g-victori-project.eu/> [Online]
- [22] Afolabi, Ibrahim & Ksentini, Adlen & Bagaa, Miloud & Taleb, Tarik & Corici, Marius & NAKAO, Akihiro, "Towards 5G Network Slicing over Multiple-Domains," in *IEICE Transactions on Communications*, Vol E100.B, 2017. Available: <https://doi.org/E100.B.10.1587/transcom.2016NNI0002>.
- [23] Taleb, Tarik & Mada, Badr Eddine & Corici, Marius & Nakao, Akihiro & Flinck, Hannu, "PERMIT: Network Slicing for Personalized 5G Mobile Telecommunications", in *IEEE Communications Magazine*, Vol. 55, pp. 88-93, 2017. Available: <https://doi.org/10.1109/MCOM.2017.1600947>.
- [24] U. Habiba and E. Hossain, "Auction Mechanisms for Virtualization in 5G Cellular Networks: Basics, Trends, and Open Challenges," in *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2264-2293, thirdquarter 2018. Available: <https://doi.org/10.1109/COMST.2018.2811395>.
- [25] R. Mijumbi, J. Serrat, J. -L. Gorricho, N. Bouten, F. De Turck and R. Boutaba, "Network Function Virtualization: State-of-the-Art and Research Challenges," in *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 236-262, Firstquarter 2016. Available: <https://doi.org/10.1109/COMST.2015.2477041>.
- [26] K. Samdanis, X. Costa-Perez and V. Sciancalepore, "From network sharing to multi-tenancy: The 5G network slice broker," in *IEEE Communications Magazine*, vol. 54, no. 7, pp. 32-39, July 2016. Available: <https://doi.org/10.1109/MCOM.2016.7514161>.
- [27] *Applying SDN Architecture to 5G Slicing*, ONF TR-526 (2016).
- [28] M. Iwamura, "NGMN View on 5G Architecture," in *2015 IEEE 81st Vehicular Technology Conference (VTC Spring)* (2015): 1-5.
- [29] *Framework of Network Virtualization for Future Networks*, ITU-T Y.3011, January 2012. [Online]. Available: <https://www.itu.int/rec/T-REC-Y.3011-201201-I>.
- [30] *Network Functions Virtualisation (NFV) Release 3; Evolution and Ecosystem, Report on Network Slicing Support with ETSI NFV Architecture Framework*, ETSI GR NFV-EVE 012 v3.1.1, December 2017.
- [31] King, Daniel and Young J. Lee. "Applicability of Abstraction and Control of Traffic Engineered Networks (ACTN) to Network Slicing." (2018).
- [32] *Study on management and orchestration of network slicing for next generation network (Release 15), Technical Specification Group Services and System Aspects, Telecommunication management*, 3GPP TR 28.801 V15.1.0, January 2018.
- [33] X. de Foy, A. Rahman, "Network slicing -3GPP Use case draft-defoy-netslices-3GPP-Network-Slicing-02", in *InterDigital Communications, LLC*, 2017.
- [34] *Open Source MANO*, ETSI [OSM \(etsi.org\)](https://www.etsi.org). [Online]

- [35] *Network Functions Virtualisation (NFV) Release 3; Protocols and Data Models; RESTful protocols specification for the Os-Ma-nfvo Reference Point*, ETSI GS NFV-SOL 005 V2.7.1, January 2020.
- [36] Canonical Juju. Juju docs <https://juju.is/>. [Online]

Chapter 7

Conclusions and Future Work

In this thesis, we focused on the disaggregated architecture of 5G Networks. The concept of a new generation universal 5G platform is adopted, where a variety of networking technologies (wireless and wired) that are integrated and jointly optimized. This infrastructure interconnects "disaggregated" compute/storage and network components through the usage of programmable HW and network softwarisation. The aim of this thesis is the development of mathematical tools, algorithms and protocols for the optimization of the infrastructure and the services it provides. All of our contributions were experimentally tested, and implemented in open-source software that operates on commodity hardware, allowing for performance measurement in actual environment scenarios and direct comparison of our frameworks with existing solutions and standards.

To begin with, we studied the challenges of future telecommunication radio access networks by exploring the concept of "resource disaggregation." This concept involves separation of components and functions that can be placed at discrete geographical locations. This which enables splitting of BBU function chains between RUs, DUs and CUs, and their connection through FH and MH links. To improve cost efficiency and energy consumption, we proposed an architectural model that leverages compute resources required for BBU function processing located both at the MEC and large-scale centralized DCs. To support this novel architecture seamlessly, high bandwidth/low latency SDN controlled optical transport networks were utilized to connect MEC domains with medium to large-scale DCs that host general-purpose servers in the optical access and metro domains. In order to optimize FH flows in this proposed model, we propose a hybrid centralized/distributed 5G network management solution. According to our proposed scheme, the centralized controller makes high-level FH connectivity decisions, while RUs make local decisions associated with the optimal FH flow selection dynamically and in a non-cooperative manner based on EGT, with the objective of minimizing their total operational expenditures. We also investigated the controller placement problem since this decision has a direct impact on the system's stability. The stability of the proposed scheme is influenced by network latency; therefore, we proposed a metric for sizing the SDN FH/BH network. Finally, a relation of a basic learning algorithm with EGT was utilized to extend the analysis and approximate more complex scenarios of asymmetrical, multi-interaction systems. With the adoption of the proposed scheme, the controller placement problem is confronted through a new perspective, that guarantees a stable 5G system that supports the optimal split option with efficiency gains in terms of cost and energy consumption. Following this, our focus shifted towards the core aspect of 5G Networks, which introduces the new CUPS. CUPS separates the signaling traffic carried by the control plane between different 5GN entities, and the user traffic forwarded by the user plane. The UPF plays a crucial role in handling a significant portion of the user plane functionality in 5G

systems, as it is responsible for forwarding actual user traffic with very strict performance requirements. Depending on the type of service required and the 5G-RAN deployment option, UPF nodes can be located either closer or further away from the 5G-RAN, redirecting traffic to MEC servers to reduce latency or to central cloud facilities. Hence, the challenge of selecting the optimal UPF elements arises since selecting an unavailable UPF computational resource may result in service blockage and delays. To tackle this issue, we proposed an EGT (Evolutionary Game Theory) model specifically designed for dynamic selection of optimal UPF elements with the objective of minimizing total service delay. We developed cost functions for the EGT model using lab measurements obtained from an open-source 5G platform hosted in an optical datacenter cloud environment. With our proposed EGT model, we can dynamically choose the most appropriate UPF element to utilize computational resources and reduce service delay.

Expanding upon the concept of disaggregated network architecture, the analysis shifted focus from 5G-networks to 6G-systems, which are anticipated to support a diverse range of services through a shared infrastructure facilitated by network slicing. 6G systems are projected to operate in a decentralized manner, which allows for applications to directly intervene in the control processes necessary for ensuring QoE. This is carried out through the use of the AF, which manages the application running on UE and the AS that supports the service. The AF plays a critical role in providing high QoE services, as it receives feedback from the application and can influence traffic routing and steering decisions. However, allowing the AF to operate without regulation may result in instability within the system. To address this issue, a fully distributed decision-making framework was designed, implemented, and evaluated theoretically and experimentally to solve the flow assignment problem in 6G systems. It has been shown that under specific conditions, the framework converges to a stable point that provides the optimal balance between QoE and cost efficiency. The cost functions utilized by the EGT model incorporate both network and compute costs, which are realistically derived through a detailed profiling process conducted on an operational 5G testbed. This profiling process enables modeling of system performance and requirements under different operational scenarios, which can practically assist in the optimized lifecycle management of provisioned services.

Finally, we practically deployed and orchestrated a sophisticated 5G multi-operator network and performed automated provisioning of multiple E2E slices. To achieve this, we used an orchestrator called OSM, which is compliant with the NFV standard. A network slice has its own dedicated resources from the access, transport, core and edge network as well as cross-domain component from various domains under the same or distinct administrations. Combining the control and user plane descriptors, we successfully deployed dynamic 5G network slices on top of a softwarized multi-operator 5G platform hosted in a containerized data centre infrastructure. Using the created VNFs we were able to perform single click deployment of the core network, and provision multiple slices with different characteristics.

In the future, we foresee extending our schemes towards integrating ML and AI approaches to the optimization of the resource allocation problems that arise in the 5G infrastructure. We plan to investigate possible dynamic provisioning and placement of compute resources and UPF elements, based on the historical traffic and power consumption data that are fed to AI algorithms in order to predict the forthcoming network demand. Through the employment of 5G MANO tools, we will concentrate on developing a 5G platform where the requested resources for the slices will be

dynamically optimized through ML techniques, and deployment of the slices will be fully automated and continuously monitored.