**INFSO-ICT-248523 BeFEMTO**

**D5.2**

Version 1.0

*Femtocells access control, networking, mobility and management mechanisms (final)*

| | |
|---|---|
| **Contractual Date of Delivery to the CEC:** | M24 |
| **Actual Date of Delivery to the CEC:** | M24 |
| **Author(s):** | Andreas Maeder;Andrey Krendzel;Cristina Peña;Emilio Mino;Frank Zdarsky;Jaime Ferragut;John Fitzpatrick;José Núñez;Josep Mangues;Luis Cucala;Miquel Soriano;Pablo Arozarena;Tao Guo |
| **Participant(s):** | NEC, TID, CTTC, UniS, PTC |
| **Workpackage:** | WP5: Femtocells Access Control, Networking, Mobility and Management |
| **Estimated person months:** | 50 |
| **Security:** | PU |
| **Nature:** | R |
| **Version:** | Version 1.0 |
| **Total number of pages:** | 183 |

**Abstract:**

This deliverable D5.2 describes the final networking schemes developed in WP5 that were chosen by the BeFEMTO project due to their suitability to the concepts and system architecture presented in WP2. In particular, it presents concrete solutions in the areas of traffic management, mobility management, network management, and security as well as some initial evaluation of some of these solutions.

# Executive Summary

This document describes the concrete solutions selected by BeFEMTO for the traffic management, mobility management, network management, and security concepts studied within Work Package 5 (WP5). The purpose of this deliverable is to describe such solutions, as well as a summary of those already presented in D5.1. In this sense, this document serves as a compilation of all the networking solutions developed towards the realization of the BeFEMTO system concept, as presented in WP2.

Following an introduction of this deliverable (Chapter 1), that also provides a general overview of the technical activities tackled within WP5 and the relationship among them, the following chapters provide a detailed description of the work activities that have been addressed:

Chapter 2 describes activities related to traffic forwarding and resource sharing.

Chapter 3 is focused on the mobility management issues.

Chapter 4 deals with the management of evolved femtocell networks.

Chapter 5 is devoted to security issues.

Chapter 6 includes technical work on revenue sharing in multi-stakeholder scenarios.

## List of Acronyms and Abbreviations

| μs | Micro-second |
|---|---|
| 3GPP | 3rd Generation Partnership Project |
| 4G | 4th Generation |
| A/D | Analog to Digital |
| AAA | Authentication, Authorization, Accounting |
| AC | Alternate Current |
| ACK | Acknowledgement |
| ACLR | Adjacent Channel Leakage Ratio |
| ACPL | Adjacent Channel Power Leakage |
| AKA | Authentication and Key Agreement |
| AP | Access Point |
| API | Application Programming Interface |
| AS | Access Stratum |
| AWGN | Additive White Gaussian Noise |
| b/s/Hz/cell | bit per second per Hertz per cell |
| BCH | Broadcast Channel |
| BeFEMTO | Broadband evolved FEMTO networks |
| BLER | Block Error Rate |
| BS | Base Station |
| BW | Bandwidth |
| CA | Carrier Aggregation |
| CC | Component Carrier |
| CDF | Cumulative Distribution Function |
| CF | Carrier Frequency |
| C-Plane | Control-Plane |
| CQI | Channel Quality Indication |
| CRC | Cyclic Redundancy Check |
| CRS | Cell Specific Reference Symbols |
| CS | Customer Service |
| CSG | Closed Subscriber Group |
| dB | Decibel |
| dBm | decibel (referenced to one milliwatt) |
| DC-HSUPA | Dual Carrier-High Speed Uplink Speed Packet Access |
| DeNB | Donor Evolved Node-B |
| DFT-OFDM | Discrete Fourier Transform- Orthogonal Frequency Division Multiplexing |
| DHCP | Dynamic Host Configuration Protocol |
| DL | Downlink |
| DSCP | Differentiated Services Code Point |
| EAP | Extensible Authentication Protocol |
| ECN | Explicit Congestion Notification |
| eNB | Evolved Node-B (LTE macro base station) |
| EPA | Extended Pedestrian A-model |
| EPC | Evolved Packet Core |
| ETSI | European Telecommunication Standards Institute |
| ETU | Extended Typical Urban model |
| EUTRA | Evolved Universal Terrestrial Radio Access |

| EUTRAN | Evolved Universal Terrestrial Radio Access Network |
|---|---|
| EVA | Extended Vehicular A-model |
| EVM | Error Vector Module |
| FAP | Femto Access Point |
| FDD | Frequency Division Duplex |
| FER | Frame Error Rate |
| FFR | Fractional Frequency Reuse |
| FI | Fairness Index |
| FRC | Fixed Reference Channel |
| FTTH | Fibre To The Home |
| FUE | Femtocell UE |
| Gb/s | Giga bits per second |
| GHz | Gigahertz |
| GPRS | General Packet Radio Service |
| GTP | GPRS Tunneling Protocol |
| HARQ | Hybrid Automatic Repeat Request |
| HeNB | Home evolved Node-B |
| HetNet | Heterogeneous Networks |
| HNB | Home Node-B |
| HO | HandOver |
| HSS | Home Subscriber Server |
| HSUPA | High Speed Uplink Speed Packet Access |
| ICIC | Inter Cell Interference Coordination |
| IMS | IP Multimedia Subsystem |
| IMT-A | International Mobile Telephony – Advanced |
| InH | Indoor Hotspot |
| IP | Internet Protocol |
| ISD | Inter Site Distance |
| ITU-R | International Telecommunication Union-Radiocommunication Sector |
| km/h | kilometre per hour |
| KPI | Key Performance Indicator |
| LFGW | Local Femtocell GateWay |
| LGW | Local P-GW |
| LIPA | Local IP access |
| LLM | Local Location Management |
| LNG | Local Network Gateway |
| LNM | Local Network Manager |
| LOS | Line Of Sight |
| LTE | 3GPP Long Term Evolution |
| LTE-A | Long Term Evolution – Advanced |
| MAC | Media Access Control |
| Mb/s | Megabits per second |
| MBMS | Multimedia Broadcast Multicast System |
| MBS | Macro Base Station |
| MBSFN | MBMS Single Frequency Network |
| MCS | Modulation and Coding Set |
| MHz | Mega Hertz |
| MIMO | Multi Input Multi Output |

| MME | Mobility Management Entity |
|---|---|
| MNO | Mobile Network Operator |
| MRN | Mobile Relay Node |
| MUE | Macro UE |
| NACK | Negative Acknowledgement |
| NAS | Non-Access Stratum |
| NASS | Network Access Support Subsystem |
| NGN | Next Generation Network |
| NLOS | Non-Line of Sight |
| NRB | Number of Resource Blocks |
| ns | Nanosecond |
| OA&M | Operation, Administration and Maintenance |
| OFDM | Orthogonal Frequency Division Multiplexing |
| OFDMA | Orthogonal Frequency Division Multiple Access |
| OLT | Optical Line Termination |
| ONT | Optical Network Termination |
| OTDOA | Observed Time Difference of Arrival |
| PCI | Physical Cell Identity |
| PDCCH | Physical Downlink Control Channel |
| PDCP | Packet Data Control Protocol |
| PDF | Probability Distribution Function |
| PDSCH | Physical Downlink Shared Channel |
| PDU | Protocol Data Unit |
| PHY | Physical (Layer) |
| $P_{loss}$ | Penetration Loss |
| PMME | Proxy MME |
| ppm | parts per million |
| PPP | Point to Point Protocol |
| PRACH | Physical Random Access Channel |
| PRB | Physical Resource Block |
| PS-GW | Proxy Serving GateWay |
| PUCCH | Physical Uplink Control Channel |
| PUSCH | Physical Uplink shared Channel |
| QAM | Quadrature Amplitude Modulation |
| QoS | Quality of Service |
| QPSK | Quadrature Phase Shift Keying |
| RACH | Random Access Channel |
| RACS | Remote Access Control Subsystem |
| RADIUS | Remote Authentication Dial in User Server |
| RAN | Radio Access Network |
| RAN4 | Radio Access Network (Working Group 4) |
| RB | Resource Block |
| RF | Radio Frequency |
| RLC | Radio Link Control |
| RLF | Radio Link Failure |
| RMa | Rural Macro |
| RMS | Root Mean Square |
| RN | Relay Node |

| R-PDCCH | Relay – Physical Downlink Control Channel |
|---------|------------------------------------------------------|
| R-PDSCH | Relay – Physical Downlink Shared Channel |
| R-PUSCH | Relay – Physical Uplink Shared Channel |
| RRC | Root Raised Cosine |
| RRH | Remote Radio Head |
| RRM | Radio Resource Management |
| RS | Relay Station |
| RSRP | Reference Signal Received Power |
| RTDOA | Relative Time Difference Of Arrival |
| Rx | Receiver |
| SCH | Synchronization Channel |
| SCTP | Stream Control Transmission Protocol |
| S-DeNB | Source DeNB |
| SDSL | Symmetric Digital Subscriber Line |
| SDU | Service Data Unit |
| SFR | Soft Frequency Reuse |
| S-GW | Serving GateWay |
| SIMO | Single Input Multi Output |
| SINR | Signal to Interference plus Noise Ratio |
| SIPTO | Selected IP Traffic Offload |
| SISO | Single Input Single Output |
| SMa | Suburban Macro |
| SNR | Signal to Noise Ratio |
| SO-CRRIM | Self-Optimizing Centralized RRIM |
| SON | Self-Organising Networks |
| SPS | Semi-Persistent Scheduling |
| SR | Scheduling Request |
| TDD | Time Division Duplex |
| T-DeNB | Target DeNB |
| TISPAN | Telecommunications and Internet converged Services and Protocols for Advanced Networking |
| TR | Time Ratio |
| TTI | Transmission Time Interval |
| Tx | Transmitter |
| UE | User Equipment |
| UICC | Universal Integrated Circuit Card |
| UL | Uplink |
| UMa | Urban Macro |
| UMTS | Universal Mobile Telecommunication System |
| U-Plane | User-Plane |
| UTC | Coordinated Universal Time |
| UTRA | Universal Terrestrial Radio Access |
| UTRAN | Universal Terrestrial Radio Access Network |
| VoIP | Voice over Internet Protocol |
| WAN | Wide Area Network |
| WCDMA | Wideband Code Division Multiple Access |
| WID | Work Item Description |
| WiMAX | Worldwide interoperability for Microwave Access |

| WLAN | Wireless Local Area Network |
|------|------------------------------|
| WP | Work Package |

## Authors

| Partner | Name | Phone / Fax / e-mail |
|---|---|---|
| **NEC** | | |
| | Andreas Maeder | Phone: +49 6221 4342-204<br>e-mail: andreas.maeder@neclab.eu |
| | Frank Zdarsky | Phone: +49 6221 4342-142<br>e-mail: frank.zdarsky@neclab.eu |
| | John Fitzpatrick | Phone: +49 6221 4342-<br>e-mail: johnfitzpat@ieee.org |
| | Marcus Schöller | Phone: +49 6221 4342-217<br>e-mail: marcus.schoeller@neclab.eu |
| **Telefonica I+D** | | |
| | Cristina Peña | Phone: +34 974 248917<br>e-mail: alcega@tid.es |
| | Pablo Arozarena | Phone: +34 91 483 28 66<br>e-mail: pabloa@tid.es |
| | Emilio Mino | Phone: +34 91 3128799<br>e-mail: emino@tid.es |
| | Luis Cucala | Phone: +34 91 3128799<br>e-mail: lcucala@tid.es |
| **PTC** | | |
| | Mirosław Brzozowy | Phone: +48 224135881<br>e-mail:Miroslaw.Brzozowy@t-mobile.pl |
| | Zbigniew Kowalczyk | Phone: +48 224136741<br>e-mail:Zbigniew.Kowalczyk@t-mobile.pl |
| **CTTC** | | |
| | José Núñez | Phone: +34 93 645 29 00<br>e-mail: jose.nunez@cttc.cat |
| | Jaime Ferragut | Phone: +34 93 645 29 00, ext. 2113<br>e-mail: jaime.ferragut@cttc.cat |
| | Josep Mangues | Phone: +34 93 645 29 00<br>e-mail: josep.mangues@cttc.cat |
| | Andrey Krendzel | Phone: +34 93 645 29 16<br>e-mail: andrey.krendzel@cttc.cat |
| | Miquel Soriano | Phone: +34 93 645 29 00<br>e-mail: miquel.soriano@cttc.cat |
| **University of Surrey** | | |
| | Tao Guo | Phone: +44 1483 689330<br>e-mail: t.guo@surrey.ac.uk |

# Table of Contents

## Table of Figures

# 1. Introduction

## 1.1 Scope

During year 1 of the BeFEMTO project, Work Package 5 (WP5) developed a number of innovative concepts in the area of traffic management, mobility management, security and network management of standalone and cooperative networked femtocells. These concepts are described in the WP5 deliverable D5.1.

The present deliverable D5.2 describes WP5 concrete solutions for those concepts. In this sense, it concludes the introduction of new concepts and solutions. For several of these solutions, D5.2 also already presents results from extensive analytical, simulation-based and/or experimental evaluation studies, a work that will be finished during project year 3, which is mainly devoted to this purpose.

While the first project year focused on BeFEMTO's standalone and networked femtocell themes, the second project year additionally addresses the mobile femtocell theme. Work in this direction is also included in this deliverable.

Finally, it is worth noting that WP5 scope is the research, development, and experimentation of novel femtocell technologies, but that detailed descriptions of the implementation of these technologies for the BeFEMTO testbeds is outside the scope of WP5, but is instead addressed by WP6 and documented in the respective deliverables.

## 1.2 Organisation and Overview

The present deliverable is organised as follows:

Chapter 2 groups the work items related to the traffic forwarding of user and control plane traffic within a network of femtocells and to the sharing of that network's forwarding resources.

Chapter 3 presents the work items related to mobility management within the four BeFEMTO themes (standalone, networked, outdoor relay and mobile femtocell).

Chapter 4 contains the work items related to the management of BeFEMTO femtocells in general and the networked femtocells in particular.

Chapter 5 is concerned with security-related work items.

Chapter 6 discusses revenue sharing between multiple stakeholders in "femtocell infrastructure as a service" scenarios both from the point of view of potential business models as well as enabling technologies.

## 1.3 Contributions

During the second project year, WP5 has made the following achievements and contributions:

Validation of the distributed routing concept for all traffic patterns (i.e., from LFGW to UEs, from UEs to LFGW, and local traffic). The evaluation carried out during the second year confirms the initial design ideas and the advantages of combining backpressure (for load balancing) and geographic routing (for steering packets to the destination), (i.e., low control overhead, only neighbour state kept at the femtocells, and distributed nature of the routing protocol). Specifically, the proposed strategy seems to cover the potential requirements of load-balancing and directionality an all wireless NoF dealing with all traffic patterns needs. The main parameters to tune the protocol have also been identified and preliminary evaluations have been carried out so as to eventually being able to generate deployment recommendations. This permits to continue doing refinements, and consequently, more evaluations of the proposed algorithm in order to improve its efficiency. For instance, extending the distributed routing protocol to manage multiple LFGWs, or improving the configuration of the degree of traffic distribution in every node.

> *A paper describing a distributed, stateless, backpressure-based routing protocol for large-scale, all-wireless networks of femtocells and its integration in a 3GPP context was presented in IEEE 73rd VTC Spring 2011 conference in Budapest.*

> *A paper presenting an initial evaluation (uplink routing) of the distributed routing protocol has been published in IEEE PIMRC 2011.*

*A paper presenting the implementation of the distributed routing protocol has been published in ICST WNS3 2011.*

*A paper presenting an extensive evaluation of the distributed routing protocol for any-to any- traffic patterns (i.e., uplink, downlink, and local routing) has been published in IEEE MESHTECH WORKSHOP 2011 (co-located with IEEE MASS).*

Analysis of the voice call capacity of long range WiFi as a femtocell backhaul solution.

*A paper presenting this analysis has been submitted to the IEEE Journal on Selected Areas in Communications.*

Numerical evaluation of the local mobility management and traffic offload solutions described in D5.1, in particular their benefits in terms of reducing signalling and data traffic load on the backhaul links and mobile core networks.

*A paper presenting this analysis has been submitted to the IEEE Journal on Selected Areas in Communications.*

Design of a QoS based call admission control and resource allocation mechanism for LTE femtocell deployments and the evaluation of this mechanism.

*A paper presenting this mechanism and its evaluation has been submitted to the IEEE Consumer Communications & Networking Conference (CCNC) 2012.*

Design of a self-organized tracking area list (TAL) mechanism for large-scale networks of femtocells. The aim of this proposal is to improve the accuracy of standard 3GPP location management schemes, while reducing the signalling traffic over the network of femtocells. This mechanism allows MMEs/P-MMEs to provide UEs with adaptive TALs depending on their mobility state and, eventually, current network conditions.

*A paper presenting this mechanism and its evaluation has been submitted to the IEEE International Conference on Communications (ICC) 2012.*

*A paper on Mobility management for large-scale all-wireless networks of femtocells in the Evolved Packet System has been submitted to a special issue on Femtocells in 4G Systems of the EURASIP Journal on Wireless Communications and Networking.*

*A book chapter on Mobility management for large-scale all-wireless networks of femtocells in the context of heterogeneous networks has been submitted for publication.*

Design of novel local mobility management schemes for networked femtocells based on X2 traffic forwarding to reduce the signalling cost to the core network. The target femtocell can use the local path for ongoing sessions without requiring switching the data path from core network for each handover. Both analytical models and simulation experiments show that the proposed schemes can significantly reduce the expensive signalling cost incurred to core network with reasonable local data delivery cost for traffic forwarding.

*A paper presenting the proposed local mobility management schemes based on X2 traffic forwarding and the performance evaluation has been submitted to IEEE Transactions on Vehicular Technology*

Development of a Forward Handover with Predictive Context Transfer (FH_PCT) scheme for fast handover failure recovery during inbound/outbound mobility. Analysis shows that the proposed scheme can significantly reduce the service interruption latency in case of a handover failure, especially when the backhaul latency is long.

Study on the implementation of networked femtonodes in an enterprise LAN, including: 1) The LAN configuration changes that are needed to support the networked femtocell group, 2) The logical connectivity of the networked femtonodes group to their corresponding femtonode subsystem, 3) The initial networked femtonode radio planning, and the effective radio coverage, 4) The networked femtonodes group radio self-configuration and mobility, synchronization, and performance.

Design of a distributed framework for Fault Diagnosis that enables local management capabilities in the networked femtocell network based on a distributed architecture, based on multi-agent approach. Fault Diagnosis is cooperatively conducted by a set of cooperation agents distributed in the different nodes. This fault Diagnosis network is tailored for the use of video services served by an enterprise femtocell network.

Study of various aspects of revenue sharing scenarios for femtocells related services, including legal, privacy and national roaming implications. Analysis of advertisement potential for Femtocells services, and interaction between revenue sharing stakeholders.

# 2. Traffic Forwarding and Resource Sharing

## 2.1 Centralized Traffic Management for Cooperative Femtocell Networks

In co-operative femto networks, management of femtocells, including cell provisioning and traffic prioritization, must be handled carefully. In addition to managing the femto traffic, it is also necessary to guarantee that femto traffic is not affected by the presence of non-femto traffic and vice versa. To provision this, either a separate IP network can be used for femtocells or the femto traffic can be overlaid onto the existing internal network and provide strategies to manage the traffic.

In this context, a potential solution in co-operative femto networks is to provision packets on centralized flow based mechanisms. In flow-based strategies, packets are forwarded based on explicit forwarding state installed in the forwarding elements, allowing the network to be "traffic engineered" for higher resource utilizations. They also allow for a finer control on how network resources are shared between flows. Depending on the classifier used for forwarding, flows-based mechanisms can handle anything from micro flows to aggregate flows, even concurrently.

The current work focuses on networked femtocells in general, and on enterprise networks in particular. The target is to design mechanisms necessary to allow resource-efficient traffic forwarding within the co-operative femto networks.



**Figure 2.1: Enterprise femto network**

The enterprise network under study is depicted in Figure 2.1. It consists of enterprise switches meshed into 3D cubical manner. Further each switch is connected to femtocells and enterprise servers which generate real-time and non-real-time traffic. One of the switches is connected to a femto gateway where all the femto traffic within the enterprise network terminates. The objective of this work is to study the various parameters that affect the real-time traffic and to provide basis for resource sharing and QoS mechanisms for multi-class femto and non-femto traffic within an enterprise network.

This section is organized as follows. The evaluations scenarios for centralized routing solution are explained in section 2.1.1. In subsection 2.1.2, we summarize the simulation setup and the modules that are going to be implemented for centralised traffic management. Finally, subsection 2.1.3 concludes this work.

### 2.1.1 Evaluation scenarios

The initial step is to understand the effects of co-existing traffic in cooperative femtocell networks. To accomplish this, a scenario for an enterprise network (Figure 2.1) is created where real-time and best-effort traffic, regardless of whether they originate from femto or non-femto, is mixed. The flow tables installed on the switches are based on the MAC addresses of source and destination and forwarding is based on a spanning tree protocol. In other words, the forwarding decisions entries are based on the shortest path between the source and the destination. This scenario will help in understanding the effects of co-existing traffic and will act as a baseline to further design the traffic management entity (TME). The results of the baseline evaluation should lead to answering following questions:

- The issues which arise due to sharing of network between the two traffic stakeholders and how it can be resolved
- How can resources be guaranteed for the real-time femto and non-femto traffic without starving the best effort traffic
- How can the operator validate that such resources or SLA are being met.

#### 2.1.1.1 Scenario1

Based on the baseline analysis, the next logical direction is to segregate the femto and non-femto traffic through VLANs and within each VLANs either provide strict priority or weighted fair queuing (WFQ) mechanism. The forwarding paths are still based on shortest path algorithms. This can be achieved in two ways which are described below.

<u>Case1:</u> **Strict Priority between real-time / non-real-time traffic and WFQ between femto and Non-femto traffic.**

This case is shown in Figure 2.1, where a strict priority mechanism is implemented within each VLAN. This approach will ensure that the real time traffic within femto and non-femto network is treated with high priority so as to minimize the latency.



**Figure 2.2: SP within VLAN and WFQ between VLANs**

This implementation can help to treat the real-time traffic within femto and non-femto domain efficiently. However, WFQ between femto and non-femto traffic might lead to an overall degradation in the performance. Such degradation can be frequent if most of the non-femto traffic is best effort. Under this context, the real-time femto traffic might have to wait for a longer duration in the queue which results into an increased latency, even more, excessive delay in the critical femtocell network synchronization traffic can suppose the loss of the connection of the local femtocell network with the central femtocell subsystem.

<u>Case2:</u> **WFQ between real-time / non-real-time traffic and Strict Priority between femto and Non-femto traffic.**

To overcome the drawback in previous method, as shown in Figure 2.3, we reverse the queuing principle and apply WFQ within the VLAN and then adopt a strict priority

**Figure 2.3: WFQ within VLAN and SP between VLANs**

This method will work well if the problem mentioned in case 1 is persistent. However, there might be an issue which arises when real-time traffic dominates both femto and non-femto traffic. This issue can be addressed if we adopt an arbitrary routing mechanism within the enterprise network which will find suitable paths for real-time femto and non-femto traffic.

### 2.1.1.2  Scenario 2

In this scenario the traffic management will be based on fixed resource allocation method. This allocation can be based on flow tuples which can constitute:  a) (src_ip, dst_ip) tuple or b) (src_ip, dst_ip, dscp). The queuing mechanism utilised will be the same as in scenario 1. However, fixed resource allocation may lead to following disadvantages:

- chronic underutilization of resources
- highly restricted bandwidth for both enterprise and femto traffic
- insufficient flexibility under conditions of increasing load.

The magnitude of these will be analysed and based on this analysis we will propose a load balancing mechanism.

### 2.1.1.3  Scenario 3

The function or entity performing the routing for flows needs to be aware of: a) the capacity on each link of the topology and b) the traffic within the network which includes the femto and non-femto traffic. Under these requirements a distributed approach for traffic engineering and routing would require a large signalling overhead to disseminate this information to all routing functions and would be more complex and thus more error-prone to implement. It therefore seems logical to take a centralized approach, in which traffic engineering and routing is performed within a single "routing controller" function that then installs paths with the forwarding entities in the network.

Based on the analysis of scenario 2, the logical direction for centralized routing in cooperative femto networks would be to introduce dynamic resource allocation mechanisms. Under this scope, in this work, we propose to implement and evaluate a centralized routing based on load-balancing architecture using Openflow switches connected to a common controller.

Openflow was created in 2008 by a team from Stanford University.  Openflow switches are like a standard hardware switch with a flow table performing packet lookup and forwarding. However, the difference lies in how flow rules are inserted and updated inside the switch's flow table. A standard switch can have static rules inserted into the switch, or can be a learning switch where the switch inserts rules into its flow table as it learns on which interface (switch port) a machine is attached. The Openflow switch on the other hand uses an external controller to add rules into its flow table. These rules can be based on more fine-granular identifier like QoS requirements, which will help in selecting the next

forwarding hop and to route the femto and non-femto traffic efficiently within cooperative femto networks (Figure 2.4).



**Figure 2.4: Centralized Routing Based on Openflow**

As evident, the load balancing strategies in an Openflow based environment has to be designed at the controller to which the switches are connected. This design is governed by the following criterions:

1) What kind of information does the controller require and where does it get it from?

   In an Openflow based solution the controller is responsible to install flow tables in the switches. To make an efficient decision, the controller requires a constant influx of flow requirements. This information can be related to topology, capacity, utilization/load, etc. In such a scenario it is necessary to decide if an external Traffic Management Entity (TME) is required in the network or it can be a part of the controller itself. Moreover, in a cooperative femto network, femto and non femto traffic may have different QoS requirements. Hence, the flow tables installed on the switches should be able to address these individual requirements.

2) How long are flows valid?

   In cooperative femto networks the traffic flow can be very dynamic and hence the flow table should be update frequently. This can be done either pro-actively or reactively, whichever is suitable to optimize the traffic flow. However, frequent changes in the flow tables will result in a lot of traffic between the switches and the controller which will be an overhead in the network. Hence the flow tables should be able to adapt to these situations without inducing delay in decision making process.

With these two basic design requirements we analyse the cooperative femto network with basic configuration and based on the analysis we design a load balancing method in the controller.

The evaluation scenarios are summarized in Table 2-1

| | Scenario 1 | Scenario 2 | Scenario 3 | Scenario 4 |
|---|---|---|---|---|
| **Routing** | Spanning Tree | Spanning Tree | Arbitrary Routing | Arbitrary Routing/ Load Balancing |
| **Flow** | Single flow | Flow based on (DSCP, VLTAG) tuple | Flow based on (Femto IP, Gateway IP) | Open switch Scenario, with flow based on TEID |
| **No. of Queue** | 1 | 4 | 4 | 4 |
| **Queuing Principle** | FIFO | SP and WFQ | SP and WFQ | Priority |

**Table 2-1 Evaluation Scenarios**

### 2.1.2 Simulation Setup

The analysis and implementation of the scenarios mentioned in the previous section are carried out in the simulation environment of NS3. The following works are in process

a) In the current implementation of NS3 only droptail queuing mechanism is implemented. This implementation is sufficient to carry out the simulations for the baseline scenario. The remaining queuing procedures are being implemented so that scenarios 1 and 2 can be validated.

b) NS3 has facility to integrate the Openflow implementation which has been tested for a group of switches connected to a single learning controller. Based on this study, a load balancing controller is going to be implemented.

## 2.2 Distributed Routing

One important concern of any MNO is to minimize the cost of deploying infrastructure for the backhaul. With the foreseen decreasing prices of the mobile data, this concern will be a key issue to the MNOs to keep profitable mobile networks. Currently, such infrastructure is mainly based on leased T1/E1 lines. The cost of leased T1/E1 lines for backhaul is motivating MNOs to switch to other more cost effective packet-based backhaul. This is even more important in zones where the Average Revenue per User (ARPU) is low.

As a potential solution, MNOs are considering is to combine a packet-based switched network with wireless technologies (more cost-effective in terms of deployment cost) for providing the backhaul. A low cost all-wireless network of femtocells (NoF) tries to cover the necessity of giving cellular coverage in such zones when the deployment of fiber/copper would be economically unfeasible (e.g., for a temporary deployment, or a fixed deployment in a rural zone). As a result, MNOs have an increasing interest in the use of other wireless technologies as backup technology, or even as primary backhaul due to the high cost and unfeasibility of deploying wired infrastructure.

In this context, we conceived and developed a distributed routing protocol suited to such challenging scenarios. If optimally tuned, this can entail several benefits for the MNO in terms of Operational Expenditure (OPEX). In this section we are interested in studying the performance of Wi-Fi as a potential backhaul solution. So, the problem we address is to what extent Wi-Fi can fulfil the requirements of network operators (local and cellular) when used in an all-wireless network of femtocells. Our target is to evaluate its feasibility and to give hints on under what conditions (number of nodes, traffic load, number

of gateways, etc.) these requirements may be fulfilled. Notice that in this context, backhaul is understood as the wireless multi-hop network interconnecting the femtocells that form the NoF.

This section is organized as follows. In subsection 2.2.1, we summarize the work done during project year 1 on the design and preliminary evaluations of the distributed routing solution explained in section 2.2.2. After that, some practical issues to be considered with regards to the distributed routing algorithm are explained in section 2.2.3. A preliminary evaluation carried out for all possible traffic patterns in an all-wireless NoF (i.e., from gateway to femtocells, from femtocells to gateway, and local traffic to the NoF) stressing the considerations explained in previous subsections is given in section 2.2.4. Finally, subsection 2.2.8 concludes this work.

## 2.2.1 Work during Year 1

During the first project year, an extensive study of the state of the art was carried out. As a result of this analysis, a practical solution based on dynamic backpressure routing was proposed. In the proposed solution, a new network element called Local Femtocell GateWay (LFGW) is introduced within the local femtocell network. The LFGW is transparently inserted into the S1 interface between the femtocells and the EPC such that neither the femtocells nor the EPC must be changed. We designed the main building blocks of the distributed routing protocol (see D5.1) to provide routing in the backhaul using inexpensive wireless technologies, whose main concepts are summarized below:

1) The use of a stateless and distributed routing protocol based on dynamic backpressure routing in order to maximize load balancing in an all-wireless NoF.

2) The use of geographic coordinates to assist dynamic backpressure routing in order to head packets to the destination, hence potentially allowing support for all possible traffic patterns (i.e., uplink, downlink, and local).

3) Transparency of the distributed routing protocol with regards to the 3GPP procedures.

Furthermore, in previous work, we carried out an evaluation of the algorithm mainly focusing on uplink traffic routing (i.e., traffic from the femtocells and to a single LFGW). The simulation results obtained in ns-3 showed improvements in terms of aggregated throughput and delay with respect to current topology-based state of the art routing protocols for similar environments. Hence, we validated the correct operation of the routing protocol in a many-to-one scenario.

## 2.2.2 The Dynamic Backpressure Routing Algorithm for an all-wireless NoF

Our distributed routing scheme is based on previous work on backpressure by Georgiadis et al. [68]. As such, it inherits some of its goals (i.e., minimization of a target function while maintaining finite queue backlogs) as well as the theoretical framework for studying it. The first year of the project served to validate the applicability of a distributed version of such concepts in an all-wireless network of femtocells. The goal for the second year is to identify the main parameters and functions of interest in a network of femtocells. This should enable designing a sound protocol for such environments as well as to study what bounds on performance can be offered to operators.

The routing problem in the NoF can be defined as the following stochastic network optimization problem: minimize the time average of objective function $y(t)$ subject to maintaining strong stability (see definition below) in the queues of the nodes, i.e.,

$$Minimize: \quad \bar{y}$$

$$Subject\ to:\ maintain\ NoF\ strongly\ stable \qquad (1)$$

The time average of the objective function $y(t)$ depends on the routing control decisions in the NoF. Specifically, $y(t)$ denotes the cost of performing routing control decisions. Let $Q_i(t)$ denote the queue backlogs of femtocell $i$ at timeslot $t$. The queue backlogs in a NoF are strongly stable if all the femtocells $i \in NoF$ satisfy the following expression:

$$\lim_{t->\infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \mathrm{E}[Q_i(\tau)] < \infty, \forall i \in NoF$$

The optimization problem in (1) can be solved by means of Lyapunov theory. Specifically, a Lyapunov function $L(t)$ defined as the sum of squares of backlog in all queues at slot $t$ is used. Basically, it is a function that grows whenever at least one data queue in the network grows. Then, we define $\Delta(t)=L(t+1)-L(t)$, as the difference of the value of the Lyapunov function in slots $t+1$ and $t$, also referred to as the

Lyapunov drift. Therefore, in order to satisfy queue stability constraints, the NoF takes routing control actions that minimize the Lyapunov drift.

The minimization of the Lyapunov drift is used to solve the queue stability constraint in (1). In order to include the objective function $y(t)$, this function is mapped to an appropriate penalty function $p(t)$.

Therefore, the optimization problem formulated in (1) becomes a minimization of the difference of two functions,

$$\Delta(t) - V \times p(t), V \geq 0$$

where $V$ is a control parameter weighting the penalty function $p(t)$ with respect to $\Delta(t)$. Hence, if $V=0$ routing actions exclusively focus on stabilizing the NoF. On the other hand, if $V>0$ there is a trade-off between minimizing the Lyapunov drift and the penalty function. Stochastic network theory results show that the time average of the objective function deviates by at most O(1/V) from optimality, with time average queue backlog bound of O(V) [28].

The objective function (and so, the penalty function) in our case targets the minimization of the time average end-to-end delay in the NoF.

The penalty function can potentially restrict the network capacity region, especially as the $V$ value grows. The network capacity region is the set of feasible input rates the NoF can service. On the other hand, while a routing protocol merely based on Lyapunov drift minimization can exploit the whole network capacity region, it can also incur into large end-to-end delays in the network. For instance, if the NoF is empty or lowly loaded, a routing protocol merely based on the Lyapunov drift minimization might make packets take a random path to the intended destination. The purpose of this penalty function is to head the packet towards the intended destination, while still maintaining queue stability in the NoF.

In the next section, we provide a summary of how to try to minimize this drift-plus-penalty expression for a NoF in a distributed way (i.e., without requiring a central entity with all the network state information).

### 2.2.2.1 HeNB Link Weight Computation

The minimization of the Lyapunov drift $L(t)$ is associated to minimizing the difference of data queue lengths between femtocells. And, the penalty function $p(t)$ focused on minimizing the average end-to-end delay is associated to a function that depends on the Euclidean distance of the femtocells to the intended destination.

More specifically, given a local femtocell $i$, and its set of 1-hop reachable neighbours $N_i$, the selected neighbouring femtocell $j$ is the femtocell that maximizes the link weight $w_{ij}$ between the local femtocell $i$ and all the femtocells $j \in N_i$, where link weight $w_{ij}$ is calculated as follows. For each link $(i, j)$, the accumulated data packets of node $i$ and node $j$ are stored in queues $Q_i$ and $Q_j$, respectively. Notice that the number of stored data packets in a queue depends on the difference between serviced packets, and endogenous and exogenous arrivals at each femtocell. Let $\Delta Q_{ij}$ denote the difference of physical queue lengths between femtocell $i$ and femtocell $j$.

With regards to the penalty function, in the link computation, we included a simple *penalty* function that reduces the weight of links with those neighbours that make the packet move away from the intended destination. In this case, we want to penalize packets that incur into excessive end-to-end delays. We do this by using the Euclidean distance in order to have a sense of proximity to the intended destination. Let $p(i,j,d)$ denote the penalty function computed as the cost to traverse the link between femtocell $i$ and femtocell $j$ to reach destination $d$. In our preliminary evaluation, we defined the following simple penalty function $p(i,j,d)$ as follows:

$$p(i,j,d) = \begin{cases} -1 & if\ i\ is\ farther\ from\ d\ than\ j\ in\ terms\ of\ euclidean\ distance \\ +1 & if\ i\ is\ closer\ to\ d\ than\ j\ in\ terms\ of\ euclidean\ distance \end{cases}$$

Moreover, the penalty function is parameterized by a non-negative control parameter $V \geq 0$. This parameter indicates the importance given to the penalty function when calculating the weight of each pair $(i, j)$ with respect to the difference of physical queue lengths. In our case, the penalty function results in two possible different values. Its value is 1 when $j$ is closer to $d$ and -1 when $j$ is farther from $d$. Therefore, the penalty function maximizes the weight when $j$ is closer to $d$, and minimizes the link weight when $j$ is farther from the destination than the local femtocell $i$. On the other hand, it is minimized when the distance to the destination is reduced by a given neighbour with respect to the local femtocell. As a result, the weight $w_{ij}$ of a link between wireless link $(i, j)$ is calculated as follows:

$$w_{ij} = \Delta Q_{ij}(t) - V \times p(i,j,d)$$

### 2.2.3 Some Issues to Consider on the Proposed Algorithm

We have identified three main issues to take into account in the preliminary evaluation of the proposed distributed routing algorithm into an all-wireless NoF. First, the trade-off proposed by the *V* parameter between delay and queue stability. Second, the impact of the parameter *V* has to determine which set of input rates are feasible or not, and under which circumstances. Third, the management of a NoF deployed with multiple LFGWs. And finally, how this distributed geographic routing algorithm could react to dead-ends (i.e., managing the case in which a femtocell does not have neighbours closer to the intended destination) in sparse NoF deployments.

### 2.2.3.1 The importance of choosing an appropriate V parameter

In this section, we will illustrate through the use of heatmaps, the importance of the *V* parameter. The *V* parameter determines a trade-off between queue stability and end-end delay. To illustrate this fact, we send a single flow from femtocell 40 to femtocell 49 in the NoF of 100 nodes in Figure 2.5. A unidirectional flow is sent from the source node in position (0, 4) of the grid towards the node in position (9, 4). As it can be observed in in Figure 2.6, Figure 2.7, and Figure 2.8, *V* determines the relative importance given to 1) the directionality of forwarding decisions towards the destination vs. 2) queue backlog reduction. That is, lower values of V make our scheme distribute more the traffic, but it makes packets to follow longer paths. On the other hand, bigger values of *V* foster shorter paths, but queue backlogs grow at nodes traversed by this flow.



**Figure 2.5: All wireless Network of Femtocells**



**Figure 2.6: Data Packet Distribution V=1**



**Figure 2.7: Data Packet Distribution (V=256)**



**Figure 2.8: Data Packet Distribution V=1024**

### 2.2.3.2 The feasibility region

The input rate matrix can be defined as the set of source-destination pairs and their respective traffic intensities injected to the NoF. An identified issue is the observation that the input rate matrix is not only important in terms of the sum of the traffic intensity, but also in terms of the combination of source-destination pairs. The input rate matrix can be inside or outside the network capacity region depending on both the traffic intensity and the combination of source-destination pairs. Therefore, both components (i.e., traffic intensity and the set of source-destination pair) determine if an input rate matrix is feasible or not.

### 2.2.3.3 Multi LFGW (as an issue to exploit to increase the feasibility region)

The introduction of several LFGWs may potentially increase the feasibility region. It is well known that a capacity of wireless links of LFGW is a potential bottleneck, as they act as interconnector of the HeNB with the S-GW. As a result, as the number of LFGWs capable of pulling packets from the NoF grows, the feasibility region might also potentially grow. The specific algorithm used for LFGW load-balancing is described in section 2.2.6.

### 2.2.3.4 Dead-Ends

This solution explores the combination of backpressure routing with the use of geographic coordinates. One issue in the context of geographic routing is the dead end problem. The dead-end occurs when there is no neighbour closer to the destination in terms of Euclidean distance. A study of interest is to see how/whether backpressure could circumvent dead-ends in the NoFs. The queue backlog gradients established by dynamic backpressure routing may provide a very simple distributed solution to the problem of dead ends in sparse deployments. In section 2.2.7, we provide an initial study by means of assuming sparse deployments with potential dead ends in the all-wireless NoF under evaluation.

## 2.2.4 Methodology and Performance Evaluation of the routing algorithm with constant-V parameter, single LFGW, and no holes

Simulations in this paper were carried out with the ns-3 simulator. The NoF model used is a grid of 100 femtocells connected by a cheap Wi-Fi backhaul based on 802.11a. In this section, we will show, first, results that we have on the evaluation of the routing protocol for all the possible traffic patterns in the NoF. We carried out five experiments in which two randomly chosen femtocells injected five different UDP CBR traffic intensities (i.e., 500Kbps, 1Mbps, 1.5Mbps, 2Mbps, 2.5Mbps, and 3Mbps) towards two randomly chosen femtocells. In particular, we carried out 24 replications for each experiment. Moreover, every experiment was repeated for 12 different $V$ parameters ranging from 1 to 500. As a result, a total of 1728 experiments were performed, in which the experiment duration was of 300 seconds. The data flows were injected during a window interval of 150 seconds. All the measurements of the parameters under evaluation were obtained during this window interval of 150 seconds. The data queue length of the femtocells was configured to 400 packets, and the queuing policy used is FIFO.

### 2.2.4.1 Network Throughput



**Figure 2.9: Average Network Throughput in the NoF.**



**Figure 2.10: Standard Deviation of Average Network Throughput in the NoF.**

Figure 2.9 plots how the average network throughput for several input rate matrices (i.e., two sources injecting 1.5Mbps, 2Mbps, 2.5Mbps, and 3Mbps to two destinations). Figure 2.10 plots the variability of throughput with respect to the average value for the previous input rate matrixes.

For low $V$ values the network capacity region is maximized, but not all packets arrive to the intended destination during the measured window interval of the network throughput, as the routes taken by the packets are close to a random walk in the NoF. This causes a variable throughput degradation that depends on the input rate matrix (see Figure 2.10).

On the one hand, this variability depends on the specific combination of the source and destination pairs. For instance, if the source femtocell of flow 1 is close to the destination femtocell of flow 2, the throughput attained by flow 2 would not be feasible for low $V$ values. Hence, flow 1 can be seen as a noise source that avoids the reception of packets by the destination femtocell of flow 2. This effect affects the offered input rate matrix, hence leading to severe throughput degradation in the NoF. On the other hand, if neither flow 1 nor flow 2 had such dependence, there are more chances for the destinations to receive data packets.

Additionally, the distance between flows in the NoF also influences the received data rate in the corresponding destinations. The closer the chosen source and destination femtocell nodes are located, the

more likely to increase the rate of the received packets at the destination, and so, the average network throughput.

As *V* grows (see Figure 2.9 and Figure 2.10), packets use paths close to the shortest path to the destination. At this stage, the algorithm provides a good trade-off between the network capacity region and direction towards the intended destination, providing a feasible solution to serve the injected input rate.

However, with the input rate matrix maintained constant, when choosing high V values the network is not able to serve the injected input rates. Specifically, Figure 2.9 illustrates that for 6Mbps and 5Mbps there is a slight decrease of the achieved network throughput. This is because the network capacity region is shrunk to an extent in which the 6Mbps and 5Mbps input rate matrix cannot be totally served. This is caused by data queue drops due to queue overflows in the femtocells, hence breaking the strong stability constraint to be maintained in the NoF. Interestingly, the data queue drops present variability depending on the combination of source-destination pairs. We have evaluated this for high *V* values. We observe that the more femtocells in common by the two flows, the more queue drops occur in the NoF. Therefore, variability depends on the amount of resources (femtocells) shared by the flows.

Another aspect to remark is that the average network throughput variability seems to grow as the absolute injected input rate in the NoF grows. The increase of the offered load emphasizes the throughput degradation pathologies previously explained (i.e., lack of direction, queue drops).

### 2.2.4.2 Flow Throughput

To evaluate the throughput attained by each flow independently, we have carried out an experiment in which the offered load in the NoF is 6Mbps (i.e., 3Mbps injected by each flow). We generated various combinations of two flows that send traffic from a random source to a random destination. Each of these combinations is identified with a flow index value in the figure. One of such two flows is arbitrarily tagged as flow A and the other one of the same flow index is tagged as flow B. We present the throughput attained by each separate flow as well as its evolution with the *V* parameter in two different graphs (for flow A and for flow B). For each flow, we have a different flow index, which represents a different random source destination femtocell pair.



**Figure 2.11: Throughput attained by Flow A in the NoF**



**Figure 2.12: Throughput attained by Flow B in the NoF**

As can be shown, in Figure 2.11 and Figure 2.12 both flows receive the offered load in most of the cases under evaluation. The differences reside in the *V* parameter under evaluation. As explained in the

previous section, we have identified two main reasons with regards to the convergence to the offered load.

An important aspect is the distance between the source and the destination, the shorter the distance, the lower the required *V* to converge to the injected flow throughput. On the other hand, the potential dependencies between flows injected to the NoF are an important aspect to take into account. For instance, in Flow index 14 there are two flows in opposite directions (flow A between nodes 97 and 11) and (flow B node 13 and 44). As these flows are sharing a slice of the NoF but they are sending data packets in opposite directions, they require a higher *V* value in order to converge to the offered load of the NoF.

For not so low *V* values (i.e., *V=25*), with respect to the throughput attained by other source-destination femtocell pairs with *V=25*, Figure 2.11 and Figure 2.12 show that a low percentage of packets arrive to the intended destination femtocell (throughput is 0.8Mbps and 1.6Mbps for flow A and B, respectively).

In order to understand what is exactly happening in this case, we plot the heatmap of average data queue backlog, which illustrates this issue for *V=*25 (see Figure 2.13). Specifically, Figure 2.13 illustrates that most of the packets are accumulated in data queues without reaching the respective destinations. Both flows are, to some extent, blocking themselves, hence avoiding that a higher percentage of packets could be received in the respective destination femtocells. Furthermore, flow A experiences a higher throughput degradation than flow B, as the distance between the femtocells is bigger in the case of flow A than in flow B (i.e., 14 hops versus 4 hops).

As *V* grows (see heatmap for *V=100* in Figure 2.14), both flows do not block themselves, and so, they can reach the intended destinations, since they tend to follow paths closer to the shortest one. As shown in Figure 2.14, the femtocells do not suffer as much accumulation in their data queue as in the case *V=25*. This is because packets have a bigger sense of direction towards their respective femtocell sinks.



**Figure 2.13: Heatmap of Average Data Queue Backlog V=25**

**Figure 2.14: Heatmap of Average Data Queue Backlog V=100**

### 2.2.4.3 End-to-end Delay



**Figure 2.15: Average End-to-end dealy in nanoseconds**

**Figure 2.16: Standard Deviation of Average End-to-end delay in nanoseconds**

**Figure 2.17: Average End-to-end jitter in nanoseconds**



**Figure 2.18: Standard Deviation of End-to-end jitter in nanoseconds**

Figure 2.15 plots the average end-to-end delay versus a range of 12 *V* values between 1 and 500. As expected, average of end-to-end delays decrease as the control *V* parameter is increased. Specifically, the function to be minimized, i.e., the average end-to-end delay attained by the NoF, decreases with *O(1/V)*.

As can be shown in Figure 2.15, as the input rate injected into the NoF grows, the NoF requires a higher *V* parameter to experience the minimum values of end-to-end delay. This means that higher *V* values are needed to converge to the minimum end-to-end delay. Therefore, as Figure 2.15 shows, for low *V* values, the end-to-end delay convergence time will increase with the input rate injected to the network, as there is a higher percentage of traffic load not served and queued at HeNBs. This effect causes an increase in the queuing delays with the input rate injected to the network as a high percentage is unable to be served due to lack of direction, and hence it is kept at HeNB data queues.

On the other hand, as Figure 2.16 depicts, variability with respect to end-to-end delays grows with the input rate, given the higher number of packets injected to the network. This is caused by the arrival of packets to the destination femtocells with higher waiting time at data queues. Moreover, standard deviation with respect to average end-to-end delays (see Figure 2.16) is bigger for low *V* values, as the randomness in the paths taken is higher. Lack of variability with respect to input rate matrix for sufficiently high *V* values, as illustrated in Figure 2.16, validates the fast convergence time of the routing algorithm.

Figure 2.17 presents the average jitter data packets experience. Similarly to end-to-end average delays, the average jitter decreases when *V* increases given the increase of direction given to the packets. As with the end-to-end delays, the jitter also decreases with *O(1/V)*. Figure 2.18 plots standard deviation against the V control parameter. The convergence time to the optimum values also increases with the input rate, given the higher number of packets the NoF has to handle.

### 2.2.4.4 Fairness study

We consider the NoF as fair if, during the interval, the number of packets exchanged by each source destination pair in the NoF is approximately the same. Hence, we see the network as fair with regards to

the throughput attained by each flow. We use the Worst Case Fairness Index (WFI) in order to compare the throughput attained by both flows injected to the network. The WFI is computed as follows:

$$WFI = \frac{\min_{1 \leq i \leq n} \{x_i\}}{\max_{1 \leq i \leq n} \{x_i\}}$$

where $x_i$ is the ratio between the measured throughput and the expected throughput of flow $i$. Since there are no losses over-the-air, the expected throughput should ideally be equal to the input rate (except for queue overflows). The WFI is appropriate to detect differences in how the network handles a small number of flows.



**Figure 2.19: Worst Case Fairness Index**

To evaluate fairness under we choose a input rate matrix when the total offered load is 6Mbps. The offered load injected by each single flow is 3Mbps. Note that as explained before, the *V* parameter determines whether the offered load is inside or outside the network capacity region (i.e., the input rate is able to be served or not). Figure 2.19 clearly demonstrates two facts. Precisely, only two of the twenty-four pair of flows under evaluation are able to get a fairness index over 0.85. The rest of the flow pairs experience between 20% and 100% of variability with respect to the *V* parameter. As can be shown in Figure 2.19, there are some cases ranging from *V = 50* up to *V = 100* in which the network experiences a degree of fairness bigger than 0.88. This happens when the network achieves a good compromise between direction and network capacity region. And secondly, the impact of the location of the source-destination flows on the fairness in the network. The WFI shows full variability for low values of the *V* parameter given the lack of direction of data packets. Due to the lack of direction by data packets, this causes the network to be unable to deliver the packets within the measured window interval, even though the network has sufficient available network resources.



**Figure 2.20: Average Queue Length per second during the measured windows interval**

### 2.2.4.5 Summary of Main Findings (summary of practical implications of Lyapunov optimization framework in the NoF)

We point out the following statements from the previous study that came out from the use of a routing control policy with a constant *V* parameter.

#### 2.2.4.5.1 *The dependence of network performance metrics on V*

We study the dependence of network performance metrics on *V*. As an extreme case, we show the limitations of making forwarding decisions in any-to-any scenarios only based on the Lyapunov drift without pre-computed end-to-end routes. This comparison also serves to show the role of the penalty function in steering packets towards the intended destination. As shown in Figure 2.9, the NoF suffers

from performance degradation for values of *V* ranging from 1 to 25. The reason for the low values of throughput and high end-to-end delay in Figure 2.9 and Figure 2.15 is the absence of pre-computed end-to-end routes and the fact of only having one finite queue at each node, which breaks the assumptions made by theoretical centralized backpressure routing. This leads packets to take close to random walks to reach their intended destination during the time interval of the measurement, resulting in high delays for delivered packets within the time interval of the measurement.

Because of this, a high percentage of data packets are still at data queues of nodes in the NoF. For instance, we observe that an algorithm taking forwarding decisions based on queue differences between neighboring nodes almost exclusively (i.e., *V=1*) just obtains a 43% and 32% of packet delivery ratio during the measured time interval for the 3Mbps and 6Mbps case, respectively. Likewise, this randomness in the path followed by packets leads to randomness in the throughput attained by each flow, which, in turn, leads to a high WFI variability (e.g., for *V = 1*, WFI ranges from 0 to 0.96).

### 2.2.4.5.2 *The objective function-backlog trade-off with finite queue Sizes*

Recall that the routing control policy is designed so that V gives relative weight to the objective function with respect to the Lyapunov drift (see equation). The aim of the distributed penalty function is to minimize the distance covered by packets to reach the destination, which is related to end-to-end delay, even more considering our simple channel model with no interference. Hence, there are only retransmissions for collisions due to simultaneous backoff exhaustion. In [68], Neely proofs that there is an objective function backlog trade-off of [$O(1/V)$, $O(V)$]. In our network setup, having an objective function related to Euclidean distance and end-to-end delay is equivalent. Based on our simulations, Figure 2.15 illustrates that the average end-to-end delay evolves with $O(1/V)$. In fact, as *V* increases, the forwarded packet makes more progress on average towards the intended destination. However, having size-limited data queues avoids experiencing high queuing delays for high values of *V* at the cost of experiencing potential queue overflows. With regards to queue backlogs, we highlight two observations. The first observation is the high average queue backlog when Lyapunov drift minimization has more weight than distance minimization. This might seem contradictory with the theoretical trend of $O(V)$. Nevertheless, the reader should notice that, in this case, all nodes in the network are potential candidates to forward a data packet without even trying to restrict potential routing loops. Conversely, when the set of potential paths is restricted (i.e., higher *V*), Neely's growth with $O(V)$ of the average queue backlogs (and so average time in queues) is satisfied up to the bound determined by the network (i.e., 400 packets). After this bound, queue drops occur, and so, the $O(V)$ trend is not satisfied.

### 2.2.4.5.3 *The dependence of network performance metrics on the relation between V and the queue size*

Figure 2.9 illustrates that for 5Mbps, and especially for 6Mbps, there is a decrease in the measured network throughput (averaged over all input rate matrices) for high values of V, to an extent in which the injected traffic cannot be served. The reason for this reduction is that as V increases, the routing policy restricts the set of candidate nodes used to forward a data packet, given the increasing importance of the penalty function with respect to the Lyapunov drift. An illustrative example of such reduction can be found in Figures 2.6 to 2.8. Restricting the set of nodes used to forward packets incurs into higher congestion of the nodes traversed by the flows, given the higher number of data packets they have to handle. This causes an increase of their queue backlogs, which may lead to queue overflows. This dependence between the queue size and *V* can be shown in other network metrics, such as fairness. Figure 2.19 shows that the network fairness depends on *V*. More precisely, only two of the twenty-four pairs of flows under evaluation are able to achieve a fairness index over 0:85 for all tested values of *V*. The rest of the flow pairs experience between 20% and 100% of variability with respect to *V*. On the other hand, and as shown in Figure 2.19, there are some cases (from *V = 100* to *V = 150*) in which the network presents a WFI higher than 0.94, no matter the input rate matrix. Furthermore, for these cases, network throughput and delay are close to the best achieved values. The reason is that these values of *V* allow attaining an appropriate tradeoff between directionality of flows and congestion at nodes, but without exceeding the queue size. On the other hand, Figure 2.20(c) shows that for some input rate matrices the network experiences a decrease in WFI for high values of *V*, given the lower number of resources (i.e., nodes) packets are traversing. In this case, there are losses due to queue overflows. Additionally, Figure 2.20(b) shows that depending on the source-destination pair, such unfairness is also observed for low values of *V* (e.g., see input rate matrix indexes 21 and 22). For instance, Figure 2.18b illustrates one of these cases. In this case, the number of nodes used to route the flow originated at node in position (2, 7) towards node in position (2,1) is reduced to an extent in which many packets cannot be delivered. Thus, for the same value of *V*, the use of resources (i.e., the number of nodes traversed) to route data packets when applying the algorithm is different for both flows.

### 2.2.5 Variable-V algorithm Solution

Results previously shown suggest the importance of a variable-*V* algorithm to avoid queue drops at the NoF. Specifically, we have seen that the traffic served could highly vary (in terms of throughput, delay, jitter and fairness) depending on the *V* parameter. In a NoF, it is expected that the input rate matrix would be variable in time. Additionally, mobility of UEs could also lead to changes in the input rate matrix in the NoF.(i.e., the HeNB injecting traffic coming from a given UE can change). Therefore, this suggests that it may be convenient to choose on-demand the more appropriate *V* value in each HeNB in order to avoid queue drops while still minimizing the penalty function (see section 2.2.2 for a more detailed explanation of the routing algorithm).

#### 2.2.5.1 The Problem

This section is devoted to show the complexity of the studied problem. In the first subsection, we describe the problem of experiencing queue overflows due to the *V* parameter taking into account a pair of HeNBs, and studying the queue evolution of one HeNB considering serviced packets by the other HeNB. In the second subsection, we extend the problem for generic NoF in which the HeNB taken as reference to study its queue evolution can have an unlimited number of HeNBs as neighbours.



**Figure 2.21: A System with one reference node and one neighbour**

**Queue Overflow due to V parameter at a HeNB with one HeNBs as Neighbour**

In order to define the problem, we study a simple case in which there are two nodes in the NoF; node *i*, and *j*. As reference node we take node *j*. The goal is to study the queue evolution $Q_j$.of node *j*. As the resulting distributed link weight computation algorithm derived in section 2.2.2.1 specifies, we determine whether node *i* will transmit to node *j* through the following expression:

$$w_{ij} = \Delta Q_{ij} - V p(i,j,d)$$

when $w_{ij} > 0$ the node *i* will transmit to node *j*. At each routing decision, two different cases can happen. If *p(i, j, d) =+1* it means that the node *j* is further from destination *d*, then there is no chance of generating a queue overflow in *j* that legacy backpressure would not experience. It can be seen as sort of additional virtual congestion in node *j*. In this case, from the point of view of transmissions from node *i* to node *j*, the decisions taken will be equivalent to those decisions that a backpressure algorithm will take as well. Thus, we assume that in terms of minimizing queue overflows, the algorithm cannot do any better than legacy backpressure algorithm (i.e., *V*=0). In other words, we assume that if a pure backpressure policy (i.e., *V=0*) cannot avoid queue overflows, a variable *V* algorithm can neither do it.

$$w_{ij} = \Delta Q_{ij} - V; w_{ij} = Q_i - Q_j - V > 0$$

Therefore, as *V>0* this means that $Q_i - Q_j > 0$ with $V > 0$. Hence the original backpressure algorithm would take the same decision. Given the previous assumption, if queue drops occur, these are not avoidable by any other routing policy.

On the other hand, if *p(i,j,d) = -1* the node *j* is closer of destination *d* than *i*. In this case, the V parameter is additive in the computation of the weights. It can be seen as sort of additional virtual congestion in node *i*. As a result of this, there are more chances that transmission from node *i* to node *j* occur that will not happen under the traditional backpressure algorithm (i.e., *V*=0). Therefore, potentially a queue overflow in node *j* that will not happen under original backpressure algorithm may happen:

$$w_{ij} = \Delta Q_{ij} + V; w_{ij} = Q_i(t) - Q_j(t) + V$$

Assuming that node *i* will transmit to node *j* implies that $w_{ij} > 0$. Consequently:

$$Q_i(t) - Q_j(t) + V_i(t) > 0$$

Recall that the main goal is to avoid queue overflows in node *j*. Therefore:

$$Q_j(t) + a_i(t) - b_j(t) < Q_{MAX};$$

Let $H$ be the range of the arrival process $a_j(t)$ number of packets transmitted by node $i$ to node $j$ during timeslot $t$. The range $H$ could vary as follows,

$$H(a_j(t)) = [0, min(Q_i(t), Q_i(t) - Q_j(t) + V_i(t))];$$

where $a_j(t)$ is the arrival process at node $j$ and it comes determined by the number of packets sent by $i$ and $b_j(t)$ which is the process describing the serviced packets by node $j$ is a variable affected by the wireless network conditions. The $V_i(t)$ parameter will then depend on the network conditions of the medium which, in turn, will define the specific value of $H(a_j(t))$ which is not known a priori in a practical system. Nevertheless, if we are assume CSMA/CA medium allows an indefinite number transmission of packets, and $Q_i(t) > Q_i((t) - Q_j(t) + V$ then $a_i(t) = Q_i((t) - Q_j(t) + V$. In this specific case, the resulting expression for $V_i(t)$ comes determined by the following expression:

$$V_i(t) < Q_{MAX} - Q_i(t) + b_j(t)$$

Therefore, we define $V_i(t)$ as the upper bound in the worst-case scenario from the point of view of node $j$ to experience queue overflows. Even in this case in which we know the maximum number of packets that can be transmitted from node $i$ to node $j$, the resulting expression is determined by the process defining the serviced packets $b_j(t)$ in node $j$. Therefore, the number of packets transmitted during the interval *[t, t+1]* in node $j$ should be also known. However, this is not practical since it is not local information, and will not be known until next timeslot *t+1*.

***Queue Overflow at a HeNB with a set of HeNB Neighbors***



**Figure 2.22: System with one reference node (i.e., node j) and a set of neighbours**

In this section, we extend the previous problem to the case in which node $j$ has a set of neighbors that are also potential senders of packets. Recall that the goal is to study the queue evolution of node $j$, and evaluate under which conditions the $V$ parameter must satisfy in order to avoid queue drops. As in the previous case, we fix our attention in the queue evolution of node $j$. Specifically, we describe the conditions any potential neighbor $k$ of node $j$ must satisfy in terms of the $V$ parameter in order to do not experience queue overflows in node $j$. To do so, we extend the conditions the $V$ parameter in node $i$ must satisfy when a set of neighbors willing to transmit to node $j$ is defined.

For simplicity, we assume node $j$ at timeslot $t$ do not experience exogenous arrivals (i.e., node $j$ is not generating packets in upper layers). Then, the range $H$ of the arrival process $a_j(t)$ at any given timeslot $t$ occurs in the $a_j(t)$ component. We are assuming that $V_j(t) <= Q_i(t)$ as in the previous case. Therefore, the range determining the number of packets arriving at node $j$ at timeslot $t$ is defined as follows:

$$H(a_j(t)) = [0, \sum_{i \in N}(Q_i(t) - Q_j(t)) + \sum_{i \in N} V_i(t)]$$

Therefore, as in the previous case we should know a priori the number of packets that node $i$ will transmit to node $j$. Nevertheless, assuming that node $i$ will transmit all its packets to node $j$, the maximum number of packets comes determined by all the data packets all neighbors $i \epsilon N$ which could potentially send packets during timeslot $t$. If we separate the arrivals of node $i$ from the rest of arrivals of the neighbors of node $j$, we can split $a_j(t)$ in two parts:

$$a_j(t) = a_{ji}(t) + a_{jN(j)-i}(t)$$

$$a_{ji}(t) = Q_i(t) - Q_j(t) + V_i(t)$$

$$a_{ji}(t) + a_{jN-i}(t) + Q_j(t) - b_j(t) < Q_{MAX}$$

Let $k$ be the number of neighbors of $j$ minus node $i$ ($N(j)-i$). Therefore, in this case the expression of the upper bound of the appropriate $V_i(t)$ value is defined as follows:

$$V_i(t) < Q_{MAX} - Q_j(t) - \sum_{k \in N-i} (Q_k(t) - Q_j(t))$$
$$- \sum_{k \in N-i} (V_k(t)) - b_j(t)$$

Therefore, it is even more difficult in this case to estimate the $V$ parameter given the increasing number of uncertainties as the number of neighbors of node $j$ increases. There are additional uncertainties with respect to the previous case. Precisely, in this case we also need to know in advance the number of packets transmitted by each one of the rest of neighbors in order to estimate $V_i(t)$.

*Some reasoning and justification of the use of an estimator*

As we have shown in previous subsections, in practice obtaining the optimal $V$ parameter at each HeNB in order to avoid queue drops while still minimizing the penalty function value at each node is a priori unfeasible. In summary, the main issue is that it is not known a priori what will happen during the timeslot interval [$t$, $t+1$] to neighbors $N$ of node $j$. Furthermore, we do not even know a priori the number of packet the local node $i$ will be able to transmit.

On the one hand, this requires to have knowledge of future events (i.e., the serviced packets during timeslot $t$ is $b_j(t)$), and also information which potentially may not be present at 1-hop neighbors. In other words, the number of serviced packets by any of the nodes in a wireless medium (e.g., CSMA/CA medium) is quite uncertain during timeslot interval [$t$, $t+1$], given issues such as unfairness and unreliability of the wireless medium.

On the other hand, recall that the goal is to have a distributed and practical algorithm to reduce the queue drops in the network, and at the same time deliver data packets to the intended destination in the NoF. To do so, we propose a distributed estimator of the $V$ parameter based on current information state available at 1-hop neighbors. Precisely, we approach the problem with a cautious and easy-to-implement estimator of the appropriate $V$ parameter in every HeNB.

### 2.2.5.2 Practical Variable-V algorithm: Aimed properties to satisfy

The goal is to design and implement a controller that auto-configures the $V$ parameter in every HeNB in the NoF. To do so, we propose to design a controller that must satisfy the following characteristics:

- The $V$ controller must avoid queue drops in the network if a pure legacy backpressure algorithm is able to avoid them (i.e., routing control policy with $V=0$).
- The $V$ controller must prioritize the penalty function whenever light loads (indicating low or null queuing) is detected in the NoF.
- The $V$ controller must be adaptive and also be able to react under varying traffic conditions (e.g., a new flow is injected in the NoF, or a new flow stops).

### 2.2.5.3 Practical Variable-V algorithm: Main Intuition behind the proposed controller

In order to autoconfigure the $V$ parameter in every HeNB, we built a virtual queue (see Figure 2.23) describing network conditions in terms of network load exploiting 1-hop information of the data queues at every timeslot $t$. In other words, we are aggregating the information gathered from 1-hop HeNBs in terms of congestion.

**Figure 2.23: Virtual data queue describing the aggregated 1-hop data queue length state**



**Figure 2.24: Distributed V-Variable Controller**

In order to build this aggregated queue and so get the new value of the $V$ parameter at a given node (see Figure 2.24), there are three basic components locally accessible to take into account:

$Q_{REF:}$ This is a constant value of the system corresponding to all the nodes in the NoF that denotes the data queue length limit allowed at a HeNB to do no experience queue drops.

There is a significant disparity in the buffer size limit used in various research platforms. Many researchers using the ns-2 simulator use a buffer size of 50 packets, the default size for the queue object in ns-2. In contrast, ns-3 network simulator uses a queue size limit of 200 packets that can easily be modified.

On the other hand, the legacy open source madwifi drivers for Atheros chipset use a device driver ring buffer of 200 packets. The new ath5k drivers divide this 200 packet buffer1 equally among four queues representing the traffic access categories as defined in Enhanced Distributed Channel Access (EDCA) mechanism. The ath9k drivers for 802.11 cards based on Atheros chipsets increase the buffer size. Specifically, they use a buffer of 512 packets, again equally divided among the four ACs. In addition, the Linux network stack introduces other layers of buffers, including the network stacks transmit queue, typically set to a default value of 1000 packets. We want to take a look at the problem from a different perspective. Given a determined buffer size maximize the offered load the network can serve with the help of the routing protocol. In this study we opt for the most common buffer size of the madwifi legacy drivers (i.e., buffer size of 200 packets) which is the default value used in the ns-3 network simulator.

$max(0,Qj(t)-Qj(t-1))$: This component summarizes the previous event experienced in the neighborhood of a given node that leads to the most increasing queuing length with respect to previous timeslot. This is usually caused by the injection of a new flow or set of flows in the network, the increase in the offered load of an existing flow, or even the failure of a given HeNB which may cause other HeNBs to increase

their load. It is calculated as the maximum increase experienced in the data queue length of a neighbor between two timeslots. The estimator uses this component to estimate the future next event that will cause an increment of data queue lengths. Basically, we are assuming that the bigger increase in data queue lengths during time interval *[t ,t+1]* is the same as the one experienced during time interval *[t-1,t]*. This component can be calculated in a distributed manner by means of storing the queue lengths of every HeNB neighbor experienced in two consecutive timeslots.

*max($Q_k(t)$):* This components describes the HeNB neighbor with maximum data queue length. Note that the HeNB experiencing the maximum data queue length in the 1-hop neighborhood could be different from the HeNB which experiences the bigger increase during the previous timeslot. For instance, this could happen if a HeNB do not have enough transmission opportunities due to the CSMA/CA medium. Thus, in this case the HeNB keeps its data queue highly loaded.

$V_i(t)$: This corresponds to the maximum number of packets that potentially can be transmitted in the aggregated virtual queue without causing a queue overflow. As a result, the sum of the previous three components corresponds to the queue limit of the HeNB in the NoF:

$$Q_{REF} = \max(0, Q_j(t) - Q_j(t-1)) + \max(Q_k(t)) + V_i(t)$$

### 2.2.5.4   Practical V-variable algorithm: Details on the controller



**Figure 2.25: Virtual Aggregated Queue Details**

As shown in pseudo code depicted in Figure 2.26, at instant *t=0* all the nodes in the network compute a *V* parameter which denotes maximum level of greediness with respect to the stress in the penalty function. In our case, this means a high degree of directionality towards the intended destination when taking routing decisions.

Furthermore, let $Q_{GREEDLIM}(t)$ a parameter calculated within the controller shown in Figure 2.24, which denotes the queue length that once exceeded forces the local node to use a routing policy equivalent to the legacy backpressure routing algorithm (i.e., *V=0*). Therefore, when the queue length is below $Q_{GREEDLIM}(t)$, the local node still has some chances of computing a local routing decision that take into account the penalty function (i.e., the *V* parameter is greater than 0).



**Figure 2.26: Pseudo Code of the variable-V controller**

On the other hand, $Q_{GREEDLIM}(t)$ also defines the maximum value the variable V can attain. This limit is calculated every time slot given the actual local vision of the maximum traffic load in the neighborhood of $i$ as described in the pseudo code of the controller shown in Figure 2.26. Precisely, it is calculated as the difference between $Q_{REF}$ and the maximum increase experienced by a data queue of a node in the 1-hop neighborhood.

If a node does not experience a data queue length bigger than $Q_{GREEDLIM}(t)$ a node can greedily transmit up to V packets which is defined by $Q_{GREEDLIM}(t),- max(Q_k(t))$. On the contrary, if a node in the neighborhood experiences data queue length bigger or equal than $Q_{GREEDLIM}(t)$ the distributed algorithm consider the network should focus on the minimization of the Lyapunov drift. In this case, the algorithm is detecting that the local neighborhood, and decides to act at the maximum degree of load balancing with the unique goal of minimizing queue drops. The local vision of the traffic can be seen as the maximum differential of queue lengths experienced by a node in a neighborhood (i.e., $max(Q_k(t))$ in Figure 2.26).

Therefore, the determination of the value of the $V$ parameter at instant $t$ comes determined by an expression merely associated to the data queue lengths of the set of neighbors:

$$Q_{jmax}(t) = \max Q_j(t) \in N_i;$$

$$V(t) = \max(0, Q_{greedlim}(t) - \max Q_{jmax}(t))$$

Therefore, the maximum $V$ values come determined by $Q_{GREEDLIM}(t)$, which defines the maximum value the $V$ parameter must have. On the one hand, $V_i(t)$ can be seen as the maximum number of packets that greedily will be transmitted without taking into account load in the neighborhood. On the other side, we established that $Q_{GREEDLIM}(t)$ is the threshold from which a pure backpressure algorithm must be executed in order to avoid queue drops. Consequently avoiding the greed transmission of more than $Q_{GREEDLIM}(t)$ $Q_{JMAX}(t)$ data packets during a timeslot interval for preventing queue drops. We are assuming that the probability of transmitting more than $V$ packets greedily is zero. If any neighbor $j$ has a $Q_j(t) > Q_{GREEDLIM}$ the legacy backpressure is the routing policy used. Note that this change is not propagated over the whole network. It is just performing legacy backpressure in the HeNBs which observe high degree of congestion in its neighborhood. It is expected that the packets will eventually find HeNBs in the network with lower congestion levels, and so bigger $V$ parameters can be used.

### 2.2.5.5 Practical Variable-V algorithm: Evaluation Methodology

We launch a CBR unidirectional data flow F1 (see Figure 2.27) of an intensity able of saturating data queue backlogs of reference node (i.e., node R in Figure 2.27). F1 is launched at instant $t$=5s until instant $t$=80s from node R towards node D of the simulation. We are considering the 10x10 grid NoF depicted in Figure 2.27. The relevant part of the NoF to describe the experiment can be shown in Figure 2.27. Precisely, node D is located at 6 hops from node R. Furthermore, we launch another CBR unidirectional data flow F2 originated in node C directed to node D from instant $t$=20s up to instant $t$=60s in the simulation. F2 also saturates data queue backlogs of node C. Therefore at instant $t$=20s in the ns-3 simulation there is change in the offered load injected to the NoF. This can cause a very different behavior in the evolution of data queue lengths in the NoF, and so illustrates a use case of the presented variable-$V$ controller. We compared the behavior of the variable-$V$ controller policy against, the routing protocol under fixed-$V$ policies (i.e., $V = 0$, $V = 50$, $V = 100$, $V = 150$, and $V = 200$).



**Figure 2.27: Network Scenario Evaluating the Variable-V algorithm**

To do so, we evaluate the behavior of some critical nodes during time participating in the routing of both flows in terms of certain metrics: data queue backlogs, queue drops, and value of the $V$ parameter. Note that data queue drops corresponding to exogenous arrivals at node 0, and 1 are not taken into account. Precisely, the focus of the variable-V algorithm is not on regulating exogenous arrivals but endogenous arrivals. The regulation of exogenous arrivals is out of the scope of this paper. It will require another controller at the transport layer in order to do some sort of shaping of exogenous arrivals.

### 2.2.5.6 Practical Variable-V algorithm: Evaluation Results

We can summarize the effect of the proposed controller with the previous described scenario (section 2.2.5.5). Precisely, there are four events E1, E2, and E3, that results in a different offered load in the all-

wireless NoF under simulation. This events occur at instant $t$=5s (E0), instant $t$=20s (E1), instant $t$=60s (E2), and instant 80(E3). E0 denotes the injection to the network of flow F1 from node R to node D. E1 denotes the injection of the second flow F2 from node C to node D. E3 denotes the stop of flow F2. E4 denotes the stop of F1.

E0: As node R launches F1at instant t=5s, node C starts experiencing queue drops (see Figure 2.29 and Figure 2.31) for high $V$ parameters (i.e., $V$=200, $V$=150, and $V$=100).



**Figure 2.28: V value over time in node R**



**Figure 2.29: Queue Length Evolution in node C (closer neighbour to the destination D)**



**Figure 2.30: Queue length evolution in node F (farther neighbour from the destination D)**

However, the variable-V controller is able to react to this change in the offered load in the network by decreasing the $V$ variable in node R (see Figure 2.28 which describes $V$ parameter evolution in node R), as it detects an abrupt change in the data queue backlog of node C. Therefore, node R is able to its actual greedy behavior emphasizing the penalty function. The $V$ parameter in node C enters in a transient stage in which, as a result of the previous change, decreases the $V$ parameter. However, note that afterwards there is a slight increase due to the fact that the node detects that it can be greedier given that the queue length of node C has decreased (see Figure 2.29). This queue has decreased due to increase in the degree of load balancing in node R, which starts to send packets to node F. Therefore, node C tries to reach its optimal value increasing its $V$ parameter as much as possible so that there are no queue drops at R.

E1: Node C launches F1: At instant t=20s there is another abrupt change in the network the one caused by event E2. In this case, the variable-V controller is able to react under these circumstances and decrease its V variable. This is shown by the behavior of node C which can be observed in Figure 2.29. As node C is announcing a full queue length neighbors react operating with the original backpressure algorithm just taking decisions based on their data queue lengths (see Figure 2.28). In this way, the zone affected by the network maximizes the load balancing in order to avoid data queue drops. Packets are load balanced until they found a less congested network and then can be routed towards the destination.

E2: Node C stops F2: At instant 60s there is a change in the network the one caused by E2. In this case, the variable-V controller is able to react under these circumstances, and augment its V variable. As node C has stopped F1 it starts announcing a not full queue length. And neighbors react operating in a zone in which the greedy penalty function can take importance (see Figure 2.28).

E3: Finally, the fourth event consisted on stopping flow F1 at instant 80. Consequently, in terms of the V parameter all the nodes affected by the set of events in the network return to its previous initial state which corresponds to the V = 200 (see Figure 2.28).

On the one hand, we can see how the proposed variable-V controller reacts under the variance in the traffic load occurring at instant t=20s (see Figure 2.28), decreasing the V parameter in node C. In contrast, as can be shown in the graph describing the data queue drops in Figure 2.31), the queue length evolution with fixed and high V values grows until the queue length in the neighbors is exceeded. In the case of fixed and low V (i.e., V = 0, and V=50), nodes maintain a lowest queues, however, they do neither get high throughput and low end-to-end delays. In contrast, the proposed variable-V controller is able to optimize both network performance metrics while maintaining acceptable data queue lengths (in this case data queue lengths are under the size limit of 200 packets during the whole simulation time).

On the other hand, the variable-V controller attains similar queue drops to that obtained by the legacy backpressure algorithm (see Figure 2.31). Specifically, Figure 2.31 shows the queue drop evolution through time in one of the nodes in the network suffering more traffic load (i.e., node C in the simulation). We can observe how the variable-V controller do not experience queue drops, which is also the case of the legacy backpressure algorithm (i.e., V=0). In contrast, for fixed-V policies the number of queue drops increases, given that the emphasis in the penalty function giving "direction" denoted by the magnitude of the V parameter causes higher data queue lengths and so queue drops.



**Figure 2.31: Queue Drops Evolution in node C through time**

With regards to the usual network performance metrics for the all-wireless NoF we have studied the evolution of throughput and end-to-end network delay of both flows (see Figure 2.32, and Figure 2.33). Interestingly, the proposed variable-V policy outperforms in this case all the fixed-V policies but V=50 in terms of throughput. On the other hand, we observe how it experiences lower delays compared with more direct strategies including V=50. The high queuing delays that high fixed-V policies experience compared with variable-V policies are caused by the fixed and high degree of direction of these policies. Note that this case is quite favorable for the V=0 case as both flows are directed towards the same destination. In this case, the delay distribution corresponds to packets that arrive to the intended destination. Additionally, they do not experience high queuing delays as V=0 minimizes the queuing delays when the destination is unique. However, note that not all the packets arrive to the destination due to the decreasing queue backlog gradient towards the destination the routing control policy has to generate. We can conclude that a variable V policy is able to achieve a good tradeoff between throughput and delay under varying offered load conditions.

Aggregated Throughput



**Figure 2.32: Achieved Throughput Evolution with time**



**Figure 2.33: Boxplots of End-to-end network delay with different routing policies**

## 2.2.6 Multi Local Femto Gateway Extension

In the solution proposed we assumed there is only one LFGW in the all-wireless NoF. However for high-scale deployments there might be several LFGWs. We propose to do per-packet load-balancing from packets of one flow. Roughly, the intuition behind the solution we propose follows: With low loads we assume all the HeNBs are directed to the predefined LFGW given by the LMM entity. If opportunistically, packets cross another LFGW (see Figure 2.34), this LFGW will also pull the data packet from the NoF.

Furthermore, if a packet which is directed through a LFGW (LFGW 1 in Figure 2.34) finds in its way a congested area it will choose randomly another LFGW (LFGW 3 in Figure 2.34). As a metric to avoid a congested area we use the data queue length of the neighbors of the current node in which the packet is located. Note that ping pongs may occur in the network if the data queue lengths highly vary in a short timescale. We are assuming average data queue lengths (i.e., the input rate vector) may not vary in such a short timescale, and hence ping pong do not occur. Additionally, we find it might be of interest to compare this policy with random LFGW at the first hop in the NoF.

**Opportunistic Packet Pulling (1)**



**Figure 2.34: Multi Local Femto Gateway**

## 2.2.7 Dead Ends

We have done a preliminary study in a sparse deployment of a NoF with potential dead ends. Two kind of holes can occur light holes (e.g., an unique HeNB find a dead-end), and big holes (e.g., a big slice of the NoF composed by a set of HeNBs is experiencing dead-ends).

### 2.2.7.1 Light Holes

To do this, we generated a hole in the NoF by deleting a femtocell of the grid NoF under evaluation. A flow was sent from femtocell 40 to femtocell 49. The reader should also take into account that a dead end can also be caused by a node generating traffic. We tagged with an "H" (hole) the extracted femtocell (i.e., femtocell 46 in Figure 2.35) from the NoF. In order to stress the importance of the dead-end, we illustrate a case in which the $V$ parameter is fixed to a high value, hence giving a high weight to quickly forwarding packets towards the destination. Therefore, the location of femtocell 46 to direct this flow is of primal importance given the weight chosen to quickly forward packets to the destination.

Figure 2.35 plots the number of packets transmitted per second during the window interval during which the flow is launched. As can be illustrated in Figure 2.35, the $V$ parameter is such that a single path is used to reach femtocell 49 ($V$=250). However, when the distributed routing algorithm detects the dead-end (i.e., no neighbour closer in terms of Euclidean distance to the destination), it is able to circumvent the hole in the position of femtocell 46.

The preliminary solution is based on decreasing the relative importance in link weight calculations of the directivity to the destination. This is done operating using legacy backpressure algorithm (i.e., $V$=0) in those femtocells in which a dead-end is detected for a given destination. The algorithm is modified so that whenever a dead end is found in the network, the $V$ parameter is decreased. Specifically, the current algorithm divides by two the $V$ parameter.



**Figure 2.35: Heatmap illustrating the light hole problem**

This preliminary study suggests that a variable *V* algorithm could overcome dead ends in sparse femtocell deployments. However, a further study should be done. For instance, it might be of interest to evaluate how a V variable distributed routing protocol deals with bigger, heterogeneous, and concave holes in a sparse NoF. This is provided in next subsection.

Big Holes



**Figure 2.36: Hole Avoidance Scenario in a sparse NoF deployment**

Bigger Holes in the NoF may occur, especially in the case of NoF sparse deployments. This, in principle, should not happen in a NoF deployed by a MNO. However it might be of interest to study in order to increase the degree of self-organization of the NoF.

For the solution, we assume there are multiple LFGWs evenly distributed in the NoF and connected by wired (e.g., fiber) infrastructure. Therefore, hole avoidance and multi LFGW solutions are tightly coupled.

Whenever, the distributed routing algorithm detects a big hole in the NoF, it will redirect data packets to another LFGW (e.g., LFGW 1 in ). The strategy to select this other LFGW may strongly vary. It could be the LFGW closer to the current location, or even a LFGW chosen randomly. Thus, the routing protocol will redirect packets to another LFGW (see Figure 2.36). To reach packet to this LFGW, the solution described in Figure 2.36 is used. Packets will eventually reach a LFGW that exploits its wired infrastructure to lead packets to the closer LFGW to the destination. Once they reach the closer LFGW they will be directed towards the destination.

## 2.2.8  Conclusions

We have studied the feasibility of a dynamic backpressure routing combined with geographic coordinates for an all-wireless NoF dealing with all possible traffic patterns (uplink, downlink and local). We have shown that the distributed routing algorithm is able to direct packets to the intended destination no matter the traffic pattern directionality. To do so, we evaluated the NoF under 24 different combinations of source destination HeNBs in a regular 10x10 grid topology.

Specifically, we have evaluated how various network performance metrics (network and flow throughput in sections 2.1.12.2.4.1 and 2.2.4.2, delay in 2.2.4.3, and fairness in section 2.2.4.4) are essentially affected by the value of the weight of the penalty function and the size of the queue at nodes, which constrains the results obtained with the theoretical framework. In this direction, we also study how the location of source-destination pairs and the tradeoff between objective function and backlog are affected by these parameters.

On the other hand, in the solution proposed, we have seen that an important aspect is the appropriate choice of the control parameter *V*. Specifically, the variability of the form of the input rate matrix injected to the NoF may cause differences in terms of achieved performance of the network performance metrics.

Specifically, we have observed that the traffic served could highly vary in terms of network and per flow throughput, delay, and jitter depending on the $V$ parameter.

Nevertheless, in a real NoF, the input rate matrix would be variable in time. Therefore, this suggests that it may be convenient to have a variable $V$ value for each HeNB. We have proposed and provided an initial evaluation of a distributed variable-$V$ algorithm, which auto-configures the more appropriate $V$ value in each HeNB depending on the traffic conditions. Initial results show that a variable-$V$ can outperform fixed-$V$ policies in terms of throughput and end-to-end delay.

Finally, we described a solution for a multi LFGW NoF in order to increase the offered load the NoF can satisfy. And, we have illustrated how the proposed distributed routing strategy jointly with assisted with the presence of multiple LFGWs connected by a wired infrastructure can assist in the dead-end problem of geographic routing and the directionality in backpressure to support all traffic patterns with a variable $V$ algorithm.

Therefore, the combination of both routing strategies (backpressure and geographic) seems to be a promising approach to cover the needs of a large-scale all-wireless NoF with support for all traffic patterns and capable of circumventing dead ends)

## 2.3 Voice Call Capacity Analysis of Long Range WiFi as a Femto Backhaul Solution

### 2.3.1 Motivation

The work in Section 2.2 addresses wireless backhauling between femtocells in a local network of femtocells. A related but different problem, which is studied in the following, is that of long-range wireless backhauling for outdoor femtocells.

Due to the high cost of deploying macro cells, both in terms of base station cost and backhauling, there are still areas that lack cellular coverage. The primary reason for this is that mobile network operators tend to plan network deployments in order to maximise population coverage and maximise Average Revenue per User (ARPU). Traditionally, mobile network operators rely on E1/T1 copper, optical fibre, and microwave links to backhaul their networks. However, these solutions can be extremely expensive to deploy and constitute a large portion of the operators Operational Expenditure (OPEX). Therefore, it is often not economically viable to deploy macro cells to cover areas with relatively low population densities or low ARPU; this is particularly true in remote, rural or third world areas. This has motivated many operators to begin looking for alternative lower cost backhauling solutions such as WiFi [43].

Although originally designed as a short range best effort wireless technology, recent studies and deployments have shown the ability of IEEE 802.11 (WiFi) based technologies to achieve much longer distances than was originally envisaged [44]. This means that using long range bi-directional Point-to-Point WiFi as a backhaul solution has the potential to be a much lower cost alternative than traditional microwave backhaul links, both in terms of equipment costs and licensing. Furthermore, many areas, particularly in developing nations, which have poor connectivity have already deployed long range WiFi solutions to provide internet and telecommunications connectivity.

For example, in 180 the authors outline their experiences in deploying long range WiFi in rural areas of India which had little or no penetration of cellular technologies. It is worth noting that in these deployments, due to the low penetration of telecommunications infrastructure, the most popular service being used was Voice over Internet Protocol (VoIP) and in 180 by the same authors they outlined the socio-economic benefits of this service. Another study by the ITU-T [48] detailed a number of areas in the Dominican Republic which are successfully using long range WiFi links to provide internet access. In [49] the authors presented the successful deployment of a 279 km link in Venezuela, and a permanent 133 km test network in northern Italy that is used for ongoing research; although this work was performed for experimentation purposes it shows the significant potential and feasibility of such deployments.

Apart from backhauling the other major factor limiting the deployment of cellular services to these areas is the Capital Expenditure (CAPEX) and OPEX of deploying macro base stations. For example, the power consumption of macro cells in certain developing countries accounts for almost 2/3 of the OPEX. This is primarily due to the power required for air conditioning.

Femtocells are small, low cost and low power base stations which utilise cellular technologies to deliver operator services to users in home and office environments. A femtocell is similar in appearance to a WiFi access point and connects to a consumer broadband connection such as Digital Subscriber Line (DSL) or cable over which it connects via a secure tunnel to the mobile operator's core network. Unlike macro base stations, femtocells do not have expensive deployment costs and due to their automated self-

configuration, interference mitigation mechanisms and relatively low power require minimal radio planning. Moreover, as they are a consumer device they do not need the expensive air conditioning systems required by normal base stations which have a significant impact on the operators OPEX. Although they are normally designed for use in homes and offices, there is now significant interest from a number of operators and manufacturers to develop and deploy outdoor femtocells which can support a larger number of calls and have similar range than a typical femtocell. The main goal of these systems is to provide a low cost and easily deployable solution to increase capacity in metropolitan areas; nevertheless it is obvious that they can also be used/adapted to provide coverage in rural and remote areas.

It is therefore clear that integrating long range WiFi backhauls or indeed existing rural broadband networks that utilise long range Wireless Fidelity (WiFi) links with the low cost of Femto hardware can greatly reduce deployment costs and therefore make it economically viable to deploy cellular coverage in areas in which it has previously not been feasible. Indeed due to the lack of fixed line infrastructure in the developing world, wireless coverage is the often only solution for providing cost effective services. It is also worth noting that due to the rapid growth in mobile phone penetration in developing countries, the demand for mobile services in such areas will continue to grow; according to a United Nations (UN) report [47] mobile phone penetration in developing nations had surpassed 57% by the end of 2009 and is continuing to grow.

In the long range WiFi deployments described earlier, solutions have been developed for many of the challenges that are encountered when using WiFi over long distances. There are however further questions that must be answered before WiFi can be used as a femtocell backhaul solution to deploy cellular services. The WiFi Media Access Control (MAC) layer was designed for short distances and therefore many of the timing values used are not sufficient for longer distances. Moreover, WiFi is inherently inefficient for transporting small payload frames such as those produced by voice calls This is further exacerbated by the large amount of overhead introduced by the Iu for Home NodeB (Iuh) protocol used to transport femto traffic to the operator core network. It is therefore not clear what capacity and call quality could be expected from such deployments. Another important point is the delay /jitter introduced in the critical synchronization traffic (heartbeat packets) with the operator's time/frequency server, located in the core network.

To date no study has been performed to analyse the capacity of WiFi as a backhaul solution for femtocell deployments. This paper performs both a numerical and simulation based analysis of the capacity of heterogeneous data rate WiFi links to act as a backhaul for multiple femtocells. Simulation models were developed in the NS-3 simulator [51] for the Adaptive Multi-Rate (AMR) voice codec incorporating Quality of Service (QoS) based codec adaptation, the Real-Time Transport Protocol (RTP), and the femto Iuh interface including GPRS Tunnelling Protocol User plane (GTP-U). Modifications to the existing simulator WiFi model were made to accurately model long range WiFi links and to take into account many common mistakes and assumptions made in previous WiFi capacity studies. The results show the capacity of long range WiFi links to backhaul femtocell deployments, specifically we show the number of simultaneous and high quality AMR voice calls that can be transmitted over the backhaul WiFi link considering both Circuit Switched (CS) and Packet Switched (PS) operational modes.

### 2.3.2  AMR Codec

Originally developed by the European Telecommunications Standards Institute (ETSI), AMR [52] is a robust narrowband voice codec designed for use in cellular systems. It is the mandatory codec defined for 3rd Generation Partnership Project (3GPP) systems including Universal Mobile Telecommunications System (UMTS) femto cells, which is the system considered in this paper.

The codec has eight source encoding rates which range from 4.75kb/s to 12.2kbps for voice payloads, a sampling rate of 8 kHz and a static frame size of 20ms is used for all rates. The codec is dynamic and based on the network conditions currently being experienced, can modify the source coding rate in an effort to adapt to channel conditions, reduce network load and mitigate any significant quality degradation. The codec also utilises Discontinuous Transmission (DTX) and Voice Activity Detection (VAD) to minimise bandwidth utilisation during silent periods. Specifically, when no user speech is detected the codec simply transmits silent packets containing comfort noise at a rate of 1.8kb/s. It should be noted that AMR achieves almost the same voice quality as the commonly used G.711(64kb/s) codec but with significantly lower bit rates.

Although originally designed for CS systems the codec can also be used for VoIP. In this case each AMR frame is encapsulated into an RTP packet and transmitted over UDP; in fact this is exactly what is done by femtocells using Internet Protocol (IP) backhauls.

**2.3.2.1 Implementation of AMR in NS-3**

Prior to this work there was no implementation of either the AMR codec or RTP available for NS-3. In order to achieve an extremely accurate simulation of femto voice traffic over WiFi, it was essential to develop these models within the simulator.

The RTP header was implemented in NS-3 according to the standard [53] and supporting all payload types defined by Internet Assigned Numbers Authority (IANA); this includes the dynamic payload type required by AMR. Validation of this implementation was done using the Wireshark network protocol analyser [21]. Specifically, the Packet Capture (PCAP) files generated from the simulator were analysed using the RTP packet dissectors available in Wireshark, each RTP header was successfully decoded and each element in the header was valid.

The AMR application was developed using the information available in [52] and [55]. It supports all eight codec modes and includes codec adaptation, DTX and voice activity generation functionality. The first component implemented was the packet structure of AMR frames. Each AMR frame consists of three elements, a header, a Table of Contents (ToC) and the speech data. The header is used to provide in-band signalling between both endpoints, it contains a Codec Mode Request (CMR) parameter which allows the receiver to request a particular source coding mode from the alternate endpoint. This request is in the form of an integer value in the range 0-8 where each value represents the requested source bit rate as shown in Table 2-2. The ToC is used to provide information about the contents of the payload, it contains a flag indicating if the payload is the last frame in a speech burst, a quality flag that can be used to indicate if the payload is damaged and a frame type parameter indicating whether the payload is AMR, AMR Wide Band (AMR-WB) or a silent packet.

A bit level accurate implementation of the AMR header and ToC was developed, however to simplify the implementation and to minimise simulation time, each payload was filled with dummy data. Although the application could be relatively easily modified to use real AMR voice payloads read from AMR encoded files, this would have no impact on the results presented in this work. The developed frame structure used octet-aligned mode as opposed to bandwidth-efficient mode, meaning that the last octet is padded with zeroes. This was done based on the observation that real handsets appear to be mainly operating in octet aligned mode. The packet structure and behaviour of the developed applications was validated by comparing the PCAP traces produced by the simulator with those produced by both a Sony Ericsson Elm and a Samsung Galaxy S running android, both operating in a real UMTS femtocell testbed.



**Figure 2.37: Flowchart of Codec Adaptation Decision Process**

**Figure 2.38: Example of Full Duplex Voice Activity generated by the AMR application**

### 2.3.2.2 E-Model based Codec Adaptation

A key feature of the AMR codec is the ability to seamlessly switch source encoding during ongoing voice calls. Each terminal continually monitors the downlink radio link quality between itself and the base station, this QI is then compared to a set of pre-defined thresholds to decide on the optimal source encoding rate that should be used [22]. The terminal then sets the CMR field of outgoing voice frames instructing the alternate endpoint to modify the source encoding rate.

As AMR was designed for use only in cellular networks the original adaptation mechanism only considered the radio conditions and assumed ideal conditions in the backhaul. However, new cellular technologies such as femtocells utilise IP backhauls which can become congested and degrade the quality of any ongoing AMR voice calls. To resolve this issue the AMR adaptation algorithm was modified to incorporate Explicit Congestion Notification (ECN) capable adaptation. During periods of congestion, as AMR voice traverses the IP network, ECN is used to mark the IP header to signal impending congestion. When a terminal receives a congestion notification it is recommended that the terminal "backs off" by reducing the encoding data rate in an effort to make network resources available and mitigate congestion. This approach works well in networks that support ECN such as enterprise networks or within the cellular operators private IP network. However, if all network segments do not support ECN, as would be the case in many home/office femtocell deployments, the terminal has no way of detecting backhaul induced call quality degradation or of determining the ideal encoding rate.

A number of other papers such as [22] and [23] developed algorithms for AMR voice codec adaptation, however these primarily focus on adaptation for improved resilience to interference and noise on the radio link and do not consider quality degradation in an IP backhaul link. In [24] the authors describe the development of an IMS testbed which utilised terminals with voice adaptation; although this work demonstrated the importance of adaptation the results only considered two AMR encoding rates, AMR12.2 and AMR4.75, with simple threshold based adaptation between both encoding rates.

In [25] the authors developed an algorithm for optimal source and channel coding selection based on the Mean Opinion Score (MOS) computed using the E-Model. Although this work provides an excellent overview of the improvements in call quality that can be attained through adaptation, it does not describe any specific implementation or describe how all of the required parameters can be obtained. Also, unlike the work presented here, the paper does not present results which demonstrate the proposed algorithm operating in a real or simulated environment.

In this work a terminal centric algorithm for call quality based AMR adaptation is proposed. Specifically, a modified variant of the E-Model is used to compute the call quality of the ongoing AMR call. Based on this the algorithm determines and continually modifies the requested source encoding rate to best match the current network conditions resulting in higher overall voice quality. Unlike many previous works, this paper provides a detailed description of utilising the E-Model in real time for AMR and the implementation of an AMR voice application in the NS-3 simulator utilising the E-Model based adaptation mechanism.

Source rate adaptation is normally done based on a combination of the radio channel conditions. However, in this work we developed a QoS based adaptation mechanism. The QoS metric used is the MOS and is computed using a real time implementation of the E-Model; a detailed description of the E-Model and using it in real time to compute the MOS for AMR is provided later. Furthermore, a more detailed overview of the E-Model based adaptation is provided in Section 2.3.2.2.

Figure 2.37 shows a flowchart of the adaptation algorithm decision process. Each VoIP application continually monitors its downlink QoS in real time based on delay, jitter, loss and the AMR codec mode. When the MOS score falls below a predefined threshold the application sends a CMR request using the AMR header of outgoing packets to the source application, on reception of this the source application will change the source coding rate to the requested mode. Correspondingly, if the computed MOS value is above a predefined threshold then the receiving application sends a CMR request to increase the coding rate to the next highest rate.

In order to prevent rapid and continual changes in the codec modes due to the small network fluctuations, a time constraint is placed on the period between codec changes; in the results presented in this paper a time constraint of 1 second is used. This constraint is not used when changing from a silent period to a voice active period. The utilised predefined MOS thresholds are shown in Table 2-2 and were determined by calculating the MOS value for each codec mode based on 0% packet loss and a 150ms end-to-end delay.

**Table 2-2: Codec Adaptation MOS Thresholds**

| Mode | Codec | MOS |
|------|-------|-----|
| 0 | AMR 4.75 | 2.46357 |
| 1 | AMR 5.15 | 2.56896 |
| 2 | AMR 5.9 | 2.77826 |
| 3 | AMR 6.7 | 2.83263 |
| 4 | AMR 7.4 | 3.03283 |
| 5 | AMR 7.95 | 3.18159 |
| 6 | AMR 10.2 | 3.46013 |
| 7 | AMR 12.2 | 3.62938 |
| 8 | AMR SID | None |

**Table 2-3: Call Rating Thresholds**

| R Value (Lower Limit) | MOS (Lower Limit) | User Perception |
|-----------------------|-------------------|-----------------|
| 90 | 4.34 | Very Satisfied |
| 80 | 4.03 | Satisfied |
| 70 | 3.6 | Some Users Dissatisfied |
| 60 | 3.1 | Many Users Dissatisfied |
| 50 | 2.58 | Nearly All Users Dissatisfied |

### 2.3.2.3   E-Model based Codec Adaptation Results

A number of simulations were performed in NS-3, however due to space constraints only a subset of the performed simulations can be presented. The presented results demonstrate the improved call quality achieved by the proposed adaptation mechanism and the ability to dynamically adapt the source encoding rates to match the current network conditions. Each set of presented results shows the call quality of a single AMR call, however in most of the simulations there were also a number of other 'background' AMR calls. Furthermore, in order to make the figures easier to interpret the DTX functionality was disabled on the call for which the results are shown; however any background calls that were present utilised DTX functionality. This has no impact on the presented results but simply means that the figure showing the AMR mode does not fluctuate between voice active modes and the silent mode (AMR8).



**Figure 2.39: Single AMR Call through 63kbps Channel with Proposed Adaptation Mechanism**

Figure 2.39 shows the results of a single call being transmitted through a band limited link. Specifically, a single AMR voice call was placed through a 63kbps CSMA link between two endpoints which is insufficient to support the call at the maximum AMR rate. This simple scenario shows the ability of the proposed algorithm to automatically adapt the source encoding rate to match the available network resources. It should be noted that without QoS based adaptation the call would simply saturate the link and suffer detrimental call quality degradation.



**Figure 2.40: Call Quality over 2Mbps Link with Random Packet Loss**

Figure 2.40 demonstrates the ability of the mechanism to increase resilience during periods of packet loss. In this simulation a 2Mbps link was used with the random packet loss model from NS-3. As can be seen, as the packet loss varies the algorithm continually and dynamically modifies the source encoding rate to reduce network load in an attempt to mitigate the loss and maintain high call quality.



**Figure 2.41: Call Quality over 2Mbps Link with Background Voice Calls and No Adaptation**



**Figure 2.42: Call Quality over 2Mbps Link with Background Voice Calls and Adaptation**

Another set of simulations were performed to demonstrate the improved call quality achieved when a large number of simultaneous voice calls traverse the same 2Mbps link. Each of these simulations was run for a period of 400s, a single voice call was started at 5s with an additional background voice call being started every 5s up to a maximum number of 32 background calls. After 200s the background voice calls are torn down, one every 5s in the same order as they were started. Figure 2.41 and Figure 2.42 show the results with and without adaptation respectively.

As can be seen the link begins to become congested after approximately 160s at which point the call quality begins to degrade. In the case were no adaptation is performed Figure 2.41 the call quality of the

ongoing voice calls drops to a very low MOS value of approximately 1.5. However, in the case were all calls utilise the proposed adaptation mechanism as shown in Figure 2.42, much higher call quality was maintained due to each endpoint dynamically adapting encoding rates to match network conditions. It can also be observed that with the proposed adaptation the call quality recovers much faster as calls begin to end and link capacity becomes available.

### 2.3.2.4  Discontinuous Transmission

The AMR application uses a simple two state conversational speech model with a 50% activity ratio per user. This is used to emulate DTX, silent suppression and VAD functionality. The model assumes that while one user is speaking, the other user is silent and producing only comfort noise packets. The algorithm randomly selects an integer value between 1 and 6 seconds to determine the duration of each voice spurt. Figure 2.38 shows an example output of the voice call speech pattern generated by the algorithm.

### 2.3.2.5  Voice Quality Assessment

A standard metric for assessing a person's perception of voice call quality is the MOS. Traditionally MOS testing has been done by asking large groups of people to listen to various voice calls and to score each call between 1 and 5. Although this solution provides very accurate results it is often not feasible to assemble large groups of people to perform such subjective tests. Furthermore, in order to be able to quickly assess the quality of ongoing voice calls or indeed to use the call quality for codec adaptation, as is the case in this work, an automated approach is required.

The E-Model algorithm is an International Telecommunications Union Technical standards (ITU-T) standardised computational model for subjective call quality assessment [56]. It is widely accepted as an accurate tool for measuring call quality and has been utilised in a number of previous works including [57][58].The E- Model operates under the assumption that perceived quality impairments are additive. By combining both codec and network impairments it produces a scalar rating of voice call quality called the R rating.

The R rating is computed as:

$$R = Ro - I s - I d - I e + A$$

where Ro is the Basic signal-to-noise ratio, Is represents impairments simultaneous to voice encoding, Id is impairments due to network transmission, Ie represents the effects of equipment such as codec distortion and A is the advantage factor.

The parameters Ro and Is are associated with the voice signal and therefore are not affected by transmission over the network. The advantage factor A is used to offset the reduced quality users may be willing to accept in certain circumstances such as in a mobile environment but in this work A is set to 0 such that a valid comparison between the presented results and other existing work can be made.

Based on the previous assumptions considering only the variable parameters the E-Model algorithm can be reduced to the following simplified equation:

$$R = Ro - I d - I e \qquad\qquad (2)$$

where Ro has a default value of 93.2. The only variable parameters are those affected by network transmission and equipment, namely Id and Ie. Id is affected by one-way delay and jitter while Ie is codec and loss dependent.

### 2.3.2.6  The Delay Impairment Parameter Id

One-way delay is defined as the time between the utterance at the mouth of the speaker to the time the signal arrives at the earpiece of the receiver, known as mouth-to-ear delay. Low delay values below 150ms have very little impact on call quality or interactivity. As delay values continue to increase above 150ms call quality begins to degrade, with delays above 400ms making a duplex call extremely difficult due to loss of interactivity. The one-way delay measurement used in the E-Model is made up of the following five parameters:

Send/Receive Delay: This accounts for send and receive side medium access delay

Propagation Delay: The end-to-end transit time across the network including time to traverse firewalls, routers etc.

Decoding/Jitter Buffer Delay: The delay introduced at the receiver due to decoding and buffering.

The Id parameter is the sum of all the individual delay parameters. Two crucial elements needed to obtain an accurate calculation of call quality is good estimates of one-way delay and jitter. The most common

approach is to use the delay calculation obtained using the Real Time Control Protocol (RTCP), while jitter is calculated using the E-Model recommended RTP jitter algorithm [35].

### 2.3.2.7  The Equipment Impairment Factor Ie

The Equipment Impairment Factor Ie takes into account how the used voice codec affects the perceived call quality; this includes the effect that varying loss levels have on the quality. For example, codecs which have high levels of inter-packet dependency are more significantly impacted in the presence of loss. The ITU-T have defined Ie values for the most commonly used non-adaptive codecs such as G.711 (Pulse Code Modulation (PCM)) and G.729, and for only the 12.2kbps encoding rate of AMR. Unfortunately no such values have been provided for other AMR codec modes.

In [60] the authors developed a combined E-Model and Perceptual Evaluation of Speech Quality (PESQ) method to compute the quality of AMR calls. A set of real AMR call measurements were used to compute the PESQ scores for varying loss rates. For each measurement the authors assumed zero delay and from this derived Ie values using the E-Model. These Ie values were used with the normally computed Id parameter to calculate MOS values giving a combined E-Model/PESQ algorithm. Although this concept is novel, it only considered the 12.2kbps mode of AMR.

In [61], the authors propose the use of a differential MOS approach to compute the Ie value. Similar to the previous paper a number of tests taken on a real Wideband CDMA (W-CDMA) network were used to compute the quality degradation of various loss rates with respect to AMR operating at 12.2kpbs. The problem with this work is that the authors do not provide the specific mappings between the loss rate and the Ie factor. In [62] the authors also used a combined E-Model/PESQ approach similar to the previous papers. The authors utilised the Ie algorithm defined in the E-Model combined with the methodology from another ITU-T recommendation [63] to determine Ie values for each AMR mode and for varying loss rates.

In this paper the same Ie equation is used and is defined as:

$$I_{e-eff} = I_e + (95 - I_e) \cdot \frac{100 \cdot P_{pl}}{\frac{100 \cdot P_{pl}}{BurstR} + B_{pl}}$$

where Ppl is the packet loss, Bpl is a codec specific loss robustness factor and BurstR is the quotient of the average burst length and is dependent upon the theoretical burst length under random loss conditions. In this paper only a single frame is sent per packet giving independent losses and hence we set BurstR = 1. Unfortunately for AMR Bpl is only defined for the 12.2kbps mode and there are currently no values defined for the other modes. Bpl is dependent upon the interpacket dependencies and the packet loss concealment scheme, and so since all AMR codec modes have a similar structure the authors recommend utilising the value defined for the 12.2kbps mode of Bpl = 10 for all codec modes. Utilising Equation (3), Figure 2.43 shows the impact of increasing loss rates on the call quality achieved by the AMR codec.

### 2.3.2.8  R to MoS Conversion

As stated previously, the output of the E-Model is a scalar rating of voice call quality which ranges from 0 100. However, this is not a commonly used metric for voice quality assessment. It is therefore possible to convert the R value to the more commonly used metric, the MOS value. This allows the results to be directly compared with other existing work and the outputs produced by most call quality assessment methodologies. Figure 2.44 shows the mappings between the R rating and MOS values, and Table 2-3 shows the commonly accepted user satisfaction thresholds as defined in ITU-T G.107 [56].

### 2.3.3  Iuh for HNB over WiFi

The stage 2 UTRAN architecture for 3G Home NodeB is used as the reference architecture for this work; a simplified version if this architecture is shown in Figure 2.45 (this also corresponds to the simulation scenario used). It is worth noting that in 3GPP parlance a femtocell is referred to as a Home NodeB (HNB) and a femto gateway is called a Home NodeB Gateway (HNBGW).

**Figure 2.43: AMR MOS Values for Varying loss Rates.**



**Figure 2.44: Mapping of E-Model R value to MOS.**

**Figure 2.45: Reference Architecture.**



**Figure 2.46: Iuh CS/PS User Plane Protocol Stack.**

The Iuh protocol stack is based upon the Iu interface used in UMTS which is the interface between the Node B and the core network. The Iuh interface is however specifically designed for use with femto cells and provides tunnelling of user and control plane messages over IP networks to the operators core network. For example, the Iuh protocol simplifies the control plane by removing the Iu protocols required for connecting over high performance links and SS7 emulation. These are replaced with a more scalable

solution which will allow the existing network to support the large number of femto cells that is envisaged. Enhancements were also added to support the seamless authentication procedures required to make femto cells a plug and play device. These enhancements also allow the femto cell to dynamically join and leave the network as may be the case if a user regularly powers down the device. Further discussion of these changes is outside the scope of this paper and interested readers are referred to [64]. The work in this paper is focused on the user plane aspects of the Iuh interface protocol stack. This is what allows user plane data to be securely tunnelled from the femtocell over open access IP networks, such as a private internet connection, to the UMTS core network. Figure 2.46 shows the Iuh protocol stack transporting AMR voice over both the CS/PS modes. Currently, the majority of voice traffic being transported over a UMTS femtocell is CS traffic, however with the increasing popularity of VoIP clients being used on mobile handsets this may change in the near future or at least become a far more significant percentage of the overall voice traffic.

In the current architecture each femtocell has a single IP Security (IPSec) tunnel to the Femto Gateway (FGW) and all traffic is aggregated into this tunnel for transport to the core network via the FGW. In normal CS operation the user's handset sends the AMR payload to the femtocell which is then encapsulated in an RTP packet for transport using UDP through the IPSec tunnel over the IP backhaul. Essentially this converts a CS voice call into a VoIP call. In the other case a VoIP client application would be used on the mobile handset, this means that the VoIP application would produce the AMR payload and an RTP and UDP header encapsulated in an IP packet. From the perspective of the femto cell this would appear as normal PS data traffic and as such would be encapsulated into a GTP-U tunnel for transmission over the IP backhaul network. As can be seen from Figure 2.46 and Table 2-6 there is a significant amount of overhead required to transport each voice frame over the backhaul. In order to transmit a 12.2kbps AMR frame in CS mode there is an overhead of 134 bytes, this overhead increases to 162 bytes in PS mode due to the addition of the GTP-U and the RTP/UDP/IP overhead produced by the VoIP application.

### 2.3.3.1 Validation of Iuh Implementation

In order to validate the accuracy of the Iuh stack developed in NS-3, Wireshark was used to analyse the stack structure and verify the correct structure of each header; this was done for both CS and PS modes. The Wireshark captures were then compared with captures obtained from a real femtocell testbed and a real UMTS network. Apart from the physical layer, the capture was identical in terms of frame and header structures.

Figure 2.47a and Figure 2.47b shows a CS and PS capture, respectively, both obtained from the simulator. Each shows the AMR/RTP packet being transmitted via Iuh over WiFi from a femtocell. It should be noted that the IPSec tunnel is not shown in these examples. The reason for this is that although the Encapsulated Security Payload (ESP) header was implemented in the simulator, a setup protocol was not and as such wireshark could not correctly dissect the frame. However, for all simulations for which results are presented the IPSec tunnels were included.



(a) CS Stack                    (b) PS Stack

**Figure 2.47: Capture of Protocol Stacks Produced by Simulator.**

### 2.3.4 Long Range WiFi Considerations

The standard timing values of 802.11 are optimised for relatively short distances; however these become problematic when operating over longer distances. Due to the longer propagation times over longer distances some 802.11 timing values are too short and result in spurious timeouts and increased collisions.

In order to maximise throughput over longer distances the following parameters must be considered:

Ack Timeout: All 802.11 frames must be acknowledged and the Ack timeout is the duration a node must wait after transmission having not received an acknowledgement, before assuming that a frame has been lost and proceeding with retransmission. If this value is too low the sending node will attempt a retransmission before the frame has had time to be successfully acknowledged and may cause a collision

with the Ack frame. Conversely if the value is too high the sending node will wait too long after a packet has been lost before sending a retransmission resulting in decreased throughput. The optimal value for the Ack timeout is calculated as the sum of a Short Interframe Space (SIFS), the time taken to transmit the Ack frame and twice the propagation delay. In 802.11a the constant values used by many cards assumes a propagation time of only 1μs which leads to a maximum range of 300m over which high throughput can be maintained. In order to maximise the throughput of a link the propagation delay needs to be set to match the distance between the link endpoints. Considering the propagation speed of radio waves is 300m/μs this means that the optimal propagation delay needs to be increased by 1μs for every 300m increase in the link distance. For example, for the 5km link being assumed in this work would give a propagation delay of (5000m 300m/μs) = 16.66μs.

DIFS: This is the duration for which a station must sense the channel is idle before proceeding with a transmission. If the Distributed Interframe Space (DIFS) is set too low a node may sense the channel is idle while another node has already begun transmission thereby causing a collision. Therefore the DIFS time must be greater than the propagation time between the furthest nodes in the network or the distance between the two endpoints in the case of a point-to-point link. The default DIFS in 802.11a is (SIFS + (2 x Slot time)) = 34μs which gives a maximum distance of (300m/μs x 34μs) = 10.2km. As a 5km link is assumed for this work no modification is necessary. However, for links above 10.2km the DIFS needs to be increased accordingly.

Slot Time: The slot time is set such that a station wanting to transmit will have received any already transmitted frames prior to it beginning its own transmission, hence mitigating collisions. The default 802.11a slot time is 9μs giving a maximum distance between endpoints of (9μs x 300m/μs) = 2.7km. For any link above 2.7km the slot time must be increased, with the optimal slot time computed as (Link Distance / 300m/μs). In the 5km link being considered here, this gives a slot time of 16.6μs which is rounded up to 17μs.

### 2.3.5 Analytical Calculation of Capacity

In this section an analytical evaluation of the AMR voice call capacity of long range WiFi to support femtocell deployments is performed. The capacity is defined as the maximum number of full duplex AMR voice calls that can be supported. Assuming no collisions, the minimum time required to successfully transmit an AMR packet over Iuh can be calculated as:

$$T_{amr} = T_{Difs} + T_{Sifs} + T_v + T_{Iuh} + T_{Ack} + T_s \cdot \left(\frac{CW_{min}}{2}\right)$$

where Tv is the amount of time take to transmit the AMR voice payload (which varies depending on the codec mode), TIuh is the time taken to transmit all the required Iuh headers, TDifs is the duration of the DIFS, TSifs is the duration of the SIFS, TAck is the amount of time required to transmit ac acknowledgement frame, Ts is the slot time and CWmin is the minimum Contention Window (CW) size used. Table 2-4 shows all elements of Tv and TIuh and as an example is computed for a link rate of 6Mbps. In this example the voice frame being transmitted is AMR mode 7 corresponding to a rate of 12.2kbps.

In the voice activity model being assumed in this work, there will be an equal distribution of both silent and voice packets in each direction of the voice call. Based on this, the average packet transmission time, assuming no retransmissions, can be written as:

$$\bar{T}_p(0) = \frac{1}{2} \times \left(2\left(T_{Diffs} + T_{Sifs} + T_{Ack} + T_s \cdot \left(\frac{CW_{min}}{2}\right) + T_{Iuh}\right) + T_{v(7)} + T_{v(8)}\right)$$

where Tv(7) and Tv(8) are the times taken to transmit a 12.2kbps and a silent frame, respectively. Based on this, the upper bound on the number of calls that can be supported is given by:

$$N_{calls} = \frac{P}{2 \cdot \bar{T}_p(0)}$$

where P is the codec frame size which is 20ms for AMR. In order to perform a realistic capacity analysis, the collision probability and retransmission rate for this particular use case must be considered. Here we use the analytical approach developed in [32] to compute the total packet transmission time; by applying their approach to this work, the total transmission time assuming n retransmissions can be written as:

$$T_n = \bar{T}_p(n+1) + \{\sum_{k=1}^{n} min(\frac{2^n(CW_{min}+1)-1}{2}, \frac{CW_{max}}{2}\} \cdot T_s + n(Ack_{timeout} \cdot T_s)$$

Based on this, assuming only two nodes and operating at saturation, the average time to successfully transmit a packet can be computed as follows:

$$\bar{T}_p(n) = \frac{\sum_{n=0}^{R} T_n \cdot P_{col}^n(1-P_{col})}{\sum_{n=0}^{R} P_{col}^n(1-P_{col})}$$

Using $\bar{T}_p(n)$ as the average packet transmission time in equation (6) gives the total number of calls that can be supported, considering both collisions and retransmissions. The values in Table 2-4 and Table 2-5 were obtained utilising the developed analytical model. These tables show the number of full duplex CS AMR and VoIP AMR voice calls that can be supported over a 5km point-to-point WiFi link. The upper bound is computed assuming no collisions while the 'with collisions' column utilised equation (8) and assumes a WiFi retry limit of 4. In this table the simulation results are also presented, these will be explained in greater detail in a later section but as can be seen there is a high degree of correlation between the analytical and simulation results.

### 2.3.6 Simulation Methodology

The NS-3 simulator is a discrete-event based network simulator written in C++. It includes a number of powerful features such as producing accurate PCAP traces from simulated scenarios and the capability to integrate with real network deployments. For these reasons it is beginning to gain traction among the research community.

**Table 2-4: AMR over Iuh Voice Capacity for 5km 802.11a link (CS Mode)**

| Data Rate | Analytical (upper bound) | Analytical (with collisions) | Simulation |
|---|---|---|---|
| 6 Mbps | 25 | 23 | 22 |
| 9 Mbps | 30 | 28 | 27 |
| 12 Mbps | 33 | 31 | 31 |
| 18 Mbps | 37 | 34 | 34 |
| 24 Mbps | 39 | 36 | 36 |
| 36 Mbps | 41 | 39 | 40 |
| 48 Mbps | 43 | 40 | 41 |
| 54 Mbps | 43 | 40 | 42 |

**Table 2-5: AMR over Iuh Voice Capacity for 5km802.11a link (PS Mode)**

| Data Rate | Analytical upper bound) | Analytical (with collisions) | Simulation |
|---|---|---|---|
| 6 Mbps | 22 | 21 | 19 |
| 9 Mbps | 27 | 26 | 25 |
| 12 Mbps | 30 | 29 | 28 |
| 18 Mbps | 34 | 32 | 33 |
| 24 Mbps | 37 | 35 | 35 |
| 36 Mbps | 40 | 38 | 38 |
| 48 Mbps | 41 | 39 | 40 |
| 54 Mbps | 42 | 40 | 41 |

**Table 2-6: Calculation of Tv for an AMR(7) Frame over Iuh at 6Mbps.**

| Parameter | Bytes | Bits | Time ($\mu$S) (PS Mode) | Time ($\mu$S) (CS Mode) |
|---|---|---|---|---|
| AMR Header | 2 | 16 | 2.66 | 2.66 |
| AMR Payload(7) | 31 | 248 | 41.33 | 41.33 |
| RTP Header | 12 | 96 | 16 | 16 |
| UDP | 8 | 64 | 10.66 | - |
| VoIP Application IP | 20 | 160 | 26.66 | - |
| GTP-U Header | 8 | 64 | 10.66 | - |
| UDP | 8 | 64 | 10.66 | 10.66 |
| Remote IP | 20 | 160 | 26.66 | 26.66 |
| IPSEC ESP | 8 | 64 | 10.66 | 10.66 |
| Transport IP | 20 | 160 | 26.66 | 26.66 |
| LLC(SNAP) | 8 | 64 | 10.66 | 10.66 |
| MAC Header | 28 | 224 | 37.33 | 37.33 |
| PLCP Header | 6 | 48 | 8 | 8 |
| Preamble | 18 | 144 | 24 | 24 |
| $T_v$ | | | 262.66 | 214.66 |

It is however still at a relatively early stage and therefore the number of implemented models is limited. Hence, as described earlier a number of models were developed specifically for the work presented in this paper. In particular, an AMR application, RTP, GTP-U and the IPSec ESP; these were required in order to accurately simulate the transport of AMR voice calls over femtocells with a WiFi backhaul.

In terms of WiFi, NS-3 already has an existing model and in [33] a detailed validation of this WiFi model was performed. The authors compared results obtained from the simulator with those from a deployed WiFi testbed called the EXTREME testbed. This allowed the authors to identify a number of disparities between the results and to identify the causes of these disparities. As mentioned previously, in an analysis of the common causes of inaccuracies in simulating VoIP over WiFi was performed. Similar to the previous paper, the authors identified these causes based on comparisons between results obtained from a number of simulators versus those obtained from a real testbed.

The findings and recommendations of both of [66] and other parameters we identified are used to improve the accuracy of the simulations performed in this work. The list below details all factors that were considered for the simulations:

Preamble: The WiFi model in NS-3 uses the long preamble (144 bits) by default, however the majority of real WiFi hardware uses the short preamble (72 bits). The NS-3 WiFi code was patched to use the short preamble by default.

RTS/CTS: This is a virtual carrier sensing implementation used in 802.11 to mitigate collisions due to hidden nodes. However, it is not used in the majority of WiFi deployments and due to the point-to-point link being considered here, there will be no hidden nodes. The RTS/CTS mechanism was therefore disabled.

ACK Frame Data Rate: All 802.11 data frames must be positively acknowledged, however for sending this 14 byte frame no data rate was defined in the standard. Although hardware manufacturers use a number of different rates the parameter is usually a user definable variable in the driver of the WiFi card. As a static point-to-point link is being used in this work, the ACK frame data rate was set equal to the channel data rate. This maximises the link capacity by minimising the amount of channel time occupied by transmitting acknowledgements.

Contention Window: Each 802.11 node must select a random back off interval from 0 to the CW size, the node must wait this number of slots before attempting to transmit on the channel. The initial CW size is set at CWmin and increases exponentially after every unsuccessful transmission up until a value of CWmax is reached; after a successful transmission the CW size is reset to CWmin This mechanism introduces randomness in channel access which reduces the collision probability between stations. The

CW values in 802.11a are CWmin = 15 and CWmin = 1023 and these values were used in both the analytical and simulation study presented in this paper.

Link Adaptation: Link rate adaptation is performed based on changes in the Received Signal Strength (RSS) and Signal to Noise Ratio (SNR) of the wireless channel. Although it is used in most wireless cards it can be disabled. Currently it is not implemented in NS-3. As the endpoints of the links in our scenario are static no link adaptation should occur or should only occur for very short periods due to short term fading and it would be disabled in any real long range WiFi deployment. For this reason having no link adaptation in NS-3 will not impact the results.

Buffer Size: This is used by the 802.11 node to buffer all received packets while waiting on the medium to become idle. There is therefore a direct trade-off between packet delay and packet loss when deciding on the size of the buffer. In NS-3 this was set statically to a size of 400 packets. However, most real systems use a smaller buffer length, for example the MadWiFi driver uses a buffer size of 50 packets. For this reason the NS-3 buffer size was changed to the same as that used in MadWiFi.

Packet Generation Offset: One problem with many VoIP simulations is that the VoIP sources start at fixed time values, this leads to synchronisation between the sources and many packets being forwarded to the physical layer at the same instant, whereas the starting time of real calls is based on a random call arrival pattern. For this reason the start time was offset with a random variable between 0-20ms; this prevents any synchronisation between sources other than what would be expected in a real environment.

Ack Timeout, DIFS & Slot Time: These variables were set according to the discussion on long range WiFi considerations in a previous section.



(a) Delay for CS Voice (95th Percentile)

(b) Jitter for CS Voice (95th Percentile)

(c) Loss for CS Voice (95th Percentile)

(d) MOS for CS Voice (95th Percentile)

**Figure 2.48: CS Simulation Results.**

The simulation topology is the same as the reference architecture shown in Figure 2.45. Each User Equipment (UE) is connected directly to a femtocell; in this work it is assumed that there is ideal radio conditions between each the UE and the femtocell (HNB) and that there is minimal delay in the core network. This is done so that only the limitations of the Iuh interface over the WiFi link are reflected in the results. Each femtocell establishes an IPSec tunnel to the HNBGW. What is not shown in the architecture figure is that the AMR voice calls are terminated at nodes which are connected directly to the

HNBGW via low delay, high bandwidth point-to-point links. It is also assumed that each femtocell can only support 4 simultaneous calls and so there is always N/4 femtocells in the topology where N is the number of ongoing voice calls.

Each UE establishes a full duplex voice call with a node located outside of the femtocell network and produces the AMR payloads in the case of a CS call or the AMR/RTP/UDP/IP payload in the case of a PS VoIP call. On reception of these payloads the femtocells/HNBGW adds the appropriate headers considering either CS or PS traffic modes.

### 2.3.7 Simulation Results

In this section the simulation based voice call capacity analysis results are presented. The simulation environment was setup based on the simulation methodology described in the previous section. Each data point in the presented results represents the 95th percentile mean of 5 individual simulation runs, with each simulation run having a duration of 200 seconds. Results are presented for both CS calls and PS VoIP calls.

As mentioned previously, Table 2-4 and Table 2-5 show the upper bounds on the predicted number of high quality calls that can be supported at varying data rates, for both the analytical study and the simulation results. As can be observed there is a very high correlation between both sets of results.

#### 2.3.7.1 CS Results

Figure 2.48a,b,c and d, show the delay, jitter, loss and MOS results for an AMR CS call respectively. As can be observed in Figure 2.48a, regardless of the data rate there is only a marginal increase in the delay experienced as the number calls increases. However, when the upper threshold on the number of calls that can be supported is reached, there is a jump in both the packet loss and the end-to-end delay; this is primarily due to buffer overflow at the transmitting station. It should be noted that due to the relatively small buffer size of 50 used in this work, the end-to-end delay has very little impact on the QoS achieved as packets are dropped before any significant delay is introduced.



(a) Delay for PS Voice (95th Percentile)

(b) Jitter for PS Voice (95th Percentile)

(c) Loss for PS Voice (95th Percentile)

(d) MOS for PS Voice (95th Percentile)

**Figure 2.49: PS Simulation Results.**

Figure 2.48d shows the E-Model based MOS value experienced as the number of voice calls increases and for varying physical layer data rates. As can be seen a high MOS value of 4.2 is achieved until the upper threshold is reached at which point the MOS value for all calls drops rapidly, primarily due to the increase in packet loss.

### 2.3.7.2 PS Results

Figure 2.49a,b,c and d, show the delay, jitter, loss and MOS results for a PS VoIP call using the AMR codec, respectively. These results have the same characteristics as the CS results presented earlier with only minor differences in the upper bound on the number of calls that can be supported. The minor difference in the number of calls that can be supported is due to the larger overhead required by the VoIP calls.

## 2.4 Local Breakout for Networked Femtocells

### 2.4.1 Introduction

Local IP Access (LIPA) is a 3GPP architecture enhancement enabling femtocell users to access services in the IP network in which the femtocell is located directly, rather than by routing traffic via the mobile network and the public IP network back into the local network. This requires traffic to be "broken out" of the 3GPP domain at or near the femtocell. A related concept is the Selective IP Traffic Offload (SIPTO) that allows mobile operators to specify the subset of traffic flows between femtocell users and Internet services that is to be broken out at or near the femtocell in order to bypass (and thus offload) the mobile network. While the solution for LIPA and SIPTO proposed in the following builds upon the 3GPP's initial solution described in [7], it extends it in two aspects:

1) It uses a single, centralized breakout point on the LFGW, rather than one on each femtocell, which allows local mobility for breakout sessions, facilitates management and control of these sessions and provides consistent breakout with legacy femtocells.

2) It allows breakout sessions to be handed-out/in to/from the macro network as well as to (re-)establish LIPA sessions from the macro network.

The latter is important because in the current 3GPP solution, when a mobile user has an ongoing LIPA session, leaves the femtocell coverage and experiences a connection loss, there is no way to re-establish the session other than to defaulting to a different mechanism (e.g. the "Remote IP Access (RIPA)"), which results in an inconsistent and often frustrating user experience.

### 2.4.2 Work during Year 1

During the first project year, the state of the art and the solution space were analysed and, based on this, a concrete solution was designed. In the proposed solution, a new network element called Local Femtocell GateWay (LFGW) is introduced within the local femtocell network. The LFGW is transparently inserted into the S1 interface between the HeNBs and the EPC such that neither HeNBs nor EPC have to be changed (see Figure 2.50). The LFGW provides functionality for local mobility management and, to support traffic breakout as defined above, is extended with a local P-GW (≈L-GW) functionality and a new, optional interface (the "Remote Access Tunnel" S-RAT) between the L-GW and the P-GW of the EPC that enables the seamless service continuity in the macro-network and only carries data in case the user continues a session in the macro network.

### 2.4.3 Performance Model

In the following we conduct a performance evaluation in order to assess the impact of the LFGW and especially of localization of signalling and LIPA on the duration of signalling sequences.

We develop a flow-level model to reflect the interaction between elastic traffic like TCP-based file transfer or web-browsing, and streaming traffic generated by applications with QoS requirements like voice over IP or video streaming. The main difference is that elastic traffic tends to use as much link capacity as possible, while streaming traffic has more or less constant bitrate requirements (at least on average). Furthermore, elastic traffic is usually volume-based, i.e. a session is terminated after a certain data volume is completely transmitted, while streaming traffic is time-based, i.e. the session ends after a certain time span. For a comprehensive overview of related literature see [8].

**Figure 2.50: Logical architecture of the HeNB sub-system with LFGW.**

Effects on packet level are modelled with an M/D/1 approximation for streaming traffic and with the assumption of full buffer utilization for elastic traffic. The approach is insofar similar to [9], [10], but with the approach that TCP and streaming traffic characteristics are determined by the flow level dynamics, not the opposite.

### 2.4.3.1 Flow-Level Model of the Bottleneck Link

From a flow perspective, best-effort and QoS traffic arrives with arrival rates $\lambda_{BE}$ and $\lambda_{QoS}$ following a Poisson arrival process, i.e. with exponentially distributed inter-arrival times (see Figure 2.51). QoS traffic terminates in the EFN and in the outside network. Best-effort traffic is terminated in the EFN from both sides, such that in a scenario without LIPA the traffic from the local services towards the UEs has to traverse the backhaul link in uplink direction first. Therefore, it is assumed that the uplink backhaul link with a capacity of $C_{UL}$ bps constitutes the main bottleneck of the system. The traffic flow then traverses the ISP network to the operator PDN-GW, from which it is again routed via the backhaul downlink through the EFN to the requesting UE. For a scenario with LIPA support, we assume that best-effort traffic is offloaded at the LFGW breakout point (i.e. the L-GW) such that it terminates directly at the local server. However, we still assume that QoS traffic is terminated in the core network or beyond, i.e. it still traverses the backhaul link symmetrically.

QoS traffic is *time-based*, meaning that traffic is generated with an average data rate of $E[R_{QoS}]$ for a certain time until the voice or video call has ended. The call (or holding) time is assumed to be exponentially distributed with mean $E[T_{QoS}]$.

In contrast, best-effort traffic is characterized by that typically a certain data volume is downloaded (e.g. a complete website). It is therefore *volume-based* (see e.g. [11]) and the data volume is assumed to be exponentially distributed with mean $E[V_{BE}]$. Since best-effort traffic utilizes the link capacity as much as possible, the instantaneous flow throughput depends on the number of concurrently active flows, $n_{BE}$. Furthermore, the congestion avoidance mechanism of TCP leads to an implicit adaptation of best-effort flow rates to available bandwidth, i.e. UDP traffic tends to displace TCP traffic in bottleneck links [12], [13]. Therefore we model the available conditional per-flow sojourn time as

$$E[T_{BE}|(n_{QoS}, n_{BE})] = E[V_{BE}] \cdot \frac{C_{UL} - n_{QoS} \cdot E[R_{QoS}]}{n_{BE}}. \tag{1}$$

Thus, the backhaul link can be modelled as two-dimensional continuous-time Markov chain (CTMC) queuing system with states $(n_{QoS}|n_{BE})$, where the steady state distribution is determined by the offered load of best-effort and QoS flows. We assume a *fluid regime*, i.e. assume that elastic traffics adapts to streaming traffic instantaneously. For a comparison of different adaptation regimes see e.g. [14]. The infinite QBD generator transition rate matrix has then the following block-diagonal structure:

$$Q = \begin{pmatrix} B_0 & A & 0 & 0 & \cdots \\ D_1 & B_1 & A & 0 & \cdots \\ 0 & D_2 & B_2 & A & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \text{ with } A = \begin{pmatrix} \lambda_{BE} & 0 & \cdots \\ 0 & \lambda_{BE} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix},$$

and

$$B_k = \begin{pmatrix} d_{k,0,0} & \lambda_{QoS} & 0 & \cdots \\ \mu_{QoS}(k) & d_{k,1,1} & \lambda_{QoS} & \cdots \\ 0 & \mu_{QoS}(k) & d_{k,2,2} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \quad D_k = \begin{pmatrix} \mu_{BE}(k) & 0 & \cdots \\ 0 & \mu_{BE}(k) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}.$$



**Figure 2.51: Flow-level model of the bottleneck link.**

The diagonal elements are defined as $d_{k,i,j} = -\left(\lambda_{QoS} + \lambda_{BE} + \mu_{QoS}(k) + \mu_{BE}(k)\right)$. The departure rates are according to the QoS and elastic traffic model defined as $\mu_{QoS}(k) = k \cdot E[T_{QoS}]^{-1}$ and $\mu_{BE}(k) = \frac{E[V_{BE}]}{c_{ul} - k \cdot E[R_{QoS}]}$. Note that the latter follows as the reciprocal of the conditional best-effort sojourn time $E[T_{BE}]$ multiplied with the number of best-effort flows $n_{BE}$, which can be eliminated from the equation. Thus it depends on the number of QoS flows only.

The steady state distribution of the system is then the vector $\bar{\pi}$ which must satisfy

$$Q \cdot \bar{\pi} = 0. \tag{2}$$

To solve the matrix equation, the infinitesimal sub-matrices are bounded by a maximum number of admittable QoS flows $n_{QoS}^{max}$ and best-effort flows $n_{be}^{max}$ such that the number of states is kept in reasonable regions. Note that we are interested in low-utilization scenarios for the QoS flows such that blocking probabilities are neglected. However, even high utilization scenarios are reasonable by assuming that in this case users will cancel the ongoing call due to bad QoS, i.e. $n_{QoS}^{max} = \left\lfloor \frac{c}{R_{QoS}} \right\rfloor$.

### 2.4.3.2 Modelling Delay

The average queue length is calculated according to the current state of the system, which is the number of QoS and BE flows. For BE flows we assume that the aggregated TCP traffic tends to occupy the full buffer space, i.e. in this case the queue length is corresponds to the buffer space. QoS flows are assumed

to have Poisson packet arrivals with deterministic packet sizes. The aggregated QoS traffic can therefore be modelled as an M/D/1/K loss system where the average packet arrival rate $\lambda_{QoS,P}^{agg} = E[n_{QoS}] \cdot \lambda_{QoS,P}$ is the aggregate of the per-flow arrival rates, which are derived from the packet size and the QoS bandwidth as $\lambda_{QoS,P} = \frac{E[R_{QoS}]}{S_{QoS}}$. The offered load follows as $\rho_{QoS,P} = \lambda_{QoS,P}^{agg} \cdot E[T_{QoS,P}]$ with $E[T_{QoS,P}]$ as packet service time. As we are only interested in relatively low utilization scenarios for QoS traffic, the mean queue waiting times of the loss system is approximated with the Pollaczek-Khinchine formula for M/D/1 waiting systems as

$$E[D_{QoS}] = E[T_{QoS,P}] \cdot \left(1 + \frac{1}{2} \cdot \frac{\rho_{QoS,P}}{1 - \rho_{QoS,P}}\right). \tag{3}$$

Finally, the average queue delay follows from the conditional probabilities that the system is empty, and that only QoS flows or that best-effort flows are present as

$$E[D_{q,bh}] = \pi_{BE}(0) \cdot E[D_{QoS}] + (1 - \pi_{BE}(0)) \cdot T_{RTT}. \tag{4}$$

where $\pi_{BE}$ is the marginal distribution of the number of active best-effort flows.

For the calculation of the signalling delay, we assume that hand-over signalling messages arrive at the queue following a Poisson arrival process. Due to the PASTA property [17], incoming messages find on average a queue size corresponding to $E[D_{bh}]$. Furthermore, the average packet size of a signalling message is much smaller than the average flow volume, and the arrival rate is low, such that the effect of the messages on the queue size is negligible.

### 2.4.4 Performance Evaluation

In the following we show some numerical examples to get an impression of the system performance. The main focus is on the impact of best-effort flows with and without LIPA support on queuing delays, signalling completion times and backhaul link capacity. Message sequence completion times have a significant influence on the system performance, since resources in the operators network (e.g. due to state and context management for forwarding tunnels) and on the backhaul link are occupied. For the calculation of the duration of message sequences, a similar approach as in [18] is used, with the difference that we also consider the layer 2 hand-over time and bottleneck queuing delays as described in Section 2.4.3.2.

### 2.4.4.1 System Assumptions

The hand-over message sequence comprises several entities in different network domains. We differentiate between messages which are sent within the enterprise domain, between enterprise and mobile operator domain, within the MNO domain, and between UE and HeNB. The average transmission delay of a message may be composed of $D_E$ for transmissions within the enterprise domain, $D_{bh}$ for the backhaul link, $D_{MO}$ for transmissions in the mobile operator domain and $D_{AI}$ for transmissions on the air interface. For the latter, we assume that HARQ with 8 processes and an initial target block error probability of 10% is used, resulting to an approximate average one-way latency of 4.8ms, see e.g. [19]. Furthermore, we consider that non-contention based layer 2 handover (association to the target eNB) is performed, with an average delay $D_{L2HO}$ of 12ms [20], [21]. Table 2-7 summarizes the assumed latency values. Note that the value $D_{t,bh}$ refers to the transmission delay only, the total link delay follows as $D_{bh} = D_{q,bh} + D_{t,bh}$.

| Link | $T_E$ | $D_{t,bh}$ | $T_{MO}$ | $T_{AI}$ | $T_{L2HO}$ |
|---|---|---|---|---|---|
| Value (ms) | 2 | downlink: 0.41, uplink 0.82 | 10 | 4.8 | 12 |

**Table 2-7: Transmission delays for different transport links.**

For the backhaul link model, we assume an asymmetric DSL connection with 10Mbps in the downlink and 5Mbps in the uplink. QoS traffic has a deterministic packet size of 512 byte with an average bandwidth demand of 128 kbps and average flow duration of 180s. Best-effort flows have an average traffic volume of 2Mbyte. Signalling traffic arrives with exponentially distributed inter-arrival times with an average packet size of 512 byte, which is used to calculate the transmission delays in uplink and downlink.

In all scenarios, best-effort and QoS flows arrive with exponentially distributed inter-arrival times according the offered load $\rho \in [0.05, 0.4]$. In scenarios without LIPA support, the offered load is equally distributed between best-effort and QoS load:

$$\rho = \rho_{\mathrm{BE}} + \rho_{QoS} = \frac{\lambda_{\mathrm{BE}} \cdot E[V_{BE}]}{C_{\mathrm{ul}}} + \frac{\lambda_{\mathrm{BE}} \cdot E[T_{QoS}]}{n_{QoS}^{\max}}. \tag{5}$$

The arrival rates $\lambda_{\mathrm{BE}}$ and $\lambda_{\mathrm{QoS}}$ follow accordingly. In scenarios without LIPA, the QoS load is in the range QoS $\rho_{QoS} \in [0.05, 0.4]$ while the best-effort load is 0, if not stated otherwise in the text.

### 2.4.4.2 Queuing Delays



**Figure 2.52: Mean bottleneck queuing delay vs. offered load.**

**Figure 2.53: Impact of best-effort and QoS load on mean bottleneck queuing delay.**

Figure 2.52 illustrates the impact of the bottleneck utilization on the average queuing delay. The curve without LIPA support increases steeply and nearly linearly with the offered load due to the dominating influence of best-effort flows on the queue size. On the other hand, in a situation with local offloading of best-effort flows, the queuing delay for a comparable offered load of QoS traffic increases only slightly. The reason is that the mean waiting times in M/G/1 queuing systems depend on the coefficient of variation of the service times, which is in our case $C_v = 0$ due to the deterministic packet size.

In Figure 2.53 the impact of best-effort and QoS load on the bottleneck queuing delay in the non-LIPA case is shown. In the first case, only the best-effort load is increased, while for the second curve the QoS load is increased. In both cases the load of the other traffic share is kept constant at 0.15, such that the maximum total load reaches 0.55. An increasing best-effort load leads to a linear increase of the queuing delay. Increasing QoS load affects the queuing mainly indirectly by increasing the average sojourn times of the best-effort flows are affected due to their reduced bandwidth share.

Summarizing, best-effort flows have a significant impact on the average queuing delay, impacting not only mobility management and related signalling, but also any kind of traffic which is delay sensitive. An average penalty of nearly 30ms on packet latency for semi-loaded scenarios would have impact on streaming traffic with small play-out buffers like voice over IP.

### 2.4.4.3 Signalling message sequences

For the computation of message sequence durations we add up the individual transmission delays of the link components of each message according to the source and destination. For example, the path delay of a message which is sent from a UE to the MME consists of the air interface delay, the intra-EFN delay, the backhaul link delay and the MNO delay. The total message sequence duration corresponds then to the sum of all link delays of all messages.

In the first scenario without LIPA support, we compare the normal S1-based hand-over and hand-in cases as described in [7] with the inter-HeNB hand-over and the S1-to-X2 hand-in. Note that we assume that a tracking area update in the inter-HeNB hand-over is not necessary.

**Figure 2.54: Handover completion times without LIPA support.**

**Figure 2.55: Handover completion times with LIPA support.**

The results are shown in Figure 2.54. We compare the S1-based inter-HeNB hand-over (circle markers) with the localized S1-based hand-over (square markers) and the S1-based hand-in (diamond markers) with the localized S1-to-X2 hand-in into the enterprise network (triangle markers).

The results show that depending on the backhaul link latency, the S1-based hand-over requires between 282ms and 395ms for completion. The reason is that in this sequence the message transactions between the HeNBs and the MME, and the intra-MNO messages dominate the overall hand-over time. Since the localized S1-based hand-over does not require any exchange of signalling towards the mobile operator, this method requires with an average of 52ms significantly less time to complete due to the low latencies in the enterprise network.

The S1-based hand-in to the EFN requires less time than the S1-based inter-HeNB hand-over since in this case the path between HeNB in the enterprise domain and the MME in the operator domain over the backhaul link is traversed only 4 times instead of 9 times as in the S1-based inter-HeNB hand-over. This also reduces the dependency on the bottleneck queuing delay which is reflected in lower steepness of the curve. However, the localized S1-to-X2 hand-over requires only two message exchanges between enterprise and mobile operator, such that an increasing backhaul link latency leads to even less impact on the overall hand-over completion time.

In the next scenario, we compare the same message sequences in the case that the LFGW supports LIPA. Here, the bottleneck queuing delay depends on the offered load of the QoS traffic only, which only has a minor impact in the low load scenario we are considering. Correspondingly, the hand-over completion times in Figure 2.55 appear nearly constant independent from the offered load. The minor increase in backhaul queuing is overlaid by the other components of the overall delay.

Summarizing the results of both scenarios from a feature point of view, both LIPA and signalling localization can reduce hand-over completion times significantly. LIPA reduces the completion time by more than 30% in the for S1-based inter-HeNB handovers. Localization leads to significant gains independently from the bottleneck link utilization: more than 80% for inter-HeNB handovers, and approx. 52% for hand-in scenarios.

In the case of PDN connection setup the situation is different, as visible in Figure 2.56. The figure shows the PDN connection setup completion times with/without LFGW support (circle and square markers) and with/without LIPA (solid and dashed lines). The PDN connection setup with LFGW requires some more message exchanges between the L-GW and the P-GW, which leads to a higher dependency on the backhaul link queuing delay. This is illustrated by the steeper curve in the no-LIPA scenario if compared to the PDN connection setup without LFGW support. However, LFGW support also requires less message exchanges in the MNO, such that the completion time for scenarios with low load (i.e. until $\rho = 0.4$) is still lower than in the case of without LFGW support. The results for the case with LIPA support corroborate this interpretation, since here the gain of the message sequence with LFGW is nearly constant at approximately 16%. As in the case of the hand-over messages, LIPA and signalling localization at the LFGW leads to a significant decrease of message sequence completion times.

**Figure 2.56: Time to complete PDN connection setup with and without LIPA support.**



**Figure 2.57: Capacity gains and per-flow throughput for local services.**

#### 2.4.4.4 Impact on Backhaul Capacity

Due to the double routing of local service best-effort traffic over the backhaul link, a significant amount of capacity is unnecessarily "wasted". Figure 2.57 shows the capacity gain of a scenario with LIPA support over an equivalent system without LIPA. The solid lines with square markers denote the downlink capacity gain, while the solid line with circle markers denotes the uplink capacity gain. On the x-axis the offered QoS load is shown, which is in the case of the non-LIPA scenario the half of the total offered load according to the scenarios in the above sections. The unused time-average capacity of the link corresponds then to $C - E[n_{QoS}] \cdot E[R_{QoS}] - (1 - \pi_{BE}(0)) \cdot C_{ul}$. Note that $C_{ul}$ is the uplink capacity and $C$ is either uplink or downlink capacity. The results show that the capacity gain increases up to 30% in case of the downlink for an offered QoS load of $\rho_{QoS} = 0.2$. The dashed line with diamond markers corresponds to the average per-flow throughput in the non-LIPA case. Due to the increasing number of best-effort flows and QoS flows it decreases from nearly the full uplink capacity to below 3Mbps. In the case of LIPA, the access to the local services would have the full capacity of the EFN, which is not shown in the above figure.

### 2.4.5 The impact of traffic offloading on network performance of MNOs

#### 2.4.5.1 Introduction

Nowadays there is enormous growth of the number of a new generation of mobile devices like various smart phones (iPhone, Android-based), laptops, netbooks, etc. in the market. At the same time, mobile networks operators are incorporating actively Internet applications and services for the mobile devices. There are thousands of web data applications and services available now (e.g. YouTube, Facebook, Spotify, IM, mobile TV, etc.) that are becoming extremely popular in the mobile user environment. According to the Cisco VNI Global Mobile Data Traffic Forecast [79], overall mobile traffic is expected double every year from 2011 onwards.

As a result of these two factors, there is an explosion of both data and signalling traffic towards the core network of Mobile Network Operators (MNOs). As a consequence, congestion situations can arise in the core network. Thus, solutions to avoid unnecessary traffic load at network nodes are needed. One such solution is to apply a traffic offloading mechanism by means of femtocells. It can solve macro core network capacity crunch avoiding future upgrades of the network infrastructure.

From a network management perspective it is very important to correctly understand performance benefits (if any) of offloading in order to develop correct deployment strategies. Failing into doing so may lead to investing large efforts in installing costly offloading infrastructures that do not bring targeted benefits. In order to pursue these objectives it is important to characterize and model user behaviour and data hotspots from the viewpoint of offloading strategy that MNOs plan to use. As a result, one must assess if the gain in network resources obtained by means of traffic offloading is enough to avoid network congestions.

Thus, the objective of this study is to analyse which are the benefits, in terms of performance, coming from the deployment of an offloading network.

In this study, we consider the traffic offloading process from the viewpoint of a MNO perspective. In this context it is important for the MNO to care about a traffic load that goes to its core network (e.g. 3G CN) after offloading. What is really needed for the MNO is to know what was before offloading and what happens after in the context of traffic parameters to adjust with them dimensioning characteristics of its network, e.g. to evaluate required system capacity after offloading.



**Figure 2.58: A view of a network implementing offloading**

A view of a network implementing offloading is presented in Figure 2.58. It contains some offloading areas within the MNO network coverage. Thus, with the advent of femtocells, some parts of the network of the operator could be offloaded when the user is at home or at an enterprise, as his traffic would be routed through the femtocell.

To describe performance of mobile networks implementing offloading an appropriate analytical model is needed. In the next subsection we propose a model for traffic of a single source that eventually reaches the network of the operator after offloading.


### 2.4.5.2   Traffic model for offloading

In the same way as in previous literature, the activity of a single source is modelled as a strictly alternating ON/OFF process [81]. In our model, we introduce the effect of offloading over such source by means of an additional strictly alternating ON/OFF process. The resulting traffic sent to the network of the MNO in a regular way (i.e., the non-offloaded traffic) from a single source is then modelled as the product of the two above processes. Therefore, the aggregation of several such resulting processes models the traffic still needing to be handled by the operator in the conventional way, which is the one to be used for the purposes of network dimensioning. The following subsections introduce the details and notation of all above processes.


### 2.4.5.2.1   Model of user activity

Thus, the activity of a single source/user is modelled as a strictly alternating ON/OFF process, where ON periods are i.i.d., OFF periods are i.i.d., and ON and OFF periods are independent. Furthermore, previous measurements have also shown that such periods follow heavy-tailed distributions (e.g., due to file size distributions, web pages) [82] (and references therein). This is also assumed in our model.

Therefore, the process the process Y (t) is a stationary binary time series $\{Y(t), t \geq 0\}$ such that [81]:

$$Y(t) = \begin{cases} 1 & \text{activity period} \\ 0 & \text{idle period} \end{cases}$$

A representation of such a process is presented in Figure 2.59. During the activity period, the user is transmitting and receiving packets. On the other hand, no traffic is exchanged with any network during idle periods, e.g., user reading time of downloaded content. Notice that Y(t) only describes source/user behaviour and is independent from the network through which the traffic is sent, i.e., no offloading considerations have been made in this process.

**Figure 2.59: User activity is modeled as a strictly alternating ON/OFF process, Y(t)**

### 2.4.5.2.2  Model of offloading periods

Previous measurements carried out in real networks have shown that smartphone connection and disconnection periods to offloading areas follow heavy-tailed distributions [80]. This observation led us to also model offloading periods for a single source as strictly alternating ON/OFF process. In the same way as before, we also assume that ON periods are i.i.d., OFF periods are i.i.d., and ON and OFF periods are independent. Therefore, the process X(t) is a stationary binary time series $\{X(t), t \geq 0\}$ such that

$$X(t) = \begin{cases} 1 & \text{user flow sent through MNO network} \\ 0 & \text{user flow sent via offloading network} \end{cases}$$

A representation of such a process is presented in Figure 2.60. During ON periods, the traffic generated by the source (if any) would be sent towards the network of the MNO in a regular way (i.e., traffic is not offloaded). On the other hand, during OFF periods all traffic generated by the source would be routed through the offloading network (e.g., smartphone under the coverage of a femtocell).



**Figure 2.60: Offloading periods are modeled as a strictly alternating ON/OFF process, X(t)**

### 2.4.5.2.3  Model of non-offloaded traffic from a single source

The traffic generated by a single source that is treated in a conventional way by the MNO (i.e., non-offloaded traffic) can be modelled as the product of the previously defined processes. That is,

$$Z(t)=X(t)Y(t).$$

Figure 2.61 represents such a process, which is also a strictly alternating ON/OFF process, with ON and OFF periods following heavy-tailed distributions and whose characteristic parameters can be derived from those of the original ones, as shown below. During ON periods, the traffic being generated by the source (i.e., user activity in ON state) is forwarded to the network of the MNO as usual. On the other hand, during OFF periods, either there is no activity from the source or traffic is being sent through the offloading network (e.g., through Wi-Fi, femtocells).



**Figure 2.61: Non-offloaded traffic from a single source is modeled as the product of two strictly alternating ON/OFF processes, Z(t)=X(t)Y(t)**

### 2.4.5.3 Problem formulation and high-level view of steps followed to solve it

The goal of this study is to evaluate the potential benefits that offloading techniques could bring to operators. In other words, their interest is to compare what is the resource consumption in their network without applying offloading and when applying offloading on a large scale. Hence, operators are interested in the behaviour of the process modelling the aggregation of several Z(t) processes. Therefore, our problem is to obtain an analytical model describing the behaviour of such an aggregated traffic and the resource consumption it entails.

A high level view of the steps needed for solving the problem follows. Since, as explained above, Z(t) has ON and OFF periods that follow heavy tailed durations, this process is long-range dependent. Therefore, the aggregation of several such processes is self-similar and can be characterized by means of the Hurst parameter ($H$). This can be done by means of the same techniques presented in [81]. In turn, parameter $H$ can be derived from the parameters characterizing the heavy-tailed behaviour of the ON and OFF periods of Z(t). Besides, such parameters can be obtained from those of processes of the original processes X(t) and Y(t). Therefore, once the parameters of the original processes are known, by following the above steps, we can characterize the behaviour of the aggregated non-offloaded traffic, and hence, the resources needed in the network of the MNO to serve it.

As an outcome of the problem solution performance bounds related to the provisioning of the resources in the network of the MNO when applying offloading can be evaluated. As an example, Figure X illustrates the amount of resources that the system needs to offer the required quality of service. In particular, the figure shows the capacity needed in the original system together with the worst and best case scenarios when we offload a 50% of the traffic.



**Figure 2.62: Bounds on resource needs vs. variance coefficient (*a*) with 50% of offloaded traffic**

As seen from Figure 2.62 offloading does not necessarily entail less resource consumption in the network of the operator. Under certain conditions, and due to an increase of the burstiness of the non-offloaded traffic (tail index α tends to 1), the amount of network resources to offer a given level of QoS can be increased.

More detailed analysis of performance evaluation issues related to resource consumption in the network of the MNO deploying an offloading strategy is expected to be presented in D5.3 "Access control solutions and evaluation". The performance analysis is based on results of an analytical framework for modelling traffic in mobile networks implementing offloading. The analytical framework including the proposed traffic model for offloading and validation of the main analytical results by means of simulations are also planned to present in D5.3.

### 2.4.6 Conclusion

Operators and vendors are interested in extending the use cases for femtocell deployments towards enterprise solutions. However, with the current architecture as specified by 3GPP, such deployments face technical and economic challenges which we addressed in this paper. Specifically, we proposed a localized mobility management function which minimizes hand-over signalling over the costly backhaul link between mobile network operator and enterprise deployment, as well as support for local IP access

(LIPA) and selective IP traffic offload (SIPTO) at a novel network entity denoted as Local Femtocell GateWay (LFGW). The proposed solution is fully standards-compatible, as it requires no modifications on the EPC.

The performance evaluation of the proposed solution shows that signalling message sequence completion times such as for handovers and for PDN connection setup can be decreased significantly, reducing resource occupation and increasing user QoS. Furthermore, especially LIPA leads to a generally lower utilization of the backhaul link such that user QoS and potentially costs can be optimized. The proposed architecture is thus a complete solution for enabling enterprise femtocell network deployments for 3GPP networks.

Operators are also interested in benefits (if any) that traffic offloading implemented by means of femtocell deployments can give for its network in terms of resource consumption. Hence, a model describing traffic that eventually reaches the network of the operator after offloading (non-offloaded traffic) and should be served on a regular basis is needed. Based on previous measurements found in the literature, we proposed a traffic model for offloading assuming that user's activity periods and periods characterizing offloading are heavy-tailed. We model them as strictly alternating independent ON/OFF processes. Therefore, the non-offloaded traffic is modelled as the product of these two processes. The presented preliminary results based on the model show that offloading does not always entail a gain in terms of network resources. The detailed results on performance bounds of the resource consumption in the network of the operator implementing offloading and the analytical framework to get these results are expected to be presented in D5.3.

## 2.5 A QoS based call admission control and resource allocation mechanism for LTE femtocell deployment

### 2.5.1 Motivation

The Long Term Evolution (LTE) standards for cellular systems have opened many new possibilities for future mobile communications. These include concepts such as advanced Self Organizing Networks (SONs), policy based network management and further integration of femtocells and femtocell networks [29][30]. The work presented in this study spans each of these three ideas. Specifically, this study addresses the problem of congestion that can occur when a large number of femtocells utilise DSL as a backhaul link. A method to enable both Quality of Service (QoS) aware Call Admission Control (CAC) and bandwidth negotiation in the backhaul links is proposed. The algorithms developed are suitable for SON deployment and can be driven by a number of parameters that can be adjusted based on operator-defined policies.

In LTE, the home base station is referred to as a Home evolved Node B (HeNB) and connects to an HeNB Gateway (HeNBGW) through the customer's fixed broadband connection, typically a Digital Subscriber Line (DSL) connection. This critical backhaul segment is often neglected in both research work and in real femtocell deployments; it is assumed to be of secondary importance to the radio link between the mobile unit and the home base station. Although the femtocells transport both voice and data, due to their real time requirements the voice connections are much more sensitive to constraints in the backhaul network. Traditionally all voice connections would utilise dedicated resources along the path of the voice call meaning that there would be little or no degradation in call quality due to limited network resources. Unfortunately, when transporting voice calls over residential internet connections, as is the case in most femtocell deployments, it is often not possible to provide guaranteed network resources.

Each voice call transported via an internet connected femtocell is essentially converted into a Voice over Internet Protocol (VoIP) call and encapsulated into a tunnel (Figure 2.63) for transmission to the femtocell gateway located in the operator's core network. 3GPP have defined AMR and AMR-WB as the mandatory voice codec for LTE deployments. Each AMR VoIP call has a data rate in the region of 5 to 24 Kbps plus overhead, depending on the codec mode that is being used. In the case of IPSec secured tunnelling that is used by all femtocell deployments the overhead becomes comparable to the size of the actual payload.

| | | | |
|---|---|---|---|
| AMR Header | 6 | | |
| AMR Data | AMR SID | AMR 12.2 | |
| | 4 | 32 | |
| RTP | 12 | | |
| UDP | 8 | | |
| IP | 20 | | |
| GTP-U | 8 | | |
| UDP | 8 | | |
| IP | 20 | | |
| IPSec ESP | 8 | | |
| PPPoE | 12 | | |
| Ethernet | 42 | | |

**Figure 2.63: Voice over S1-U**

Upon egress from the femtocell each VoIP packet is tagged as voice using the Differentiated Services Code Point (DSCP) field of the IP header, this allows the DSL Access Multiplexer (DSLAM) to identify it as voice and provide low-latency prioritization. In this work it is assumed that each Mobile Network Operator (MNO) operating femtocells utilizing the Internet Service Provider's (ISP) network can be individually identified on a per flow basis. This could be done by providing specific VLANs for each MNO using the ISP network, therefore allowing the DSLAM to use the VLAN tag of each packet to identify the MNO to which it belongs.



**Figure 2.64: Femtocell deployment and architecture.**

At the DSLAM all voice traffic is forwarded through the Expedited Forwarding (EF) queue. The EF queue provides the highest service level and is utilised for low latency and low bandwidth intensive services such as voice; however in typical deployments it has a fixed and relatively limited bandwidth. As the number of FAP deployments increases there will be a growing volume of voice traffic being backhauled through existing DSLAMs. It is therefore likely that as the number of calls increases the EF queue on DSLAMs will become increasingly congested and at a certain level may exceed the maximum allocated bandwidth. Indeed it is also likely that the ISPs may begin to limit the amount of voice traffic that each MNO can place in the EF queue. This will inevitably lead to increased packet delays, jitter and loss at the DSLAM with a corresponding impact on the voice call quality. Therefore, the quality of the call depends on both the radio link from the customer's mobile device to their HeNB and the level of congestion on the backhaul link.

In this work it is assumed that each MNO will have Service Level Agreements (SLA) in place with the ISP to provide a maximum and minimum bandwidth at each DSLAM dedicated for their VoIP traffic.

Although the MNO could take the typical telecommunications response and overprovision the dedicated resources on each DSLAM this would have increased costs. It is therefore of interest for the MNO to minimise the maximum bandwidth that is reserved on each DSLAM while maintaining high voice call quality for their customers.

## 2.5.2 Related Work

There have been a number of other QoS based CAC mechanisms. In particular, an End-to-end Measurement based Admission Control (EMBAC) mechanism for VoIP was described in [31]. Here the authors use fake VoIP streams as probes to emulate real VoIP traffic. This is also similar to previously published work on QoS based handover mechanisms for VoIP in which the authors utilise VoIP like probes to estimate call quality on candidate handoff networks [32]. Our concept differs from these in two respects. Firstly, we locate the objective voice quality measurements in an intermediary node (i.e. HeNBGW) to avoid imposing the requirement for mobile clients to perform call quality measurements. Secondly, our call admission decisions utilise measurements taken from actual ongoing VoIP calls rather than introducing voice probing streams to aid admission decisions. In [33] the authors motivate the need for CAC for voice calls over Packet Switched Networks (PSN). They emphasize that there is a need of a scalable and dynamic solutions to address the situation where VoIP calls are forwarded through a congested network in which no over provisioning exists.

The mechanism proposed in this study is a combination of traffic policing and traffic shaping. A lot of previous work has been done on both of these concepts and an overview of each is described in [34]. Both concepts imply the existence of a SLA between the ISP and the users (e.g. MNOs). An example of network specific metrics specified in such agreement is the CIR (Committed Information Rate), committed burst, and excess burst rate.

In order to compute the Mean Opinion Score (MOS) of the VoIP calls a modified variant of the E-Model algorithm [35] is used; specifically, our implementation of the E-Model has been adapted to support AMR voice. The E-Model computes a MOS for each VoIP call based on a number of network metrics including delay, loss and jitter. Accurate end-to-end network delay can be difficult to obtain unless the nodes between which it is calculated are time synchronized. In [36] the motivation and constraints for maintaining time synchronisation between femtocells and the HeNBGW is presented. It is shown that rigorous requirements for clock and frequency synchronization have to be met, NTP [37] or PTP [38] being the recommended synchronisation protocols. Based on this the HeNBGW maintains regular NTP message exchange with each femtocell to ensure high clock synchronisation accuracy. In this work we assume that such clock synchronisation exists. Indeed we analysed network traces from real femtocell deployments to confirm that this clock synchronisation is used.

## 2.5.3 Architecture Description

Figure 2.64 shows the network architecture of a typical femtocell deployment scenario. This was used as the reference architecture on which the proposed control mechanism and simulations were developed. It is comprised of multiple households and small offices that have femtocells provided by the same MNO and connected through the same DSLAM to the MNO's core network. It is assumed that each femtocell is connected to the broadband router.

Information gathered from multiple DSLAM vendor datasheets showed that each DSLAM is capable of serving approximately 1000 households or small offices. As each femtocell can support an average of 4 user devices this means that potentially each DSLAM may have to transport up to 4000 simultaneous voice calls.

The other major parts of the MNO's core network are also shown but for the simulations these are assumed to have little or no impact on the call quality and, as such, are modelled with fixed low delays. The alternate call endpoints are also assumed to be connected via low delay links that have no impact on call quality.

Based on the current network architecture we determined two ways of creating the required interface between the MNO and the ISP as depicted in Figure 2.64.

**Option 1** would be a proprietary interface developed by the ISP to allow the MNO to have limited dynamic resource allocation capabilities in the DSLAM.

**Option 2** would utilise an interface being defined by both the broadband forum and 3GPP for fixed/mobile convergence [39]. This PCRF-BPCF interworking would allow the MNO to push resource allocation requests to the fixed access network via the S9 interface. It would then be the responsibility of the BPCF to perform an admission control decision on the resource request and update the resource allocation in the DSLAM accordingly.

### 2.5.4 QoS Based Call Admission Control

#### 2.5.4.1 Call Quality Monitor

The proposed call admission control and dynamic resource allocation mechanism is based upon the quality of ongoing voice calls passing through the HeNBGW. In order to compute the call quality a *Voice monitor* application residing in the HeNBGW was developed. The role of this module is to maintain a list of ongoing calls, measure the real time voice call quality, determine when problems occur and employ the admission control and dynamic resource allocation mechanisms to restore/maintain the quality.

The *Voice monitor* uses the CIR value from the SLA to calculate a *nominal number of calls*. That is the maximum number of calls that can be supported with the default bandwidth allocation provided by the ISP to the MNO. It is computed as the ratio between the CIR [bps] and the call bit rate plus tunnelling overhead (using AMR's highest mode).

#### 2.5.4.2 MOS calculation

When a new packet arrives, the monitor calculates its delay, jitter and packet loss. The delay is calculated as the difference between the current time and the time when the packet was sent. The jitter is calculated using the recommended RTP jitter algorithm [40]. The RTP sequence numbers are used to compute a moving average of packet loss with a window size of 100 packets. With the formula provided by the E-model [40] the transmission factor R and the average MOS is calculated.

Keeping track of each MOS value calculated could arise in scalability issues as normally around 50 of such values are obtained per second per voice call. Thus, we average all MOS values obtained in a time window and obtain an average MOS value.

The monitor uses the individual average MOS values of each call to compute an *average MOS* value across all calls passing through the HeNBGW. We will use the term *average MOS* further in this study to refer to the above mentioned value. The *average MOS* is then used as the basis of the mechanism performing call admission control and resource reservation in the DSLAM.

An alternative approach to that described in this work would be to allocate resources the moment call requests are received in the HeNBGW without requiring any call quality monitoring. However, although we assume that there is an EF queue per MNO or subscriber this queue may be shared by other EF traffic sources in the HAN network. As such fixed resource allocation would have no way to determine if the allocated resources were sufficient to support high quality voice.

| Quality rating | MOS |
|---|---|
| Best | 4.34 - 4.50 |
| High | 4.03 - 4.34 |
| Medium | 3.60 - 4.03 |
| Low | 3.10 - 3.60 |
| Poor | 1.00 - 3.10 |

**Figure 2.65: Map between Quality Rating and MOS**

#### 2.5.4.3 Call Admission Control algorithm

Figure 2.66 shows a flowchart describing the decision process. The quality of all calls passing through the HeNBGW is continually monitored in real time and the *average MOS* value represents the impact of the associated DSLAM's EF queue.

There is a possibility that some calls may suffer quality degradation due to issues in the Home Area Network (HAN) in which the femtocell is installed. Our mechanism detects this situation and flags those calls with HAN problems.

The CAC and dynamic resource allocation mechanism is triggered periodically. When triggered, it is first checked whether each call's average MOS value is outside of the 95% confidence interval range. If so, those falling in that category will be filtered out and flagged as having HAN problems. A new *average MOS* value for all calls without HAN problems is then computed.

If this new *average MOS* value is less than 3.8 then no new calls will be accepted until the resource allocation in the DSLAM can be increased and the backlog of packets in the queue is cleared. If the current bandwidth used in the EF queue is lower than the committed burst then the VoIP monitor will request more bandwidth.

If the *average MOS* is higher than 3.9 then new call requests will be accepted. If it is higher than 4.0 and the number of simultaneous calls has decreased then the level of bandwidth in the EF queue is decreased, but no lower than the committed information rate (CIR).

The average MOS thresholds were decided considering Figure 2.65 obtained from [35]. The proposed mechanism attempts to maintain all calls in the *Medium* to *Best* quality range.

#### 2.5.4.4 Bandwidth monitoring in the EF queue and traffic shaping

The DSLAM monitors the traffic volume forwarded through the EF class [41] by each MNO. The calculation of the utilised bandwidth in the EF queue is done periodically and is based on the AMR inter-packet departure time, i.e. 20 milliseconds.

When, for a certain MNO, the CIR is reached any new VoIP packets arriving in the EF are tagged as noncompliant. Noncompliant packets are stored until the next scheduled release time. However, the buffer has a limited storage size and thus traffic policing needs to be employed once the buffer has reached its limit. Any packets arriving when the buffer is full are therefore dropped.



**Figure 2.66: Call Admission Control and Dynamic Resource Allocation Mechanism**

### 2.5.5 Simulation setup

The proposed solution was implemented and validated through simulations using the Network Simulator 3 (NS3). NS3 has been designed to overcome many of the existing network simulators' problems. It succeeded in many aspects [42], the most important being simulation time and scalability.

The simulated scenario is depicted in Figure 2.67. It consists of User Equipments (UEs) which generate VoIP traffic and are associated to Femto Access Points (FAPs). For ease, we used a wired link between the UEs and FAPs, as the radio link is not the scope of this study. We simulated 50 FAPs, each having 4 UEs associated. An additional FAP with only one UE associated was used to create a HAN problem scenario. All FAPs are connected through DSL lines to a DSLAM. The links are dimensioned as best case DSL lines, that is a speed of 100 Mbps and a link delay of 1 ms. For the particular HAN problem scenario, the link speed was decreased to 73 kbps.

The DSLAM forwards VoIP packets through its egress EF queue to the VoIP monitor (HeNBGW). This network part symbolizes the ISP's core network plus additional links needed to reach the MNO's HeNBGW. The VoIP monitor calculates the metrics used by the mechanism described in the previous section and forwards the packets through the MNO's core network to the terminating UEs.

In our simulation, simultaneous VoIP calls are consecutively generated with an increment 0.5 seconds plus a random variable in the range [0, 0.1] seconds with a granularity of 10 microseconds. The random variable is used to mitigate any potential synchronisation between VoIP sources which would results in unrealistic packet arrivals. A simulation time of 300 seconds was chosen and we simulated 201 VoIP calls with call durations ranging from 90 to 290 seconds. The voice calls are torn down in the same randomised manner as they are created.



**Figure 2.67: Simulation Scenario**

In order to accurately emulate DSLAM behaviour, DSLAM recommendations from Cisco were used [34]. These provide recommendations on a number of DSLAM parameters including the buffer size used for traffic shaping which is based on link speed. From DSLAM datasheets a typical egress link speed of 10Gbps was obtained and so this value was used in the simulation. For this link speed a buffer size of 45,000 packets or 11,520,000 bytes, assuming 256 bytes as a typical packet size is recommended.

Given these parameters the nominal number of calls is calculated using the values from Figure 2.65. From this the maximum size a VoIP packet passing through the DSLAM is determined to be 176 bytes assuming that the mode AMR 12.is being used. Based on this the data rate per call is computed as:

$$176\frac{Bytes}{packet} \times 8\frac{bits}{Byte} \times 50\frac{packets}{Sec} = 70.4kbps$$

In the results presented in this study a CIR value of 10 Mbps and a Committed Burst Value of 15Mbps is assumed. In order to impose the Committed Burst Value limits and perform traffic shaping, the EF queue bandwidth for each MNO is checked every 10ms which represents half of the AMR packet inter-departure time. In this way no unnecessary delaying of packets is introduced.

Using a CIR value of 10 Mbps in the EF queue, the maximum number of supported calls can be calculated as:

$$\frac{10Mbps}{70.4kbps} \approx 140\ Calls$$

For ease of implementation, only two modes of the AMR speech codec, i.e. AMR 12.20 and AMR SID were implemented. The switching between these two is specified by the speech activity parameter which is set as a random value between 30% and 80%. This emulates realistic Voice Activity Detection and silence suppression functionality of the AMR codec. The VoIP monitor in the HeNBGW computes an average MOS over all ongoing voice calls every 0.5 seconds; this was chosen as a reasonable trade-off between maintaining high quality voice calls and reducing the level of overhead.

## 2.5.6 Results

This section presents a number of simulation results obtained using the setup described in the previous section. In order to obtain a baseline for the results a simulation without the proposed QoS and CAC mechanisms in place was performed. The results of this simulation are depicted in Figure 2.68 and show the *average MOS* and the number of ongoing calls over time.

As can be seen, without CAC all calls are accepted and voice packets are forwarded from the EF queue as soon as they arrive; this is unless the imposed bandwidth limitation has been reached in which case the packet is queued. This queuing means introducing extra delay and this delay increases until the buffer is full at which point the buffer will overflow resulting in dropped/lost packets. At this point the MOS of all ongoing calls degrades rapidly. From Figure 2.68 it can be determined that the maximum number of calls that the DSLAM can accommodate is approximately 130, given our specific simulation setup and assumptions.

The fact that all voice calls are degraded draws attention to the difference between circuit switched and packet switched voice traffic. In circuit switched networks dedicated time slots are allocated for voice traffic at call setup time, while in packet switched networks voice packets from different voice calls potentially share the same packet forwarding capacity at each network router. In other words, accepting all requests overloads the forwarding queue resulting in call quality degradation of all calls utilising the same queue. This further highlights the need for dynamic CAC.

Figure 2.69presents results for simulations done using the proposed QoS and CAC mechanisms enabled. It can be seen that when the overall MOS drops, the feedback mechanism rejects any new call requests and the overall MOS for all calls is restored to a high level. Figure 2.70 shows a plot of the bandwidth requests made by the algorithm in the FGW and granted by the DSLAM.

It can be seen that for each moment when the average MOS is below the threshold (3.8) the requested bandwidth on the DSLAM is increased up to the maximum committed burst. This extra bandwidth is released only when the average MOS has returned to a high value and the number of simultaneous calls has decreased. Since any extra reserved bandwidth would involve a cost to the MNO, increased bandwidth requests are only made when necessary, thereby maintaining a trade-off between the extra bandwidth requested and the number of accepted calls.

A third scenario was also simulated in which we introduced a femtocell whose voice traffic is being degraded in the HAN. In this case the algorithm was able to detect and remove that call from the measurement data used to make resource and CAC decisions.



**Figure 2.68: Overall MOS and number of online sessions versus time, without the CAC and dynamic resource allocation mechanism**

**Figure 2.69: Overall MOS and number of online calls versus time, with the CAC and dynamic resource allocation mechanism**



**Figure 2.70: DSLAM's EF queue bandwidth usage versus time**

# 3. Mobility Management

Mobility management mechanisms to support seamless handover with minimal signalling cost and enable efficient location management for femtocells are developed in this section. In subsection 3.1, a Local Femto GateWay (LFGW) with proxy MME and proxy S-GW functionalities is introduced to enterprise femtocell networks to reduce the handover latency. The complete work has been reported in D5.1. A local location management architecture is developed and a self-organized tracking area List mechanism for large-scale networks of femtocells is proposed in subsection 3.2 to reduce the location signalling cost. Subsection 3.3 proposes mobility management schemes based on X2 traffic forwarding chain for networked femtocells to mitigate the heavy signalling overhead to the core network due to the frequent inter-femto handover. A fast handover failure recovery mechanism is presented in subsection 3.4 to reduce the service interruption latency in case of a handover failure during inbound/outbound mobility. In subsection 3.5, basic 3GPP relay architecture alternatives for supporting mobile femtocells have been compared and a latency analysis has been given. Finally, handover and configuration performance of networked femtocells in an enterprise LAN are demonstrated in subsection 3.6 based on the proof-of-concept implementation

## 3.1 Local Mobility Management

The work on this item was finished during project year 1 and can be found in D5.1.

## 3.2 Local Location Management

### 3.2.1 Introduction

This section describes a Local Mobility Management (LMM) scheme in the context of a large-scale, all wireless network of femtocells (NoF). Traditionally, mobility management has been divided into handoff management and location management. The former focuses on keeping ongoing communications active when a User Equipment (UE) performs a handover between neighbouring femtocells. The latter deals with updating databases that store information about UE location within the network of femtocells. LLM schemes guarantee that new communications towards a destination UE can be effectively established by previously requesting its location to such databases.

Standard 3GPP location management mechanisms have been designed with macrocell scenarios in mind. Therefore, their performance in large-scale networks of femtocells is far from optimal due to the overhead generated by frequent handovers and cell reselections. Thus, NoF scenarios require specific location management mechanisms in order to track UEs efficiently whilst keeping location signalling traffic under control.

In the context of large-scale, all-wireless networks of femtocells, location management schemes provide a mechanism that enables network entities to map a subscriber's identity to the identity (and, subsequently, the location) of the HeNB where the UE is currently camped on. On the one hand, standard 3GPP identifiers (such as the International Mobile Subscriber Identity (IMSI) and/or the Serving Temporary Mobile Subscriber Identity (S-TMSI)) are used to identify users within the cellular network. On the other hand, the serving HeNB address varies as the destination UE moves throughout the network of femtocells. This address is used by the underlying transport network to route packets towards the destination UE.

In order to avoid interoperability problems with current HeNB-GWs and signalling traffic overload of the Evolved Packet Core (EPC), local mobility management mechanisms should only affect the network elements in the network of femtocells (i.e., the local network) and not the network elements of the EPC.

### 3.2.2 Work during Year 1

The work carried out during the first year of the project focused mainly on the definition of a LLM scheme (along with its associated protocol architecture) that was capable of providing location management functionalities in the context of a large-scale, all-wireless network of femtocells. This location management scheme is based on VIMLOC [26], a distributed, wireless mesh network-oriented location management mechanism in which location information is distributed across all HeNBs in the network of femtocells.

In order to integrate native 3GPP location management procedures with VIMLOC, modifications to the protocol architecture of HeNBs were needed. As a result of this, a 2.5 protocol layer (also referred to as geosublayer) was inserted between the network and access layers. One of the main functions of the

geosublayer was to intercept all 3GPP control-plane location management messages in order to trigger the corresponding VIMLOC procedures in the network of femtocells.



**Figure 3.1: The 2.5 Layer (Geosublayer) in the LLM Protocol Architecture.**

### 3.2.3 Open Issues

VIMLOC relies on native 3GPP user location mechanisms in order to determine the IP and geographic addresses of the HeNB where the destination UE is currently camped on. In order to do so, 3GPP location mechanisms must be able to determine, with the granularity of a single femtocell, the current location of the destination UE within the NoF. Since standard 3GPP location mechanisms can only determine the location of a UE with the granularity of a Tracking Area (TA), a first approach to the problem of fine-grained user tracking could be to reduce the size of all TAs in the NoF to that of a single femtocell. However, this solution is far from optimal, as it requires UEs to perform Tracking Area Update (TAU) procedures every time they reselect/handover between neighbouring femtocells.

In static/semi-static scenarios, UEs perform little or none reselections/handovers between neighbouring femtocells. Thus, reducing the TA size to that of a single femtocell does not have a negative impact on the amount of location signalling traffic associated with Tracking Area Update procedures. Furthermore, as S1-AP Paging messages are always sent to serving HeNBs, signalling load associated with paging procedures is kept to a bare minimum.

In a dynamic scenario, UEs may move rapidly throughout the network of femtocells. This involves frequent reselections and handovers between neighbouring femtocells which, in turn, increases signalling traffic associated with Tracking Area Update procedures. Thus, reducing the size of TAs to that of a single femtocell in dynamic scenarios is not optimal in terms of signalling traffic.

### 3.2.4 Proposed Solution

In this section we propose a 3GPP-compliant self-organized Tracking Area List mechanism for large-scale networks of femtocells. Our scheme can be deployed in legacy MMEs by means of a software update and does not require modifications to the UE.

Prior to the description of the self-organized mechanism, we need to characterize the mobility pattern of a UE in order to estimate the average TAU arrival rate in a TA ($\bar{\lambda}_{TAU}$). This is done in subsection 3.2.4.1. A detailed description of the self-organized TAL mechanism, along with an analytical model for its performance evaluation, is provided in subsection 3.2.4.2.

#### 3.2.4.1 Mobility Model

In [72], the author provides a comprehensive overview on mobility models for cellular networks. In this solution we have assumed a Markov-based mobility model for a 2D hexagonal topology in order to facilitate mathematical tractability.

The concept of slotted time is implicit in Markov-based mobility models. At the end of each timeslot, the UE remains in the current cell with probability $p$ or, alternatively, transits to an adjacent cell with probability $\frac{1-p}{6}$. Variations in the UE speed can be modelled by modifying the value of the parameter $p$.

Let us consider that the cell topology in Figure 1 is a TA. We define the concept of cell ring as a group of neighbouring cells where a UE can be found in a certain timeslot. The first ring is formed by a single cell located in the centre of the TA. The second ring is formed by all one-hop external neighbours of the first ring. The third ring is formed by all one-hop external neighbours of the second ring, and so on.

The concept of inner and vertex cells is also depicted in Figure 3.2. By definition, vertex and inner cells belong to the outermost ring of a Tracking Area. A vertex cell provides three exit points from the TA, while an inner cell provides two. This will be considered during $\overline{\lambda}_{TAU}$ calculation.



**Figure 3.2: Markov-based mobility model in a hexagonal cell topology.**

Markov-based mobility models can be mathematically described by discrete-time Markov chains (DTMCs). States in the DTMC represent the cell rings in the Tracking Area, as described in Figure 3.2. Analogously, transitions between states are associated with the probability of a UE staying in the current cell ring or moving (outwards or inwards) to an adjacent one. Figure 3.3 depicts the DTMC of a Markov-based mobility model in a hexagonal cell topology formed by N cell rings.



**Figure 3.3: DTMC of a Markov-based mobility model with N cell rings.**

The values $p\{stay,i\}$, $p\{next,i\}$, and $p\{back,i\}$ correspond to the probabilities of staying at, moving outwards from, and moving inwards from ring $i$, respectively. According to the cell topology in Figure 3.2, these probabilities can be calculated as:

$$p\{stay,i\} = p + 2p'$$

$$p\{next,i\} = p\{vertex,i\} \cdot 3p' + p\{inner,i\} \cdot 2p'$$

$$p\{back,i\} = p\{vertex,i\} \cdot p' + p\{inner,i\} \cdot 2p'$$

Where:

$$p' = p\{UE\ crossing\ a\ cell\ boundary\} = \frac{1-p}{6}$$

$$p\{vertex, i\} = \frac{1}{i-1}$$

$$p\{inner, i\} = 1 - p\{vertex, i\} = \frac{i-2}{i-1}$$

Once all transition probabilities have been determined, we can calculate the one-step transition probability matrix of the DTMC ($\mathbf{P}_{mob}$). $\mathbf{P}_{mob}$ describes the evolution of the mobility model, as it contains all transition probabilities between DTMC states. We use $\mathbf{P}_{mob}$ to calculate the steady-state probability vector of the DTMC ($\boldsymbol{\pi}_{mob}$), i.e., the vector that contains the probabilities of finding a UE in each one of the cell rings in the TA. It can be proved that:

$$\boldsymbol{\pi}_{mob} = [\pi_{mob}(1) \ ... \ \pi_{mob}(N)] = \boldsymbol{e} \cdot (\boldsymbol{P}_{mob} + \boldsymbol{E} - \boldsymbol{I})^{-1}$$

Where $\pi_{mob}(i)$ is the probability of a UE being in the $i$-th cell ring of the Tracking Area, $e$ is a row vector of all ones, $\boldsymbol{E}$ is a matrix of all ones, and $\boldsymbol{I}$ is the identity matrix.

Once $\boldsymbol{\pi}_{mob}$ is known, we can calculate the probability of initiating a Tracking Area Update procedure in a TA formed by N cell rings as:

$$p\{TAU, N\} = \pi_{mob}(N) \cdot p\{next, N\}$$

Where $\pi_{mob}(N)$ is the probability of being in the outermost ring of a Tracking Area formed by N cell rings and $p\{next, N\}$ is the probability of moving outwards from that ring.

Once $p\{TAU, N\}$ is known, we can calculate the average TAU arrival rate in a Tracking Area formed by N cell rings $\bar{\lambda}_{TAU,N}$). In order to do so, let us consider the TAU Request arrival diagram in Figure 3.4.



**Figure 3.4: TAU Request arrival diagram.**

The average number of TAU arrivals in $k$ timeslots can be calculated as:

$$\bar{N}_{TAU,k} = \sum_{i=1}^{k} p\{TAU, N\} = k \cdot p\{TAU, N\}$$

If the timeslot duration is $\Delta$ seconds, the average TAU arrival rate in a Tracking Area formed by N cell rings is:

$$\bar{\lambda}_{TAU,N} = \lim_{k \to \infty} \left( \frac{\bar{N}_{TAU,k}}{k \cdot \Delta} \right) = \frac{p\{TAU, N\}}{\Delta} \ [arrivals/s]$$

### 3.2.4.2  Self-Organized Tracking Area List Mechanism

The self-organized TAL mechanism is built on top of the standard 3GPP TAU procedure. This is done to comply with 3GPP Technical Specifications. First, the MME monitors the arrival rate of TAU Request messages from the UE in order to determine its mobility state. Secondly, the MME updates the UE-specific Tracking Area List by increasing, keeping, or reducing the number of rings in the TAL according to the mobility state and the paging arrival rate. Finally, the MME sends the new TAL to the UE in the TAU Accept message.

The self-organized mechanism combines static and dynamic TAL management depending on the UE mobility state. Thus, TALs are kept static until the location signalling traffic reaches a certain threshold. Past this activation point, the MME enables dynamic TAL management in order to reduce the overall location signalling traffic in the network. To further understand this behaviour, the concept of stages must be introduced. In stage 1, UEs are registered to Tracking Area Lists formed by a single cell ring. In stage 2, UEs are registered to TALs formed by two cell rings. In stage 3, TALs are formed by 3 cell rings, and so on.

In order to design an analytical model for the self-organized mechanism, some assumptions must be made. First, we assume that the NoF is a hexagonal cell structure formed by concentric rings, as described in Figure 3.2. Secondly, we assume that all TAs in the NoF are formed by a single femtocell. This allows

MMEs to treat femtocells as TAs when managing Tracking Area Lists. Finally, we assume that UEs generate TAU Request messages according to a Poisson process with rate $\bar{\lambda}_{TAU,i}$, where $i$ denotes the number of cell rings in the $i$-th stage of the self-organized TAL mechanism.

On the MME side, two mobility management timers (T1, T2) have been introduced in order to determine the mobility state of the UE. Timers are a common mechanism in 3GPP systems. Both UEs and MMEs use them to trigger signalling procedures, manage transitions between Radio Resource Control (RRC) states, monitor UE activity, release voice and data connections, control authentication protocols, etc. Some examples of 3GPP mobility management timers are T3412 (to trigger periodic TAU procedures), T3311 (to restart the attach procedure with the network) or T303 (to clear a voice call) [73].

In each stage of the self-organized TAL mechanism, T1 and T2 are initialized to different values, namely *T1(i)* and *T2(i)*, where $i$ denotes the stage. Thus, *T1(i)* controls transitions from the current to the next stage, while *T2(i)* controls transitions from the current to the previous stage. All possible transition events (*next*, *stay*, *previous*) are illustrated in Figure 3.5. The stage transition algorithm is described below:



**Figure 3.5: Transition events in the self-organized TAL mechanism.**

- If the MME receives a TAU Request message before *T1(i)* expires, the mechanism transits to the next stage. This corresponds to a UE that is moving too fast for the TAL size in the $i$-th stage.

- If the MME receives a TAU Request message after *T1(i)* expires and before *T2(i)* expires, the mechanism remains in the current stage. This corresponds to a UE that is moving at well-suited speed for the TAL size in the $i$-th stage.

- If the MME receives a TAU Request message after *T2(i)* expires, the mechanism transits to the previous stage. This corresponds to a UE moving too slowly for the TAL size in the $i$-th stage.

After executing the stage transition algorithm, the MME updates the corresponding Tracking Area List and sends it back to the UE encapsulated in the TAU Accept message, as described in [7].

The self-organized TAL mechanism can be mathematically described by a DTMC. States in the DTMC correspond to the stages in the mechanism. Analogously, transitions between states are associated with the probability of staying in the current stage or moving (forward or backwards) to an adjacent stage. Figure 3.6: DTMC of a self-organized TAL mechanism with N stages. depicts the DTMC of a self-organized TAL mechanism with N stages.



**Figure 3.6: DTMC of a self-organized TAL mechanism with N stages.**

In order to characterize the DTMC, we need to calculate *p{next,i}*, *p{previous,i}*, *p{stay,i}* for all stages in the system. Since TAU Request arrivals in each stage follow a Poisson process with arrival rate $\bar{\lambda}_{TAU,i}$, the transition probabilities can be calculated as:

$$p\{next, i\} = 1 - p\{no\ TAU\ arrivals\ during\ T1(i)\} == 1 - e^{-\lambda_{TAU,i} \cdot T1(i)}$$

$$p\{previous, i\} = p\{no\ TAU\ arrivals\ during\ T2(i)\} == e^{-\lambda_{TAU,i} \cdot T2(i)}$$

$$p\{stay, i\} = 1 - p\{next, i\} - p\{previous, i\}$$

Once all transition probabilities have been determined, we can calculate the one-step transition probability matrix of the DTMC ($\mathbf{P}_{self}$). As shown in Section 3.2.4.1, the steady-state probability vector ($\boldsymbol{\pi}_{self}$) can be calculated as:

$$\boldsymbol{\pi}_{self} = \begin{bmatrix} \pi_{self}(1) & \dots & \pi_{self}(N) \end{bmatrix} = e \cdot \left( \boldsymbol{P}_{self} + \boldsymbol{E} - \boldsymbol{I} \right)^{-1}$$

The self-organized TAL mechanism is fully characterized by $\boldsymbol{\pi}_{self}$, as it contains the probabilities of finding the system in each one of the N stages.

### 3.2.5 Performance Evaluation

In this section we evaluate the performance of the self-organized TAL mechanism against that of a conventional (static) TAL mechanism. For both schemes, we define the following *signalling cost function* per UE:

$$C_{tot} = p\{paging\} \cdot \bar{N}_{cells} \cdot c_p + p\{TAU\} \cdot c_{tau}$$

Where $p\{paging\}$ is the probability of a paging arrival in a timeslot, $\bar{N}_{cells}$ is the average number of cells in the Tracking Area List (static or dynamic) at a given instant, $c_p$ is the signalling cost of a single paging operation, $p\{TAU\}$ is the probability of a TAU arrival in a timeslot, and $c_{tau}$ is the signalling cost of a single TAU operation. In cellular networks, a TAU operation generates significantly more signalling traffic than a paging operation. In this paper we have assumed a signalling ratio $\alpha = \frac{c_{tau}}{c_p} = 10$,, which is a common value in the literature [74], [75], [76]. By normalizing the expression of $C_{tot}$ to $c_p$ we obtain the *normalized signalling cost function*:

$$C_{norm} = \frac{C_{tot}}{c_p} = p\{paging\} \cdot \bar{N}_{cells} + \alpha \cdot p\{TAU\}$$

We want to evaluate $C_{norm}$ as a function of the UE speed for both the conventional and self-organized mechanisms. As described in 3.2.4.1, speed variations can be modelled by modifying the value of *p* in the Markov-based mobility model. Thus, low speeds correspond to values of *p* closer to 1, while high speeds correspond to values of *p* closer to 0.

The performance of the self-organized TAL mechanism depends on the values of T1 and T2 in each stage of the system. In order to find the timer values that minimize the overall location signalling traffic, we have used a sequential quadratic programming solver (SQP). The output of the SQP is a pair of vectors:

$$\mathbf{T1}_{opt} = \begin{bmatrix} T1_{opt}(1) & T1_{opt}(2) & \dots & T1_{opt}(N) \end{bmatrix}$$

$$\mathbf{T2}_{opt} = \begin{bmatrix} T2_{opt}(1) & T2_{opt}(2) & \dots & T2_{opt}(N) \end{bmatrix}$$

Where $\mathbf{T1}_{opt}$ and $\mathbf{T2}_{opt}$ contain the values of timers T1 and T2 that minimize $C_{norm}$ in each stage of the self-organized mechanism.

Numerical values for paging arrival rates ($\lambda_p$) and signalling ratio ($\alpha$) have been taken from the literature [74], [75], [76]. The static TAL size in the conventional mechanism ($N_{stat}$) has been set to 2 cell rings. Analogously, the maximum number of rings in a dynamic TAL ($N_{dyn}$) has been set to 4. The maximum femtocell transmission radius (*r*) is 200 m [77]. We have assumed pedestrian and vehicular users in an urban NoF scenario. Therefore, UEs travel at 0-50 km/h. Finally, the timeslot duration can be derived as:

$$\Delta = \frac{2r}{v_{max}}$$

Figure 3.7 summarizes the numerical assumptions considered in the performance evaluation of the self-organized TAL mechanism.

| Parameter | Description | Values |
|-----------|-------------|--------|
| $\lambda_p$ | Paging arrival rate [*pagings/h*] | [0.1 – 10] |
| $\alpha$ | Signalling ratio | 10 |
| $N_{stat}$ | Number of rings in a static TAL | 2 |
| $N_{dyn}$ | Maximum number of rings in a dynamic TAL | 4 |
| $v$ | UE speed [*km/h*] | 0 - 50 |
| $r$ | Femtocell transmission radius [*m*] | 200 |
| $\Delta$ | Timeslot duration [*s*] | 28.8 |

**Figure 3.7: Numerical assumptions for performance evaluation.**

Figure 3.6 shows the impact of UE speed on the normalized signalling cost function ($C_{norm}$) for different paging arrival rates. In general, dynamic Tracking Area Lists generate less location signalling traffic than static TALs for medium- to high-speed UEs. This reduction is significantly higher when UEs are subject to moderate paging. Since the cost of a single TAU operation is tenfold that of a paging operation, the self-organized TAL mechanism aims at minimizing $C_{norm}$ by reducing the probability of TAU arrival for each UE.



**Figure 3.8: Normalized signalling cost function vs. probability of UE staying in current cell for different paging arrival rates.**

The intersections of the two curves in each figure determine the activation points of the self-organized mechanism. Thus, at speeds where static TALs generate less location signalling traffic than dynamic TALs, the self-organized mechanism keeps the TAL size constant. Once the activation point has been reached, the MME enables dynamic TAL management, hence reducing the overall location signalling traffic in the network. This switching strategy yields a significant reduction in location signalling traffic per UE, as shown in Figure 3.7.

| $\lambda_p$ | Reduction |
|---|---|
| 0.1 | 39.45% |
| 1 | 33.53% |
| 5 | 13.21% |
| 10 | 4.45% |

**Figure 3.9: Reduction In Location Signalling Traffic per UE.**

In commercial deployments, mobile network operators design static TAL layouts according to internal network planning criteria. The number of cells per TA and TAs per TAL depends on multiple parameters, such as user density, traffic patterns, UE mobility, etc. Figure 7 compares the performance of dynamic TAL management against that of static TAL management for different TAL sizes in a specific network scenario (other scenarios show a similar generic behaviour). The arrows in the figure correspond to the activation points of the self-organized TAL mechanism. Depending on the scenario and the network planning decisions, the use of dynamic TAL sizes may benefit UEs for different speed ranges. In any case, our mechanism is designed to improve the performance of the static (planned) layout by enabling dynamic TAL management if doing so reduces location signalling traffic.



**Figure 3.10: Normalized signalling cost function for different static TAL sizes.**

### 3.2.6 Conclusions

In these sections we have proposed a self-organized Tracking Area List mechanism that adapts the size of UE-specific TALs to the mobility state and the paging arrival rate of each terminal. Our scheme is particularly suitable for large-scale networks of femtocells, where handovers and cell reselections happen more frequently than in macrocell deployments. In addition, the self-organized TAL mechanism is fully compliant with 3GPP Technical Specifications, which facilitates its implementation in a commercial scenario. We have proposed an analytical model based on discrete-time Markov chains to evaluate the performance of the proposed mechanism against that of the conventional (static) TAL mechanism. The model shows how the self-organized mechanism improves the performance of the conventional mechanism in terms of location signalling traffic for different paging arrival rates. Furthermore, analytical results show that the proposed mechanism can generate up to a 39% less location signalling traffic per UE than the conventional mechanism.

## 3.3 Mobility Management for Networked Femtocells Based on X2 Traffic Forwarding

### 3.3.1 Introduction

The small coverage and massive deployment of femtocells provide new challenges for mobility management, especially for inter-femto Handover (HO) scenario, which can be often found in shopping mall, high street stores and etc. Currently, 3GPP adopts a scheme for inter-femto HO similar to the

scheme used for inter-macro HO [67]. After a mobile establishes the radio connection with the target cell, the target cell will inform the core network entities to switch the data path. This scheme performs well for handover in conventional macrocell networks. Since the coverage of a macrocell is normally up to several kilometers, it is of little chance that several HOs will occur during one session. The data path switch after each HO can reduce the transmission latency over the data path. However, when a mobile moves between femtocells, HOs will be much more frequent than macrocell scenarios due to the small coverage of each femtocell. Further considering the large scale deployment of femtocells, the frequent data path switch operations will cause significant signalling load to the core network entities. Furthermore, different from macrocells, the signalling traffic has to go through the Internet backhaul resulting in more HO latency.

In this activity, we propose novel local mobility management schemes for networked femtocells based on traffic forwarding. The target femtocell can use the local path for ongoing sessions without requiring switching the data path from the core network for each HO. A traffic forwarding chain will be established from the original local anchor point to the current serving femtocell. Since the local traffic forwarding may increase the end-to-end communication latency and consume the local resource, a threshold of the forwarding chain should be defined to balance the trade-off between the path switch cost and traffic forwarding cost. The first scheme, namely *Traffic Forwarding with Cascading Path* (TF_CP), cascades the target femtocell to the previous source femtocell via the local path after a handover. The second scheme, namely *Traffic Forwarding with Shortest Path* (TF_SP), implements local path switch if the target femtocell has a shorter path to the original local anchor point than the cascading path. We develop analytical models based on Markov chains and conduct simulation studies for model validation and performance comparison.

### 3.3.2 Traffic Forwarding with Cascading Path

In this proposed scheme, the target HeNB does not have to send a *Path Switch Request* to the EPC each time a HO occurs. When an UE moves across the boundary of the covering areas of two neighbouring HeNBs, it disassociates with the source HeNB and associates with the target HeNB. After the synchronization between the UE and the target HeNB is completed, the target HeNB will not send the *Path Switch Request* to the EPC as long as the length of the forwarding chain does not exceed a predefined threshold $K$. Note that other criteria may be considered to determine whether or not to trigger traffic forwarding such as the end-to-end latency. For simplicity, the forwarding chain length in terms of hops is considered in this work since this length is closely related to the end-to-end latency and total resource required for local traffic forwarding. The target HeNB will still send the *UE Context Release* message to the source HeNB to inform success of HO and trigger radio and control-plane resource release. However, the resource for data forwarding at the source HeNB will remain reserved for the UE and the data received from the EPC will be forwarded to the target HeNB along the forwarding chain. If the threshold of the forwarding chain is exceeded, the normal data path switch operation will be applied. The EPC will switch the data path from the head of the forwarding chain to the target HeNB and send an "end marker" along the old path until the tail of the forwarding chain. The resource for data forwarding on a forwarding node will be released after the "end marker" is received. After the data path switch operation, the target HeNB becomes the new local traffic anchor point, i.e. the head of a new forwarding chain if any. Figure 3.11 shows the data paths of the TF_CP scheme during mobility, where the threshold of the forwarding chain $K = 2$.

The target HeNB may have been already on the forwarding chain. To remove the loop, a *Forwarding List* including all the node identities on the current forwarding chain and their orders is included in the *SN Status Transfer* message and sent from the source HeNB to the target HeNB during the HO execution phase. The target HeNB can check the forwarding chain status for the UE. If the target HeNB finds that it has been already on the forwarding chain, it will reset the forwarding chain length and send an "end maker" along the rest of the old forwarding chain to release the resource for data forwarding. In case that the UE moves out from a femtocell and moves into a macrocell, the HO signalling has to reach the MME and the S-GW is the mobility anchor point. Therefore, the forwarding chain will be reset to zero and the data path at the S-GW will be switched from the HeNB to the eNB. For the proposed traffic forwarding scheme, the mobility is transparent to the EPC, i.e. the EPC thinks that the head of the forwarding chain is the current serving cell for the UE. The proposed scheme is also transparent to the UE. Therefore, the proposed scheme can be easily fit into the current 3GPP scope, since no upgrade is required from the EPC side or the UE side.

Due to the self-deployment nature of the HeNBs, a HeNB may be switched off or failed when it is on a forwarding chain. Therefore, a mechanism is needed to fast resume the sessions in this case. In the TF_CP scheme, the next HeNB on the forwarding chain will detect the failure of its neighbour (The implementation of the failure detection is left to the manufacturer's discretion.) and send a Path Switch

Request message to the core network entity. The core network entity will switch the data path to the next HeNB on the forwarding chain after the failed one. The forwarding list after the recovery will be updated. During this process, the data packets sent along the old path may be lost and the upper layer mechanisms will be responsible for the packet loss recovery if needed.



**Figure 3.11: Data paths of TF_CP scheme**

### 3.3.3 Traffic Forwarding with Shortest Path

In the TF_CP scheme, the forwarding chain is formed by simply cascading the HeNBs along the trajectory of the UEs. In many cases, the UEs may move around the local anchor point in its surrounding area resulting in the possibility of finding a shorter path compared to the original forwarding chain. In the TF SP scheme, the HeNB will advertize its neighbor list within the local network and thus, each HeNB will possess the network topology information and calculate the per-pair shortest path. In HO completion phase, the target HeNB will compare the forwarding list with the shortest path to the local anchor point. If the length of the shortest path is less than the length of the forwarding chain, the target HeNB will initiate a local path switch operation. A *Local Path Switch Request* will be sent to the local anchor point along the shortest path. If the request is admitted, the local anchor point will respond with a *Local Path Switch Request Ack*. During the message exchange process, the resource required on the new path will be requested and allocated. The data path will be now along the shortest path from the local anchor point to the target HeNB. An "end marker" will be send from the local anchor point along the old forwarding chain to release the resource for data forwarding. And during the local path switch process, the current forwarding chain will continue forwarding the data to the target HeNB until the "end marker" is received, and thus, there is no data loss during this process. An example is shown in Figure 3.12 to illustrate the data path of the TF SP scheme, where a UE is moving along the trajectory HeNB1 →HeNB2 → HeNB3 →HeNB4 and HeNB1 is the current local traffic anchor point. In the TF_CP scheme, a local path is formed following the trajectory of the UE. For example, when the UE arrives at HeNB3, the end-to-end data path from the EPC should be (1)+(2)+(3) as shown in Figure 3.11, thus resulting in a forwarding chain of length 2. On the other hand, the TF_SP scheme will implement a local path switch operation resulting in a forwarding chain of length 1.

In case of a HeNB being switched off or failed, the next HeNB on the forwarding list can detect the failure of its neighbour and calculate a new shortest path to the precedent HeNB of the failed one on the forwarding chain. A local recovery mechanism is implemented by sending a *Local Path Switch Request* message to the precedent HeNB of the failed one on the forwarding chain along the new shortest path. In addition, all the neighbour HeNBs of the failed one will advertise the failure information within the local network such that the corresponding affected per-pair shortest paths will be recalculated. The total length

of the forwarding chain after the local recovery may exceed the threshold. However, the current serving HeNB, i.e. the tail of the forwarding chain, will initiate the local path switch or the core path switch after the network topology and per-pair shortest paths are updated. The threshold is allowed to be temporarily relaxed since the session continuity has the first priority.



**Figure 3.12: Data paths of TF_SP scheme.**

### 3.3.4 Analytical Model

In this section, the performance of the TF_SP scheme and the standard 3GPP scheme are analysed for a grid network topology. It is difficult to model the loop removal of the TF_CP scheme since it depends on the users' trajectory history. Therefore, it will be studied in the next section via simulations. However, the recovery mechanisms of both proposed schemes will be analytically compared in this section. For clarification, Table 3-1 summarizes the parameters used in this section.

**Table 3-1: Parameter used in analysis**

| Parameter | Notation |
|---|---|
| $\tau$ | Time slot duration of Markov chain |
| $\lambda$ | Session arrival rate |
| $\mu$ | Session departure rate |
| $m$ | UE mobility rate |
| $K$ | Threshold of forwarding chain |
| $D_{X2}$ | Transmission latency over X2 |
| $D_{S1}$ | Transmission latency over S1 |
| $D_{HeNB}^c$ | C-plane processing latency at HeNB |
| $D_{HeNB}^u$ | U-plane processing latency at HeNB |
| $D_{HeNB\,GW}$ | UE context retrieval and processing latency at HeNB GW |
| $D_{detect}$ | HeNB failure detection latency |
| $r_p$ | Packet arrival rate during a session |

In the grid network topology considered as shown in Figure 3.13, each HeNB has four neighbours and a UE can move randomly from the current cell to one of its four neighbours with equal probability. The grid-like femtocell deployment has been widely used in the femtocell-related performance analysis, which can be used to model the office or terraced house environment. Let $S_0^{(1)}$ denote the local anchor point and $S_i^{(j)}$ ($1 \le i \le K, 1 \le j \le 4i$) represent the cell location at which the HeNB has the shortest path of length $i$ (ring $i$) towards the local anchor point, where $K$ is the threshold of the forwarding chain. A location state aggregation method similar to [78] is used here to model the UE mobility. The original cell

location states that the UE can reach in Figure 3.13 can be aggregated according to the grid symmetry as shown in Figure 3.14. When a UE moves out of the coverage of the local anchor point $S_0^{(1)}$, no matter the direction of the UE movement, the forwarding chain will be increased by 1. Therefore, the states $S_1^{(1)}$, $S_1^{(2)}$, $S_1^{(3)}$, $S_1^{(4)}$ can be aggregated to a single state $S_1^{(1)}$. When $i = 2$, the states $S_2^{(1)}$, $S_2^{(3)}$, $S_2^{(5)}$, $S_2^{(7)}$ can be aggregated to a state $S_2^{(1)}$ since the UE at these states has ¾ probability to increase its forwarding chain and 1/4 probability to decrease it, while the states $S_2^{(2)}$, $S_2^{(4)}$, $S_2^{(6)}$, $S_2^{(8)}$ can be aggregated to another state $S_2^{(2)}$ since the UE at these states has 1/2 probability to increase its forwarding chain and 1/2 probability to decrease it. Similarly, when $i = K$, the $4K$ states can be aggregated to the $\left\lceil \dfrac{K+1}{2} \right\rceil$ states, where $\lceil x \rceil$ denotes the smallest integer not less than $x$.

The evolution of a UE's activity is modeled as a stochastic process that occurs in a sequence of discrete steps. The duration of a time slot is $\tau$ time unit. It is assumed that the UE state can only change at the end of each time slot and only one change is allowed at a time. Clearly, the system evolution can reflect the real-life characteristics when the time slot is sufficiently small. Based on this assumption, a Discrete-Time Markov Chain (DTMC) model is developed to model the evolution of the UE state when the traffic forwarding scheme is used, as shown in Figure 3.15. The state Sidle denotes the state at which a UE has no ongoing sessions. The state $S_i^{(j)}$ ($0 \le i \le K, 1 \le j \le \left\lceil \dfrac{i+1}{2} \right\rceil$) represents the aggregated state at which a UE with ongoing sessions has a forwarding chain of length $i$. When a UE is at the state Sidle, the probability of a session arriving (incoming or outgoing) during a time slot $\tau$ is $P_\lambda$. Thus, the current HeNB becomes the local anchor point and the UE moves to the state $S_0^{(1)}$ with this probability and stays at the idle state with the probability $1 - P_\lambda$, at the end of a time slot. When a session is initiated at the UE, the probability of a session departing during a time slot $\tau$ is $P_\mu$. Thus, at each state $S_i^{(j)}$, the UE may move to the state $S_{idle}$ with this probability at the end of a time slot. During a time slot, there is a probability $P_m$ that the UE will move to another cell. If the session is still ongoing and the UE decides to move to another cell, the length of the forwarding chain will change. When the UE is at the state $S_K^{(j)}$, the further forward movement will trigger the forwarding chain to be reset.



**Figure 3.13: Grid femtocell network.**



**Figure 3.14: Grid femtocell network after state aggregation**

**Figure 3.15: State transition diagram of the TF_SP scheme.**

Let $\pi_{idle}$ and $\pi_i^j$ ( $0 \le i \le K, 1 \le j \le \left\lceil \dfrac{i+1}{2} \right\rceil$ ) denote the stationary probability distribution of the UE being at the state $S_{idle}$ and $S_i^{(j)}$, respectively. Based on Figure 3.15, the balance equations can be derived and the stationary probability distribution of the Markov chain can be solved. In the following, we derive the parameters used in the Markov chain model: the probability of a session arriving, the probability of a session departing, and the probability of a UE moving to another cell during a time slot. Assuming that the sessions arrive as a Poisson process with rate $\lambda$, the duration of a session has an exponential distribution with mean $1/\mu$, and the cell residence time also has an exponential distribution with mean $1/m$ ($m$ is also called mobility rate), we obtain:

$$P_\lambda = \lambda\tau,$$

$$P_\mu = \mu\tau,$$

$$P_m = m\tau.$$

Since the proposed scheme shares the same radio access and session initiation/termination procedures with the standard 3GPP scheme, they are not taken into account here for comparison purpose. Therefore, the considered cost here is the signalling cost during HO and the data delivery cost. The cost is calculated based on processing latency and transmission latency. To show the performance improvement compared to the 3GPP scheme even when a HeNB GW is deployed, the HeNB GW is assumed to be the mobility anchor in the mobile core network. And thus, the data path will be switched at the HeNB GW without involving the MME and S-GW processing. Let $D_{HeNB}^c$ and $D_{HeNB\_GW}$ denote the signalling processing latency at HeNB and HeNB GW respectively. Let $D_{X2}$ and $D_{S1}$ denote the transmission latency over the X2 and S1 interface respectively. If a UE moves forward to a HeNB, the threshold of the forwarding chain is not reached, and the forwarding chain is one of the shortest paths to the local mobility anchor, there is no path switch during the HO procedure and the forwarding chain will be cascaded. The signalling cost in this case is:

$$C_{nps} = 4(D_{X2} + D_{HeNB}^c),$$

which only accounts for the transmission latency of HO Request, HO Request Ack, SN Status Transfer, and UE Context Release message over the X2 interface and their processing latency at HeNBs. If a UE moves backward to the previous HeNB, the target HeNB will find that it is already on the forwarding chain. It will send an "end marker" along the old forwarding chain back to itself and buffer the data packets received from the precedent HeNB on the forwarding chain. After it receives the "end marker", it can transmit the buffered data packets to the UE. The signalling cost in this case is:

$$C_{bps} = 4(D_{X2} + D^c_{HeNB}) + 2(D_{X2} + D^c_{HeNB}).$$

Compared to $C_{nps}$, the extra cost is needed for the "end marker" transmission and processing. If the threshold of the forwarding chain is not reached but the target HeNB finds that there is a shorter path to the local mobility anchor, the target HeNB will initiate a local path switch procedure. Given the considered grid network topology, this situation will happen when the source HeNB is on ring $i$ while the target HeNB is on ring $i-1$. Thus, the signalling cost is given as:

$$C^i_{lps} = 4(D_{X2} + D^c_{HeNB}) + 2(i - 1)(D_{X2} + D^c_{HeNB}) + i(D_{X2} + D^c_{HeNB}),$$

where the second term accounts for the message transmission and processing of the *Local Path Switch Request* and *Local Path Switch Request Ack* along the shortest path, and the third term accounts for the "end marker" transmission and processing along the old forwarding chain back to the target HeNB. When the threshold of the forwarding chain is exceeded, the target HeNB will trigger a core path switch procedure similar to the standard 3GPP path switch operation with the signalling cost:

$$C_{cps} = 4(D_{X2} + D^c_{HeNB}) + 2D_{S1} + D_{HeNB\,GW}$$
$$+ D_{S1} + D^c_{HeNB} + (K + 1)(D_{X2} + D^c_{HeNB}),$$

Where $2D_{S1} + D_{HeNB\_GW}$ denotes the path switch signalling over the S1 interface and the processing latency at the HeNB GW, $D_{S1} + D^c_{HeNB}$ denotes the "end marker" transmission and processing in a standard 3GPP path switch operation, and $(k+1)(D_{X2} + D^c_{HeNB})$ denotes the "end marker" transmission and processing along the local forwarding path back to the target HeNB. Thus, the total signalling cost per time slot can be expressed as:

$$C^c_{total} = (\pi^{(1)}_0(1 - P_\mu)P_m C_{nps} + \sum_{i=1}^{K-1} \pi^{(1)}_i(1 - P_\mu)P_m(\frac{3}{4}C_{nps} + \frac{1}{4}C_{bps})$$

$$+ \sum_{i=2}^{K-1} \sum_{j=2}^{\lceil \frac{i+1}{2} \rceil} \pi^{(j)}_i(1 - P_\mu)P_m(\frac{1}{2}C_{nps} + \frac{1}{4}C_{bps} + \frac{1}{4}C^i_{lps})$$

$$+ \pi^{(1)}_K(1 - P_\mu)P_m(\frac{1}{4}C_{bps} + \frac{3}{4}C_{cps})$$

$$+ \sum_{j=2}^{\lceil \frac{i+1}{2} \rceil} \pi^{(j)}_K(1 - P_\mu)P_m(\frac{1}{4}C_{bps} + \frac{1}{4}C^K_{lps} + \frac{1}{2}C_{cps}).$$

Let $D^u_{HeNB}$ denote the U-plane processing latency at HeNB. The expected total data delivery cost can be expressed as:

$$C^u_{total} = \sum_{i=0}^{K} \sum_{j=1}^{\lceil \frac{i+1}{2} \rceil} \pi^{(j)}_i(D_{S1} + i(D_{X2} + D^u_{HeNB})).$$

Similarly, the total signalling cost and data delivery cost of the standard 3GPP procedure can be derived by setting the forwarding chain threshold to be null.

Finally, we give an analysis for the comparison of the recovery mechanisms of the TF_CP and TF_SP scheme. The metrics considered here are the recovery latency and the packet loss during the recovery process. In the TF_CP scheme, when a HeNB on the forwarding chain is switched off or failed, the next HeNB on the forwarding chain can detect the failure of the precedent HeNB after a duration $D_{detect}$.

Then the next HeNB will send a *Path Switch Request* to the mobile core network to initiate a core path switch operation. Thus, the recovery latency of the TF_CP scheme is given as:

$$D_{recovery}^{TF\_CP} = D_{detect} + 2D_{S1} + D_{HeNB\,GW}.$$

Let $r_p$ denote the packet arrival rate during a session. The total lost packets are denoted as:

$$P_{loss}^{TF\_CP} = r_p(D_{detect} + D_{S1} + D_{HeNB\,GW}).$$

In the TF_SP scheme, after the next HeNB on the forwarding chain detects the failure of its precedent HeNB, it will send a *Local Path Switch Request* to the precedent HeNB of the failed one on the forwarding chain. Considering $S_1^{(1)}$ in Figure 3.14 as the failed HeNB and $S_0^{(1)}$ as the precedent HeNB of the failed one, the forwarding chain section starting from $S_0^{(1)}$ can be either $\{ S_0^{(1)}, S_1^{(1)}, S_2^{(1)} \}$ or $\{ S_0^{(1)}, S_1^{(1)}, S_2^{(2)} \}$ the recovery latency can be expressed as:

$$D_{recovery}^{TF\_SP} = D_{detect} + \frac{\pi_2^{(1)}}{\pi_2^{(1)} + \pi_2^{(2)}} 8(D_{X2} + D_{HeNB}^c)$$
$$+ \frac{\pi_2^{(2)}}{\pi_2^{(1)} + \pi_2^{(2)}} 4(D_{X2} + D_{HeNB}^c).$$

In case that a buffering mechanism is used at the precedent HeNB of the failed one, i.e. unsuccessfully transmitted packets will be buffered for certain period, the packet loss can be ignorable. Otherwise, it can be denoted as:

$$P_{loss}^{TF\_SP} = r_p(D_{detect} + \frac{\pi_2^{(1)}}{\pi_2^{(1)} + \pi_2^{(2)}} 4(D_{X2} + D_{HeNB}^c)$$
$$+ \frac{\pi_2^{(2)}}{\pi_2^{(1)} + \pi_2^{(2)}} 2(D_{X2} + D_{HeNB}^c)).$$

### 3.3.5  Performance Evaluation

In this section, we evaluate the performance of the proposed schemes and the 3GPP scheme via the analytical models and discrete-event simulations. The default parameter setting used in the evaluation is listed in Table 3-2. In the following, we may vary the values of some parameters to show their effects. The 3GPP related parameter values are based on [86]

**Table 3-2: Parameter setting**

| Parameter | Value | Parameter | Value |
|-----------|-------|-----------|-------|
| $\tau$ | 0.001 min | $D_{X2}$ | 5 ms |
| $\lambda$ | 0.01 /min | $D_{S1}$ | 50 ms |
| $\mu$ | 0.1 /min | $D_{HeNB}^c$ | 4 ms |
| $m$ | 1 /min | $D_{HeNB}^u$ | 1 ms |
| $K$ | 2,4,6 | $D_{HeNB\,GW}$ | 15 ms |
| $r_p$ | 60 /sec | $D_{detect}$ | 50 ms |

The Internet backhaul plays an important role in differentiating the femtocells from the traditional macrocells. Figure 3.16 and Figure 3.17 shows the relative cost ratio between the proposed schemes and the 3GPP scheme as a function of $K$ under various backhaul latency $D_{S1}$. The relative cost ratio is used here for better illustration, defined as the ratio of the cost of the proposed schemes to the cost of the 3GPP scheme under the same parameter setting. In terms of the relative signalling cost ratio as shown in Figure 3.16, we have the following observations: Firstly, more signalling cost saving can be achieved when the backhaul cost becomes higher. Although we do not explicitly model the congestion at the core network entity resulting from the large processing requests from HeNBs, the backhaul cost can somewhat reflect this congestion condition. The proposed schemes can significantly reduce the processing load for the EPC to manage a great number of HeNBs. Secondly, the signalling cost of both proposed schemes is reduced

with the initial increase of $K$ since the longer forwarding chain reduces the number of core path switch operations. Up to a certain value of $K$, the TF_SP scheme achieves more signalling cost saving than the TF_CP scheme because the local path switch operations reduce the chance to reach the forwarding chain threshold. With the further increase of $K$, the TF_CP scheme has less signalling cost than the TF_SP scheme because the local path switch operations incur too much local signalling cost to the local mobility anchor. However, the signalling cost saving comes at the cost of the increased data delivery cost as shown in Figure 3.17. With the increase of $K$, both proposed schemes will incur more data delivery cost compared to the 3GPP scheme without local traffic forwarding. But when the Internet backhaul latency is much higher than the local forwarding latency, which is the typical case in a networked femtocell scenario, the local traffic forwarding only introduces marginal extra cost. Opposite to Figure 3.16, the TF_SP scheme has a little higher data delivery cost than the TF_CP scheme up to a certain value of $K$, but has lower data delivery cost than the TF_CP scheme with the further increase of $K$. Therefore, based on the requirements of the $K$ value, signalling and data delivery cost, a better scheme can be selected. The cell residence time is another factor differentiating the femtocells from the traditional macrocells.



**Figure 3.16: Effect of forwarding chain threshold and backhaul latency on relative signalling cost ratio.**

**Figure 3.17: Effect of forwarding chain threshold and backhaul latency on relative data delivery cost ratio.**

Given the small coverage of a femtocell, frequent handover may happen when users move within the local femtocell networks. Figure 3.18 and Figure 3.19 shows the normalized cost of the proposed schemes and the 3GPP scheme as a function of cell residence time $1/m$. The cost shown in the figures is normalized with the maximum cost value of the 3GPP scheme in the figure being 1. It is clearly shown in Figure 3.18 that significant signalling cost saving can be obtained by the proposed schemes when the cell residence time is small, i.e. frequent HO happens. When a user stays connected with a HeNB for a relatively long period, the signalling cost generated by the 3GPP scheme is not that significant. In addition, a long forwarding chain will be preferred under high mobility to reduce the signalling cost. Figure 3.19 shows that the data delivery cost of the proposed schemes will reduce with the increase of the cell residence time since a UE will traverse less HeNBs during its communication session. The data delivery cost of the 3GPP scheme is independent of the cell residence time as the core path will be switched at each HO.

**Figure 3.18:** Effect of cell residence time on normalized signalling cost.



**Figure 3.19:** Effect of cell residence time on normalized data delivery cost.

The impact of the session duration on the normalized cost is shown in Figure 3.20 and Figure 3.21. The session duration will not affect the signalling cost and the data delivery cost of the 3GPP scheme. When the session duration is sufficiently short, the proposed schemes with different forwarding chain thresholds achieve a similar performance for both costs. Because a communication session will be likely to end before a core path switch operation is triggered. As the session duration increases, both the signalling cost and the data delivery cost of the proposed schemes will increase. And the schemes with a higher threshold of the forwarding chain will have less signalling cost but more data delivery cost. For all the above analytical results, the simulations are run under the same setting for a sufficient long time to obtain the statistical simulation results. It is clearly shown that the results derived from the analytical models closely match the simulation results. Therefore, the proposed analytical models accurately capture the behaviours of the studied schemes and can be used to derive the optimal threshold of the forwarding chain for minimizing the handover cost based on the specific requirements and environments.



**Figure 3.20:** Effect of session duration on normalized signalling cost.



**Figure 3.21:** Effect of session duration on normalized data delivery cost.

The recovery latency and lost packets during the recovery process of the proposed schemes are listed in Table 3-3. It is assumed that there is no buffering mechanism used at the HeNBs in the TF_SP scheme. It is shown that the TF SP scheme using local recovery has less recovery latency and lost packets than the TF_CP scheme using core recovery. And both schemes can recover the communication session within a short period and only incur limited packet loss. Therefore, the proposed schemes are suitable for the self-deployed femtocells.

**Table 3-3: Comparison of recovery mechanisms.**

| Schemes | Recovery Latency | Lost Packets |
|---------|------------------|--------------|
| *TF_CP* | 165.00 ms | 6.90 |
| *TF_SP* | 100.63 ms | 4.52 |

### 3.3.6  Conclusion

Given the small coverage, massive deployment, long Internet backhaul of femtocells, mobility management in femtocells faces new challenges. We propose two local mobility management schemes based on traffic forwarding for networked femtocells to reduce the interfemto handover cost. Instead of implementing the path switch operation at the EPC for each handover, a local traffic forwarding chain is constructed to reuse the old Internet backhaul path. In this way, the processing load at the EPC can be significantly reduced at the reasonable local data delivery cost distributed over the local network. Analytical and simulation studies show that remarkable signalling cost saving can be achieved compared to the 3GPP scheme, especially when the core path switch cost and the mobility rates are high. In particular, the TF_SP scheme has less signalling cost and a little more data delivery cost than the TF_CP scheme when the threshold of the forwarding chain is small while it has more signalling cost and less data delivery cost when the threshold is large. Based on the specific session and cost requirements, mobility pattern and topology availability, an appropriate threshold and traffic forwarding scheme can be selected. In addition, both schemes can recover the communication session within a short period and only incur limited packet loss in case that a HeNB on the forwarding chain is switched off or failed. As a final remark, the proposed schemes are transparent to the EPC and the UEs and no upgrade is required from either side. The modifications are only needed for the femto base stations. Given the deployment of the femtocells is still at the early stage and the standardization activities are ongoing, the proposed schemes can be easily incorporated into the current standard.

## 3.4  Inbound/Outbound Mobility Optimization

### 3.4.1  Introduction

In the 3GPP LTE system, a UE-assisted hard handover procedure is used. For inbound handover from macrocell to femtocell and outbound handover from femtocell to macrocell, the HO signalling has to go through the Internet backhaul, the latency of which may be considerably large and uncertain compared to the HO signalling latency inside the mobile core network for inter-eNB handover. Furthermore, the femto base stations are normally deployed indoors. The radio signal strength may suffer an abrupt degradation when the UE is crossing the doorway due to the wall penetration loss. Finally, the interference from the neighbouring cells under the co-channel deployment may reduce the decoding probability of the HO signalling messages. The above factors together may cause frequent Radio Link Failures (RLFs) occurring during a HO procedure, i.e. the UE may lose the connection to the current serving cell before a HO is completed. In the current specification, when the UE detects a RLF, it starts the RLF timer (T310). Upon expiration of the RLF timer, the UE will search for a suitable target cell to perform the Radio Resource Control (RRC) connection re-establishment. If the target cell has been prepared by the source cell (i.e. RLF occurs after the source cell sends the HO request but before the HO command is sent out to the UE), the re-establishment is successful and the session communication will resume. If the target cell has not been prepared (The source cell was not able to decode the measurement report from the UE or the prepared target cell is not the one selected by the UE during the RRC connection re-establishment), the UE transitions from RRC_CONNECTED state to RRC_IDLE state and attempts to establish a new connection from the scratch. This will incur a long service interruption time, extra signalling load to the core network and buffered data loss at the source cell. In order to avoid the UE entering RRC_IDLE state, the RLF timer is normally set to be conservatively long to allow enough time for the target cell preparation. Thus, the service interruption time due to RLFs may be relatively long even when the RRC connection re-establishment is successful.

Given the problems described above, our target is to investigate inbound/outbound mobility problem, in particular, the impact of Internet backhaul latency which does not appear in traditional mobility scenarios. The small coverage of femtocells in an indoor environment and the large delay of handover preparation may increase the probability of RLFs during handover. We investigate fast HO failure recovery schemes to minimize the service interruption time without compromising the resource utilization, while the impact on the current 3GPP standard should be minimal to facilitate practical system upgrade.

### 3.4.2 Work during Year 1

In year 1, we studied the LTE HO procedure and the RLF problems during inbound/outbound HO. In particular, we proposed a UE-based forward handover procedure with predictive context transfer for fast RLF recovery. Figure 3.22 shows the message sequence chart of forward handover with predictive context transfer for fast RLF recovery. In case of a RLF occurring before the HO command is successfully decoded by the UE, the UE can identify the latest measurement report successfully received by the source cell by examining the ARQ information. From the UE's side, the standard procedure defined in 3GPP for RLF recovery is applied but without the long T310 timer. After a RLF detection latency $T_{U\_det}$, the UE will initiate the RRC connection re-establishment procedure, autonomously select a target cell and contact it via Random Access Channel (RACH) access. The RLF detection is based on the consecutive number of the out-of-sync indications from lower layer. If any in-sync indication is received within $T_{U\_det}$, that means the connection with the source cell is recovered and the normal transmission resumes. If the selected target cell has enough resource to admit the UE, a Random Access (RA) response is sent back to the UE with the information required for connection establishment. After a cell update process, the UE and the target cell are synchronized and ready for data transmission. From the source cell's side, after a similar RLF detection latency $T_{S\_det}$ ( The detailed implementation for RLF detection at base station is left to operators' discretion), the source cell will select $N$ potential target cells based on the latest received measurement report to send the UE's context (downlink data buffered at the source cell can be sent as well depending on the requirements of the services). Meanwhile, a copy of the UE's context is kept at the source cell in case that the target cell selected by the UE does not match the potential target cells predicted by the source cell. After the RRC connection re-establishment is completed at the UE, if its context is already available at the target cell, the data transmission will be resumed immediately. Otherwise, a context fetch request will be sent from the target cell to the source cell to get the UE's context. A timer can be used to avoid unnecessary context fetch request, i.e. the request will be sent if no UE's context is received before the timer expires. The timer mechanisms will also be used to remove the unused context in the source cell and other unused potential target cells.



**Figure 3.22: Message sequence chart of forward handover with predictive context transfer in case of a handover failure.**

### 3.4.3 Interruption Latency Analysis

The U-plane interruption latency is analysed in this section for two cases: (1) A RLF happens before a measurement report is sent to the source cell; (2) A RLF happens after the measurement report is sent to the source cell but before the HO command is received by the UE. In the first case, the target cell is not

aware of the HO and not prepared, while in the second case the target cell has been prepared for the upcoming HO upon receiving the HO request from the source cell.

The detection of a physical layer problem is based on the counter of the out-of-sync indications. If a consecutive number of out-of-sync indications are received from the L1, the T310 timer will be triggered in the 3GPP RLF recovery procedure. The out-of-sync detection latency $T_{det} = N \times T_{ind\_period}$, where $N$ is the counter number and $T_{ind\_period}$ is the period for lower layer measurement. The T310 timer should be given a relatively high value compared to the HO preparation delay. This can enable the target cell has been prepared when the UE tries to connect with the target cell. The value of the T310 timer is generally configured by the operator based on the extensive drive tests. In this work, we assume that $T_{310} = 2 \times T_{backhaul}$, where $T_{backhaul}$ is the expected backhaul latency from the source cell to the target cell. The RRC connection re-establishment latency $T_{RRC}$ consists of radio layer access (DL synchronization +RACH procedure) and RRC signalling. The interruption of the 3GPP RLF recovery procedure for the two cases is respectively denoted as follows:

$$\text{Case (1): } T_{int}^{3GPP} = T_{det} + T_{310} + T_{RRC} + T_{idle-connected}$$

$$\text{Case (2): } T_{int}^{3GPP} = T_{det} + T_{310} + T_{RRC}$$

In case (1), since the target cell is not prepared before a RLF happens, there is no UE context available on the target cell. The RRC connection re-establishment to the target cell will fail and the UE will switch its state from RRC_CONNECTED to RRC IDLE and start the cell selection and RRC connection establishment procedure. The latency for the state switch is denoted by $T_{idle-connected}$, which accounts for the radio layer access and RRC signalling between the UE and the selected cell as well as the connection request/setup procedure between the eNB/HeNB to the MME.

The latency of the proposed scheme for the two cases is respectively denoted as follows:

$$\text{Case (1): } T_{int}^{FH-PCT} = \max(T_{det} + T_{backhaul},\ T_{det} + T_{RRC})$$

$$\text{Case (2): } T_{int}^{FH-PCT} = \max(T_{det} + T_{backhaul},\ T_{det} + T_{RRC})$$

Here the RLF detection latency at both UE and base station is assumed to be same.

According to the latency value of each component given in Table 3-4, the numerical results are given in Figure 3.23. It is shown that the proposed Forward Handover with Predictive Context Transfer (FH_PCT) scheme can significantly reduce the service interruption latency in case of a handover failure, especially when the backhaul latency is long. In addition, the interruption latency of the proposed scheme is same for both cases, i.e. no matter whether the failure happens due to missing measurement report or HO command.

**Table 3-4: Component latency of HO failure recovery**

| Latency | $N$ | $T_{ind\_period}$ | $T_{backhaul}$ | $T_{RRC}$ | $T_{idle-connected}$ |
|---|---|---|---|---|---|
| Value (ms) | 5 | 10 | 10-80 | 50 | 100 |

**Figure 3.23: Interuption latency in case of a handover failure**

## 3.5 Backhaul for Mobile Femtocells

### 3.5.1 Re-analysis of the basic 3GPP fixed relay architecture alternatives for enabling mobile relay

Relay is one of the new key features that has been defined in 3GPP LTE-Advanced (3GPP Rel' 10) to improve the coverage of high data rates, group mobility, temporary network deployment, etc. [84]. Up to now, 3GPP LTE-A have considered only non-mobile (fixed) relaying when the Relay Node (RN) is connected through *wireless backhaul* to a Donor eNB (DeNB) by means of a modified version of the E-UTRA radio interface (the Un interface) and becomes a part of the fixed access network.

In the report 3GPP TR 36.806 [83], it is concluded that the architecture Alternative 2 (proxy S1/X2) has more benefits in compare with other three alternatives and it is selected as the baseline architecture for Rel'10. But, the conclusion was done keeping in mind the fixed relay. In the context of Mobile Relay Node (MRN) changing its *backhaul link* from Source-DeNB to Target-DeNB the other alternatives considered in [83] can have more advantages. For this reason, in the next subsections we re-analyse the 3GPP fixed relay architecture alternatives presented in [83] for enabling MRN.

#### 3.5.1.1 Comparison of the baseline fixed relay architecture Alt 2 (Proxy S1/X2) with the relay architecture Alt 1 (Full-L3 relay) in the context of relay mobility handling

First of all, we compare the baseline relay architecture Alt 2 with the architecture Alt 1 (Full-L3 relay) since they provide two different concepts. If in the architecture Alt 2 the RN S1/X2 interfaces are terminated in the DeNB and the S-GW/P-GW functionality serving the RN is in the DeNB, then in the architecture Alt 1 the RN is transparent for DeNB and the S-GW/P-GW (RN) entity is separated from the DeNB.

*User plane aspects*

User planes of the architecture Alt 2 and the architecture Alt 1 for the fixed relay are presented in the following figure, a) and b) respectively [83].

**Figure 3.24 User planes of the architecture Alt 2 (a) and the architecture Alt 1 (b)**

As was mentioned yet, one of specific features of the architecture Alt 2 is that the S-GW/P-GW functionality related to the RN is located within of the DeNB. It allows avoiding a situation when packets traverse via two S-GW/P-GWs as it happens in the architecture Alt 1 (see item (b) of the above figure). In fact, the S-GW/P-GW (RN) functionality integrated in the DeNB was considered in [83] as an enhancement to optimize the routing path that had redundancy in the Alt 1 (for fixed relay).

Another feature of the Alt 2 is that the home eNB GW functionality is added in into the DeNB that results in the "Proxy S1/X2" architecture. The RN looks like a cell under the DeNB in the architecture [83]. As it is seen from Figure 3.25 (a) illustrating the packet delivery process for the architecture Alt 2, a GTP tunnel per UE bearer is puzzled from two parts. The first part of the GTP tunnel goes from the S-GW/P-GW serving the UE to the DeNB and the second part goes from the DeNB to the RN. In the DeNB one-to-one mapping of these two parts of the UE bearer GTP tunnel is performed.

In the architecture Alt 1 the S-GW/P-GW functionality serving the RN is out of the DeNB, i.e. the U-plane packets are delivered to/from a UE via two GWs, namely, the S-GW/P-GW serving the RN and S-GW/P-GW serving the UE. The packet forwarding path is longer in this architecture that can leads to large latency, especially when a UE under a fixed RN makes a handover to an eNB [83]. This unnecessary back and forth traffic forwarding (as well as signalling exchanges) was one of the main reasons why this architecture was not selected as the baseline one for the fixed relay. But, in the context of mobility handling, the separation of the S-GW/P-GW (RN) from the DeNB can brings some benefits in the sense that the S-GW/P-GW (RN) can be considered as a stable "IP anchor point" to support the mobile RN (MRN). Especially, taking into account that one more specific feature of the architecture Alt1 is transparence of the DeNB for the full L3 RN. That is, a GTP tunnel per UE bearer is entire without any switching in the DeNB and goes from the SGW/PGW serving the UE directly to the RN as it is illustrated in Figure 3.25 (b).

**Figure 3.25: Packet delivery steps, Alt 2 (a) and Alt 1 (b)**

*Control plane aspects*

Control planes of the architecture Alt 2 and the architecture Alt 1 for the S1 interface are presented in Figure 3.26 a) and b) correspondingly [83]. A similar logical structure can be applied for the X2 interface. In this case, the S1-AP interface is changed into X2-AP and the MME serving the UE is substituted into the target eNB.

A specific feature of the architecture Alt 2 in the context of control plane aspects is that the S1-AP messages are sent in two steps: between the MME serving the UE and the DeNB and between the DeNB and the RN. They are encapsulated by SCTP/IP in the MME serving the UE and transferred over an EPS data bearer of the RN where the S-GW/P-GW functionality for the RN's EPS bearers is incorporated into the DeNB [83].

In the architecture Alt 1, the S1-AP signalling messages are sent in one step directly between the MME serving the UE and the RN. The DeNB and the S-GW/P-GW of the RN play the role of user plane transport nodes from the signalling traffic point of view. That is, the S1 signalling messages sent between the RN and the MME serving the UE are mapped on user plane EPS bearers of the RN.

S1 interface relations and signalling connections of the architecture Alt 2 and the architecture Alt 1 are illustrated in Figure 3.27 a) and b) correspondingly [83].

**Figure 3.26 Control planes of architecture alt 2 (a) and alt 1 (b)**



**Figure 3.27: S1 interface relations and signalling connections - Alt 2 (a) and Alt 1 (b)**

In the architecture Alt 2 (Figure 3.27a), there is a single S1 interface relation between the RN and the DeNB, and there is one S1 interface relation between the DeNB and each MME (serving the UEs) in the

MME pool. That is, the RN has to maintain only one S1 interface (to the DeNB), while the DeNB maintains one S1 interface to each MME in the respective MME pool. All S1 signalling connections between RN and MMEs are processed by the DeNB for all UE-dedicated procedures. All non-UE-dedicated S1-AP procedures are terminated at the DeNB, and handled locally between the RN and the DeNB, and between the DeNB and the MME.

In the architecture Alt 1 (Figure 3.27b), there is one S1 signalling relation for each connected UE on the given S1 interface between the RN and the MME serving the UE. That is, the RN has to maintain one S1 interface relation to each MME in the respective MME pool and the S1 interface and the signalling connections go through the DeNB transparently.

### 3.5.1.2 Relay architectures for supporting mobile relay moving from S-DeNB to T-DeNB

Taking into account the user and control plane aspects, the relay architectures for supporting MRN when it moves from S-DeNB to T-DeNB based on the architectures Alt 2 and Alt 1 are presented in Figure 3.28a and Figure 3.28b correspondingly.

a)

b)

**Figure 3.28: Relay architectures for supporting MRN a) based on Alt 2 and b) based on Alt 1**

As was mentioned above, the S-GW/P-GW functions serving the relay in the architecture Alt 2 are embedded in the DeNB. The S-GW/P-GW (RN) entity handles the IP packets flowing to/from the relay and the attached UEs. When mobile Relay Node (MRN) changes its backhaul from S-DeNB to T-DeNB, the user plane IP connectivity that was established between the MRN and S-DeNB is missed. Since the S-GW/P-GW (RN) includes functionality for the IP address allocation of the Un interface for the MRN, the MRN needs to obtain a new IP address and establish the IP connectivity with T-DeNB. It requires to perform again the attach procedure for relay operation to establish the backhaul with T-DeNB, i.e. phase II of the relay attach procedure should be repeated. This phase of the attach procedure (see 3GPP TS 36.300 [85]) includes relay backhaul reconfiguration that is completed by OAM, relay-initiated S1/X2 setup (S1/X2 interface of the MRN is released with S-DeNB and initiated again with T-DeNB), and S1/X2 eNB configuration update. Definitely, the repetition of the attach procedure can bring additional and very essential latency in the handover process.

Besides, there is one more problem because of the IP connectivity is lost when the MRN is detached from the S-DeNB. It is related to interworking between the MRN and its OAM. The relay should send alarms, traffic counter information to the OAM and receive commands, configuration information and software upgrades from the OAM [85]. The reference RN-OAM architecture is presented in Figure 3.29 [85].



**Figure 3.29: RN-OAM architecture**

Thus, when the MRN goes from the S-DeNB to the T-DeNB the connection between the MRN and its OAM is interrupted. But, alarms in the MRN generate bursts of high-priority traffic and have to be transported in real time. Moreover, configuration messages from the OAM to the MRN are delay-sensitive [85]. Thus, the mobile relay architecture based on the architecture Alt 2 is not appropriate to support these requirements.

In the architecture Alt 1, the S-GW/P-GW of the relay is separated from the DeNB. When the MRN moves from S-DeNB to T-DeNB, the S-GW/P-GW of the MRN serves as anchor, i.e. IP connectivity is handled through retaining a stable "IP anchor point" in the network which allows for not having to change the IP address of the MRN at all. As a result, there is no need to perform the network re-entry at T-DeNB, i.e. to re-initiate the phase II of the relay attach procedure. Besides, interworking between the MRN and its OAM is not interrupted. The separated S-GW/P-GW (MRN) can bring one more benefit related to the data path switching (see next subsection for details).

### 3.5.1.3 MRN handover preparation, execution, and completion phases (Alt 2 and Alt 1)

In this subsection we concentrate on the handover procedure when the mobile RN (MRN) moves between S-DeNB and T-DeNB for the architectures Alt. 2 and 1 correspondingly.

Since the relay supports a subset of the UE functionality [85] in order to be wirelessly connected to the DeNB, the legacy UE handover procedure can be re-used as a basis for the MRN handover procedure, but taking into account new features that the Un interface can bring and the specific features of each architecture described before.

Thus, when the MRN acts as a UE it has to support the NAS and the RRC protocols towards the network. The relay has the following relay-specific functionalities related to specific features of the Un interface [85]:

- the RRC layer of the Un interface has functionality to configure and re-configure an MRN subframe configuration through the MRN reconfiguration procedure (e.g. DL subframe configuration and an RN-specific control channel) for transmissions between an MRN and a DeNB. The MRN may request such a configuration from the DeNB during the RRC connection establishment, and the DeNB may initiate the RRC signalling for such configuration. The MRN applies the configuration immediately upon reception;

- the RRC layer of the Un interface has functionality to send updated system information in a dedicated message to an MRN with an MRN subframe configuration. The MRN applies the received system information immediately;
- the PDCP layer of the Un interface has functionality to provide integrity protection for the user plane. The integrity protection is configured per Data Radio Bearer (DRB).

To support Public Warning System (PWS) towards UEs, the MRN receives the relevant information over S1. The MRN should hence ignore DeNB system information relating to PWS.

The MRN handover procedure like the legacy UE handover procedure includes three phases: handover preparation, handover execution, and handover preparation.

Figure 3.30 illustrates the handover procedure when the MRN moves between S-DeNB and T-DeNB for the architecture Alt 2.



**Figure 3.30: MRN handover procedure - Alt 2.**

Below is a description of the MRN handover procedure for the architecture Alt 2. It is based on the legacy UE handover procedure described in [85] taking into account the relay-specific and the architecture-specific aspects.

1) Handover preparation phase:

- The MRN issues a HANDOVER REQUEST message to the S-DeNB passing necessary information to prepare the handover at the target side. In particular, it can contain a list of preferred T-DeNBs among other system and RRM information.
- The S-DeNB makes decision based on the HANDOVER REQUEST message to hand off the MRN.
- The S-DeNB reads the ID of the preferred T-DeNB(s) from the request message, then it finds the T-DeNB(s) and forwards the HANDOVER REQUEST message to the T-DeNB(s) with all necessary information to prepare the HO of the MRN. The selected T-DeNB configures the required resources. It also identifies from the HANDOVER REQUEST message that the handover is required for a MRN.
- The Admission Control procedure is performed by the T-DeNB based on the received E-RAB QoS information to increase the likelihood of a successful handover, if the corresponding resources to support the MRN with the attached UEs can be granted by T-DeNB.
- If handover for the MRN is accepted, the T-DeNB sends a HANDOVER REQUEST ACKNOWLEDGE message to the S-DeNB and the information for the MRN about

reconfiguration of *the relay backhaul* towards the T-DeNB is passed to the S-DeNB to be sent to the MRN as an RRC message to perform the handover procedure.

The handover preparation phase is finished at the T-DeNB. As soon as the S-DeNB receives the HANDOVER REQUEST ACKNOWLEDGE, or as soon as the transmission of the handover command is initiated in the downlink, data forwarding may be initiated.

2) Handover execution phase:

- The T-DeNB generates the RRC message to perform the handover, i.e. RRCConnectionReconfiguration message including the mobilityControlInformation, to be sent by the S-DeNB towards the MRN. The MRN receives the RRCConnectionReconfiguration message for *the relay backhaul* reconfiguration specified for the Un interface. The MRN is commanded by the S-DeNB to perform the handover based on radio resource information of the T-DeND. This RRC information about needed resources should be specified for the Un interface. The MRN is detached from S-DeNB and starts gaining synchronization at the T-DeNB.
- The S-DeNB sends the SN STATUS TRANSFER message to the T-DeNB to convey the uplink PDCP SN receiver status and the downlink PDCP SN transmitter status of E-RABs for which PDCP status preservation applies taking into account the Un interface specific functionality.
- After receiving the RRCConnectionReconfiguration message including the mobilityControlInformation (specified for the Un interface), the MRN completes synchronization and reconfiguration of the *relay backhaul* to T-DeNB. The *backhaul* between the MRN and the T-DeNB is reconfigured based on the obtained information from the RRCConnectionReconfiguration message.
- The T-DeNB performs UL allocation towards the MRN.
- When the *relay backhaul* is successfully established between and the T-DeNB, the MRN sends the RRCConnectionReconfigurationComplete message to the T-DeNB to confirm the handover.
- The MRN configuration is downloaded from the OAM system through the T-DeNB.
- The S1/X2 setup is initiated between the MRN and the T-DeNB.
- T-DeNB updates its configuration on the S1/X2 interfaces.

3) Handover completion phase

- For UEs served by the certain MRN, the RMN sends a PATH SWITCH REQUEST message to the T-DeNB. The T-DeNB forwards the message to the MME serving the UEs to inform that the UEs has changed cell. The MME serving the UEs issues USER PLANE UPDATE REQUEST to the S-GW/P-GW serving the UEs.
- The S-GW/P-GW of the UEs switches DL path and sends packet data to the T-DeNB.
- The S-GW/P-GW of the UEs sends one or more "end marker" packets on the old path to the S-DeNB to release resources related to the U-plane.
- The S-GW/P-GW of the UEs issues a USER PLANE UPDATE RESPONSE to the MME of the UEs.
- The MME of the UEs send a PATH SWITCH REQUEST ACKNOWLEDGE message to the T-DeNB. The T-DeNB forwards the PATH SWITCH REQUEST ACKNOWLEDGE message to the MRN.

  *If UEs attached to the MRN are served by different S-GWs/MMEs the above items of the handover completion phase are completed for each S-GW/MME serving the UEs (see* Figure 3.30*)*

- By sending the MRN CONTEXT RELEASE message, the T-DeNB informs the S-DeNB about success of the handover procedure. It triggers the release of C-plane and radio related resources by the S-DeNB. The T-DeNB sends this message after the PATH SWITCH REQUEST ACKNOWLEDGE messages are received from the MME(s) serving UEs attached to the MRN.
- Upon reception of the MRN CONTEXT RELEASE message, the S-DeNB can release radio and C-plane related resources associated to the MRN context. Any ongoing data forwarding may continue.

Figure 3.31 illustrates the handover procedure when the MRN moves between S-DeNB and T-DeNB for the architecture Alt 1.

**Figure 3.31: MRN handover procedure - Alt 1.**

The handover preparation phase is similar as for the Alternative 2. Handover execution phase is also similar excepting last three items, i.e. it is not needed in this case for the MRN to download configuration from the OAM and initiate the S1/X2 setup and it is not needed for the T-DeNB to update its configuration on the S1/X2 interfaces.

The handover completion phase for the Alternative 1 is different and it is described below.

- The T-DeNB issues a PATH SWITCH REQUEST message to the MME serving the MRN to inform that the MRN has detached from the S-DeNB.
- The MME serving the MRN sends USER PLANE UPDATE REQUEST to the S-GW/P-GW serving the MRN.
- The S-GW/P-GW of the MRN switches DL path and sends packet data to the T-DeNB.
- The S-GW/P-GW of the MRN sends one or more "end marker" packets on the old path to the S-DeNB to release resources related to the U-plane.
- The S-GW/P-GW of the MRN issues a USER PLANE UPDATE RESPONSE to the MME of the MRN.
- The MME of the MRN sends a PATH SWITCH REQUEST ACKNOWLEDGE message to the T-DeNB.
- By sending the MRN CONTEXT RELEASE message, the T-DeNB informs the S-DeNB about success of the handover procedure. It triggers the release of C-plane and radio related resources by the S-DeNB. The T-DeNB sends this message after the PATH SWITCH REQUEST ACKNOWLEDGE is received from the MME of the MRN.
- Upon reception of the MRN CONTEXT RELEASE message, the S-DeNB can release radio and C-plane related resources associated to the MRN context. Any ongoing data forwarding may continue.

Thus, the handover procedure for the MRN moving between S-DeNB and T-DeNB based in the architecture Alt 1 seems more simplified than one based on the architecture Alt 2. It does not include the network re-entry procedure for the MRN at T-DeNB. Moreover, having the separated S-GW/P-GW (MRN) serving as the MRN stable IP anchor point leads to the situation when the data path of the MRN is switched only from T-DeNB to S-DeNB over the S-GW/P-GW of the MRN (under control of the MME of the MRN) handling the user plane for all UEs attached to the MRN. In the case of the architecture Alt

2, the user plane is handled by the S-GW/P-GW of a UE. The UEs attached to the MRN can be connected to different the S-GW/P-GWs (UE) that can be under control of different MMEs (UE). It can lead to more complex and time-consuming procedure of data path switching. As a result, latency of the entire handover procedure is also increased in the Alt 2.

### 3.5.1.4 Comparison of the baseline fixed relay architecture Alt 2 (Proxy S1/X2) with the relay architectures Alt 3 and 4 in the context of relay mobility handling

In accordance with 3GPP TR 36.806, the distinguished feature of the architecture Alt 3 is that RN bearers are terminated in DeNB [83]. The distinguished feature of the architecture Alt 4 is the user plane of the S1 interface is terminated in DeNB. However, in the context of mobile relay handling the alternatives have the same crucial problem as the architecture Alt 2. Namely, the S-GW/P-GW (MRN) functionality serving the RN in the architecture Alt 3 and Alt 4 is embedded into the DeNB. As a sequence of this, there is no "IP anchor point" handling IP connectivity for MRN performing the handover procedure in these alternatives. Thus, they are not not able to handle continuous interworking with OAM, require to repeat the attach procedure for RMN when it connects to T-DeNB, and generate multiple data path switching during handover procedure.

Definitely, some enhancements should be done to support the relay mobility in the alternatives 2,3, and 4. Some basic ideas for these enhancements based on the concept of proxy mobile IP or dynamic mobility anchoring to maintain the OAM signalling in DeNB were mentioned in [90]. But all these improvements entail considerable complexity in the DeNB and lead to essential changes in the current 3GPP specifications (new interfaces should be defined) [90]. The architecture Alt 1 (a full L3 relay) does not require such kind of improvements.

### 3.5.2 Performance evaluation of different relay architecture scenarios (latency analysis)

In this subsection we focus on the quantitative comparison related to the performance of the architecture Alt 2 (Proxy S1/X2) and the architecture Alt 1 (Full-L3 relay) as they provide two different concepts. The architecture Alt 2 is selected as the baseline architecture for the fixed relay in Rel'10 [83]. The architecture Alt 1 is the most appropriate architecture to handle the mobile relay based on our previous observations. Namely, in this subsection we analyse the handover procedure duration in fixed and mobile networks for handover scenarios based on these two alternatives. The first handover scenario is related to a fixed relay network when a UE under an RN moves from the RN to the neighboring DeNB. The second handover scenario is related to a mobile relay network when the MRN moves from S-DeNB to T-DeNB.

**Table 3-5 First handover scenario: Latency values used for evaluating mobile relay alternatives 1 and 2**

| Description | Alt. 1 | Alt. 2 |
|---|---|---|
| One-way transmission from RN to UE (2 ms) and processing in UE (2 ms) (HO command) | 4 ms | |
| Resource Reservation in T-eNB | 5 ms | |
| Synchronization time at UE | 20 ms | |
| One-way transmission from UE to T-eNB (2 ms) and processing in T-eNB (2 ms) (HO complete) | 4 ms | |
| One-way transmission from T-eNB to MME/S-GW serving UE (5 ms) and processing in the MME/S-GW (2 ms) (PS_Req) | 7 ms | |
| One-way transmission from MME/S-GW (UE) to T-eNB (5 ms) and processing in T-eNB (2 ms) (PS_Req_Ack) | 7 ms | |
| One-way transmission from RN to T-eNB and processing in T-eNB (2 ms) (HO_Req) | 20 ms | 13 ms |
| One-way transmission from T-eNB to RN and processing in RN (2 ms) (HO_Req_Ack) | 20 ms | 13 ms |
| One-way transmission from RN to T-eNB and processing in T-eNB (2 ms) (SN Status Transfer) | 20 ms | 13 ms |
| One-way transmission from T-eNB to RN and processing in RN (2 ms) (Context Rls) | 20 ms | 13 ms |
| Total | 127 ms | 99 ms |

The handover procedure of the first scenario for the two alternatives is considered in [83] (Section 4.2.4.3). Based on the feasibility study completed in [86], the delay budget for the UE handover in a fixed relay network was provided in [87]. With reference to the delay budget, values of latency for the scenario (when a UE moves from the fixed RN to the neighboring eNB) are presented in Table 3-5 for the alternatives 1 and 2 correspondingly. It is assumed that one-way transmission time between RN and DeNB is 2 ms is the same as the transmission time between a UE and a RN/eNB [87].

As seen from Table 3-5, the architecture Alt 1 takes about 130 ms to hand over a UE from a fixed RN to a neighboring DeNB, while the architecture Alt 2 takes just under 100 ms. This is because of traffic signalling path in the Alt 1 is longer than in the Alt 2. It goes in the Alt 1 from the RN to T-eNB via D-eNB and SGW/PGW serving the RN while in the Alt 2 only via D-eNB (see the corresponding handover procedures for the two alternatives in [83] (Section 4.2.4.3)).

Values of latency for the second scenario when the mobile relay (with attached UEs) moves from S-DeNB to T-DeNB are presented in Table 3-6 for the alternatives 1 and 2 respectively. Handover preparation, execution, and completion phases of this scenario for the two alternatives were described in detail in Section 3.5.1.3. It is assumed based on information provided in [87], [88], [89] that OAM configuration download for the MRN takes 7 ms, S1/X2 setup initiation between the MRN and the T-DeNB is around 17 ms, the T-DeNB configuration update takes 2 ms.

**Table 3-6 Second handover scenario: Latency values used for evaluating mobile relay alternatives 1 and 2**

| Description | Alt. 1 | Alt. 2 |
|---|---|---|
| One-way transmission from MRN to S-DeNB (2 ms) and processing in S-eNB (2 ms) (HO_Req) | 4 ms | |
| One-way transmission from S-DeNB to T-DeNB (5 ms) and processing in T-DeNB (2 ms) (HO_Req) | 7 ms | |
| Admission Control performed by T-DeNB | 5 ms | |
| SN transfer status from S-DeNB and T-DeNB | 5 ms | |
| One-way transmission from T-DeNB to S-DeNB (5 ms) and processing in S-DeNB (2 ms) (HO_Req_Ack) | 7 ms | |
| Synchronization and reconfiguration of the relay backhaul between MRN and T-DeNB | 20 ms | |
| One-way transmission from T-DeNB to MME serving MRN (5 ms) and processing in the MME (2 ms) (PS_Req) | 7 ms | |
| One-way transmission from MME serving MRN to T-DeNB (5 ms) and processing in the T-DeNB (2 ms) (PS_Req_Ack) | 7 ms | |
| One-way transmission from T-DeNB to S-DeNB (5 ms) and processing in S-DeNB (2 ms) (Context_Rls) | 7 ms | |
| One-way transmission from MME serving MRN to S-GW/P-GW serving the MRN (5 ms) and processing in the S-GW/P-GW (2 ms) (UP_update_Req) | 7 ms | - |
| One-way transmission from S-GW/P-GW (MRN) to MME serving MRN (5 ms) and processing in the MME (2 ms) (UP_update_Rsp) | 7 ms | - |
| The S-GW/P-GW serving MRN switches DL path | 5 ms | - |
| OAM MRN configuration | - | 12 ms |
| MRN-initiated S1/X2 setup | - | 24 ms |
| S1/X2 T-DeNB configuration update | - | 2 ms |
| One-way transmission from MRN to T-DeNB (2 ms) and processing in the T-DeNB (2 ms) (PS_Req) | - | 4 ms |
| One-way transmission from T-DeNB to MRN (2 ms) and processing in the MRN (2 ms) (PS_Req_Ack) | - | 4 ms |

| | | |
|---|:---:|:---:|
| One-way transmission from T-DeNB to MME serving UE (5 ms) and processing in the MME (2 ms) (PS_Req) | - | 7 ms |
| One-way transmission from MME serving UE to T-DeNB (5 ms) and processing in the T-DeNB (2 ms) (PS_Req_Ack) | - | 7 ms |
| One-way transmission from MME serving UE to S-GW/P-GW serving the UE (5 ms) and processing in the S-GW/P-GW (2 ms) (UP_update_Req) | - | 7 ms |
| One-way transmission from S-GW/P-GW (UE) to MME serving UE (5 ms) and processing in the MME (2 ms) (UP_update_Rsp) | - | 7 ms |
| The S-GW/P-GW serving UE switches DL path | - | 5 ms |
| Total | 88 ms | 136 ms |

As seen from Table 3-6, the latency to perform the MRN handover procedure from S-DeNB to T-DeNB for the Alternative 1 is about 1.5 times less than for the Alternative 2. This is because the handover procedure for the Alt 2 requires much more operations when MRN performs the network re-entry at the T-DeNB. Moreover, note that in the case of the Alternative 2, it is supposed in Table 3-6 that all UEs attached to the MRN are under the same S-GW/MME serving the UEs. If they are connected to different S-GWs/MMEs it can lead to more DL data path switching. That is, last six operations (41 ms) in Table 3-6 should be repeated several times depending on the number of different S-GWs/MMEs (UEs). Definitely, it can significantly increase the handover procedure duration for the Alternative 2.

Figure 3.32 illustrates the handover procedure duration in fixed and mobile networks for the handover scenarios based on the alternatives 1 and 2.



Figure 3.32: Handover procedure duration in fixed and mobile networks based on the alternatives 1 and 2.

Thus, the Alt 2 shows much better performance in a fixed relay network whilst the Alt 1 shows much better performance in a mobile relay network.

### 3.5.3 Conclusion on the baseline mobile relay architecture

During the study we have re-analysed the basic 3GPP relay architecture alternatives presented in [83] for enabling mobility of relay node. In particular, mobile relay architectures based on these alternatives were proposed and mobile relay handover procedures for these architectures were described for comparative analysis.

It was observed that for mobile relay changing its backhaul link (Un interface) from Source-DeNB to Target-DeNB the architecture Alternative 1 (a full L3 relay) has more benefits than other architecture architectures.

In particular, in contrast to other alternatives, the architecture Alt 1 supposes the stable IP anchor point (S-GW/P-GW) supporting IP connectivity for MRN. As a sequence, mobile relay architecture based on the Alternative do not require the time-consuming relay re-attach procedure when backhaul link is re-established to T-DeNB and it handles interworking between the MRN and OAM without connectivity interruption that is crucial issue for normal operating conditions of the MRN. Moreover, the MRN handover procedure for this architecture does not require the downlink data path switches related to different S-GWs/MMEs serving UEs. In this sense, group mobility for UEs attached to the MRN is supported when the MRN moves from S-DeNB to T-DeNB.

The latency analysis related to the MRN handover procedure confirms the analytical observations in accordance to which the architecture Alt 1 shows better performance in a mobile relay network.

It is concluded as a result of this study that the architecture Alt1 is more suitable for mobile relay handling then other architecture alternatives and can be selected as baseline architecture in a mobile relay network.

## 3.6 Mobile Femtocells based on Multi-homing Femtocells

### 3.6.1 Introduction

This section discusses an alternative approach to Mobile Relays, which is to turn a normal femtocell into a Mobile Femtocell. The Mobile Femtocell works in an overlay manner by establishing the IP connectivity to its mobile operator's Femtocell Gateway over a mobile cellular access. By extending such overlay solution to multiple IP backhaul links, i.e. by making the femtocell multi-homing capable, several mobile access links can be bundled for higher reliability, quality, and capacity. It is even possible to use mobile accesses of different vendors.

The topic of Mobile Femtocells is approached in two steps:

1. First, we propose an approach for making normal femtocells multi-homing capable. This has applications even if the femtocell is not mobile, as shall be explained in the following.

2. Second, we study the effect of bundling multiple backhauls established via the macro-cellular backhaul into one aggregate link on a Mobile Femtocell's performance.

The remainder of Section 3.6 focuses on making femtocells multi-homing capable, while the second step will be addressed in a later deliverable.

For cases when a single, fixed femtocell backhaul cannot meet the capacity, quality or reliability requirements, it would be desirable to enable femtocells to switch or load balance between multiple backhauls transparently to connected mobile devices. There are multiple use cases for such multi-homing:

- Capacity pooling or fail-over between multiple fixed broadband (e.g. DSL) interfaces or between fixed and mobile access provided by the macro cellular network,

- flow-mobility between fixed broadband and mobile access depending on flows' QoS requirements,

- exploiting a potentially better radio channel and transmit power budget when relaying via the femtocell to the macrocell rather than directly communicating between an indoor device and the macrocell, and

- exploiting LTE macro cell capacity to provide backhauling to 3G femtocells / devices.

All of the above use cases could be implemented with specific support by the mobile devices.

This work proposes modifications to the existing 3GPP femtocell protocol stack [92] to enable transport layer multi-homing for femtocell deployments. Specifically, this work proposes replacing the existing protocol combination of GPRS Tunnelling Protocol User plane (GTP-U) over User Datagram Prototcol (UDP) to transport user flows between the femtocell and femtocell gateway with SCTP. SCTP has a number of key features which can be used to enhance femtocell functionality; these include transport layer multi-homing, multi-streaming, concurrent multi-path transfer and partial reliability. It is worth

noting that every femtocell and femtocell gateway already utilises SCTP for transporting control plane data and as such it would not require any major modifications to existing femtocell protocol stacks.

### 3.6.2 Key SCTP Features

This section details the SCTP features that are relevant for enhancing Femtocell functionality with multi-homing capability.

#### 3.6.2.1 Multi-Streaming

An SCTP association utilises multi-streaming to separate user flows. Each stream is a logical uni-directional communication channel between the two endpoints with independent delivery including message segmentation and reassembly. Essentially this means that there is no Head of Line (HOL) blocking between streams, allowing different user flows with potentially different requirements to be multiplexed over the same transport layer connection without interfering. It is worth highlighting that SCTP can also provide unordered delivery similar to UDP thereby preventing HOL blocking within a stream and reducing latency. As will be detailed later, this work proposes utilising these streams to provide the same functionality as multiple GTP-U tunnels.

#### 3.6.2.2 Multi-Homing and Dynamic Address Reconfiguration

A key feature of SCTP is the ability to allow a single transport layer association to span multiple Internet Protocol (IP) addresses, thereby providing increased resilience to link failures. As has been shown in previous work, SCTP can seamlessly handover between available addresses with no interruption to any ongoing service including voice calls. Indeed the handover time is simply equal to the end-to-end delay between the two nodes, in this case the femtocell and femtocell gateway. Due to the initial design requirements for SCTP, replacing SS7 in all IP core networks, the initial variant of SCTP did not allow addresses to be added/deleted from an association after setup. This meant that all addresses had to be specified during association initialisation. The Dynamic Address Reconfiguration (DAR) extension [93] solved this by providing mechanisms to dynamically add/delete addresses from an association thereby creating a far more flexible protocol. Indeed it is this extension which enabled much of the work on utilising SCTP as an endpoint-based mobility solution. In the context of the multi-homed femtocell scenario being considered in this work, the extension would mitigate load on the cellular network by allowing the femtocell to dynamically remove the cellular link when it is not required (e.g. such as periods of low activity or when no User Equipment (UEs) are connected to the cell).

#### 3.6.2.3 Concurrent Multi-Path Transfer

The Concurrent Multipath Transfer (CMT) extension allows simultaneous transmission across multiple interfaces in multi-homed SCTP associations. It introduces mechanisms to manage the side-effects that can occur when spanning the same association across multiple interfaces, particularly when links have diverse characteristics. These are issues such as packet reordering, increased acknowledgements and retransmission behaviour. The CMT extension would allow more fine granular control over network offload while maintaining a high quality of experience for end users. For example, if the fixed broadband network over which the femtocell is being backhauled cannot meet the strict delay requirements of voice calls these may be transported via the cellular network, while non-real time traffic can transported over the broadband network.

#### 3.6.2.4 Partial Reliability Extension

The Partial Reliable SCTP (PR-SCTP) extension [94] allows the sender to control the level of reliability required when transmitting each message. Specifically, it allows the sender to configure how aggressively SCTP should perform retransmissions. This done by specifying a lifetime value for each message, after this period expires SCTP stops attempting to transmit the message. A nice advantage of this feature is that it can be configured individually for different services. For example, real time services such as voice would utilise relatively short lifetime values to minimise retransmissions while less time constraint services such as streaming video could use higher lifetime values to reduce packet loss. The extension also defines a new chunk type, FORWARD TSN, that simply informs the receiver to move its acknowledgement point forward, thereby skipping any messages that have not yet been received and preventing the receiver from reporting them as missing. This extension offers a number of advantages. The most relevant here being the ability to provide unordered and unreliable transport, similar to UDP, while still being able to utilise other key protocol features such multi-streaming and multi-homing.

### 3.6.3 Multi-homed Femtocell Architecture

In this section an overview of the changes required to the architecture and protocol stack are detailed. The only required architectural change is the inclusion of an additional interface at both the femtocell and the femtocell gateway. Clearly the femtocell requires the addition of an uplink cellular radio to enable

relaying of the UE traffic via the mobile network. However, the femto gateway also requires an extra interface. Using SCTP each sending node can only specify from which particular interface a message is sent based on the destination IP address. As such the routing tables in each node must be configured such that the destination address of each interface in the alternate node is only reachable via a specific interface on the transmitting node.



**Figure 3.33: SCTP Enhanced Femto Protocol Stack**

### 3.6.3.1 GTP-U/UDP to SCTP

In order to successfully utilise SCTP as a replacement for GTP-U over UDP, all required fields must be mapped to equivalent SCTP fields or at the very least meet the same functionality. The key functionality for GTP-U in a femtocell network is to separately tunnel individual user flows between the femtocell and femtocell gateway. It should be noted that each user may have multiple tunnels depending on the number of active Packet Data Protocol (PDP) contexts. To provide this functionality the two key parameters involved in a GTP-U tunnel are the Tunnel Endpoint Identifier (TEID) and the optional sequence number. The TEID is a 4 Byte field that individually identifies the tunnel into which the user flow data is encapsulated. The sequence number field is specific to each TEID and allows for packet reordering to be done. Although the sequence number field is optional it is used in most femtocell deployments.

This work proposes replacing the GTP-U TEID and GTP-U sequence number with the SCTP stream identifier and stream sequence number, respectively. The SCTP stream identifier is 2 bytes in size compared to the 4 bytes that TEID supports, this obviously significantly reduces the number of active sessions that are possible. However, the stream identifier can support over 65,535 connections and it is unlikely that a single femtocell will have more than this number of simultaneous connections. The per-stream sequence number space of SCTP is also smaller, but sufficiently dimensioned to not become a performance bottleneck over a last-mile access.

| Bits | 0-7 | 8-15 | 16-23 | 24-31 |
|------|-----|------|-------|-------|
| 0 | Source Port | | Destination Port | |
| 32 | Verification Tag | | | |
| 64 | Checksum | | | |
| 96 | Chunk N Type | Chunk N flags | Chunk 1 length | |
| 128 | TSN | | | |
| 160 | Stream Identifier | | Stream Sequence Number | |
| 192 | Payload Protocol Identifier | | | |
| 224 | Chunk N data | | | |

**Figure 3.34: SCTP Packet Structure**

| Bits | 0-2 | 3 | 4 | 5-7 | 8-15 | 16-23 | 24-31 |
|------|-----|---|---|-----|------|-------|-------|
| | GTP-U | | | | | | |
| 0 | Version | P | T | Spare | Message Type | Total Length | |
| 32 | TEID | | | | | | |
| 64 | Sequence Number | | | | | Spare | |
| | UDP | | | | | | |
| 96 | Source Port Number | | | | | Destination Port Number | |
| 128 | Length | | | | | Checksum | |
| 160 | Data | | | | | | |

**Figure 3.35: GTP-U/UDP Packet Structure**

### 3.6.3.2 Cellular-Broadband Handover

Standard SCTP provides a path fail-over mechanism which is only performed after a maximum number of failed retransmissions. Clearly for the femtocell case this would lead to non-seamless transition between interfaces resulting in disruption to any ongoing services. As described earlier, the DAR extension solves this by allowing each endpoint to selectively switch between paths thereby providing seamless handover capability. The handover delay is simply be the end-to-end delay between the femtocell and femtocell gateway and as has been demonstrated in previous work [95] this can be utilised to seamlessly handover voice calls without any interruption. Figure 3.36 shows the association setup and handover of a voice call between a multi-homed femtocell and femtocell gateway. Unlike UDP the association setup in in SCTP requires the exchange of four messages. This is done to prevent denial of service attacks. Although this adds a small amount of additional overhead, connections between the femtocell and femtocell gateway will be long lived and hence this additional overhead is negligible. During association initialisation the femtocell and femtocell gateway exchange the list of addresses through which they are reachable. In this case and with reference to the considered architecture shown in Figure 3.33, it is assumed that IP1 in each node is the fixed broadband network and IP2 is on the cellular network. The initial primary is set as IP1 on both nodes, this defines the specific network over which they are communicating (e.g. the fixed broadband link) with the secondary being set at IP2 (e.g. cellular network).

A voice call is then initiated over the femtocell via the fixed broadband network toward the femtocell gateway. At some point a handover is triggered by the femtocell, however it should be noted that this handover could also be triggered by the femtocell gateway. The exact trigger is outside of the scope of this work. However this could utilise a number of metrics. For example, the voice call quality could be continually monitored and if it degrades below some predefined threshold a handover to the alternate network could be triggered. Other metrics such as transport layer retransmissions, packet loss and delay jitter could also be used.

The femtocell sends an ASCONF: SetPrimaryIPAddress chunk to the femtocell gateway specifying it's second address. On reception of this the femto gateway sets the cellular network IP address as the primary address through which the femtocell is reachable. At the same time the femtocell locally sets the cellular network IP address of the femtocell gateway as the primary address through which is reachable. Both of these actions result in a full handover from the fixed broadband to the cellular network.

**Figure 3.36: SCTP Signalling Diagram**

### 3.6.4  Overhead Analysis

#### 3.6.4.1  Common & Chunk Header

As can be seen by comparing Figure 3.34 and Figure 3.35 the common SCTP header consists of a source/destination port pair, verification tag, and a 32 bit checksum giving a total header size of 12 Bytes. This is 4 Bytes larger than the UDP common header at 8 Bytes. Unlike UDP each SCTP packet can contain multiple data chunks and hence this can have an impact on the total overhead.

Each data chunk contains a chunk header consisting of a 1 Byte chunk type field, a 1 Byte chunk flags field, a 2 Byte chunk length field, a 4 Byte Transmission Sequence Number (TSN), a 2 Byte stream identifier, a 2 Byte stream sequence number and a 4 byte payload protocol identifier. This introduces an overhead of 16 Bytes per data chunk. In order to compute the worst case overhead, it is assumed that each packet contains only one data chunk; this results in a total SCTP overhead of 28 Bytes or 20 Bytes more than required by UDP to transmit the same data payload.

As shown in Figure 3.35 the GTP-U header is 12 bytes. Hence the increase in overhead drops to 8 bytes per data payload when both the GTP-U and UDP are replaced with SCTP. This difference in overhead reaches a break-even point when more than three data chunks are transmitted in a single SCTP message. This ability of SCTP to bundle multiple data chunks into a single message means that when more than three data chunks are transmitted in the same SCTP message, SCTP becomes a more overhead efficient transport mechanism.

SCTP can bundle as many data chunks as required to fill a message to the Maximum Transmission Unit (MTU); this is controlled by a user definable bundle time-out value which specifies the maximum amount of time that a data chunk is delayed while waiting on other data chunks with which it can be bundled; clearly for voice traffic this timeout should be quite low. However, considering the high number of packets produced by voice calls and the small size of these packets, it is clear that chunk bundling has the potential to greatly improve overhead efficiency. The advantage of this is that during periods of low traffic levels SCTP is slightly less efficient, however as traffic levels increase this will allow for higher levels of chunk bundling thereby making SCTP more overhead efficient.

#### 3.6.4.2  Selective Acknowledgements

SCTP utilises a delayed Selective Acknowledgement (SACK) mechanism which informs the alternate endpoint of the last successfully received message before any gaps in the received data and any subsequently received or duplicate messages; these are determined based upon the TSNs of received frames. This mechanism is used to prevent head of line blocking and to optimise retransmission by only retransmitting missing messages. These acknowledgements and retransmissions clearly introduce

additional overhead compared with UDP which has no such mechanisms. As previously discussed, retransmissions can be controlled utilising the PR-SCTP extension thereby mitigating additional overhead.

The SCTP implementation guidelines [96] recommend that each receiver should acknowledge at least every second received packet or within 200ms of receiving a chunk; however this is not a requirement. The requirement simply states that any implementation must not allow the maximum delay to exceed 500ms. Furthermore, SCTP allows the bundling of control chunks with data chunks, this means that SACKs can be sent with ongoing data traffic. In the case of voice calls being transmitted via a femtocell, a voice packet is transmitted in each direction every 20ms. Therefore the number of required packet transmissions can be minimised as SACK chunks will always be bundled with voice data chunks. Assuming no packet losses, each SACK chunk is 16 Bytes and using the SCTP guidelines it must be sent with every second data packet this will introduce an average additional overhead of 8 Bytes per SCTP message. If the guideline is ignored and only the maximum delay requirement is considered, the additional overhead on a voice call is negligible at less than one Byte per packet.

### 3.6.4.3  Heartbeats

Each multi-homed SCTP endpoint must maintain information on the reachability state of each peer endpoint via all available IP addresses. SCTP achieves this using a keep alive mechanism which involves periodically transmitting a 12 byte heartbeat chunk to all peer endpoint addresses not currently be used as the association primary. In response to this heartbeat chunk each endpoint must respond with a 12 byte heartbeat acknowledgement chunk.

The default periodicity of the heartbeat mechanism is quite high at 30 seconds and therefore adds negligible additional overhead. However, it is a configurable parameter and can be completely disabled. For a multi-homed femtocell utilising a mobile radio link it would be more desirable to use PDP context information rather than rely on the SCTP heartbeat mechanism and as such it may be disabled in such scenarios.

### 3.6.5  Conclusion

This work proposed a mechanism to enable multi-homed femtocells by utilising SCTP as the user plane transport protocol between femtocells and femtocell gateways in place of the current tunnelling protocol combination of GTP-U/UDP. This work also discussed the key benefits of such an approach and discussed the trade-off in terms of overhead. It was shown that SCTP's bundling of data chunks allows fully amortizing additional overhead, making it even more efficient than GTP-U/UDP in cases of high traffic loads.

In the following project year, this multi-homed femtocell approach will be studied applied and analysed in the context of Mobile Femtocells.

## 3.7  Deployment, Handover and Performance of Networked Femtocells in an Enterprise LAN

### 3.7.1  Introduction

This section presents the investigation results on 3G networked femtonodes, deployed in an enterprise LAN with the objective of study the limitations of the analysed implementations and derive some conclusions that could of interest in the development and specification of new products or concepts as the LFGW,as central point of an enterprise networked femtonodes group. Many investigations and simulations are currently carried out on networked femtonodes solutions but the final assessment, the selection of the best solutions and the some problems only arise with real world experiences. The integration of a network of femtonodes in an enterprise LAN is a field quite new, justified by operators' interest of femtonode deployments in the enterprise.

In this study is assumed that the connection of individual standalone femtonodes to the femtonode subsystem through individual ADSL routers or ONT fibre connections don't present major problems, as it is the usual method of current deployments, and therefore not individual femtonodes test have been realized.

Conversely to the usual standalone femtonodes, used in residential deployments, it were used enterprise grade femtonodes, using the enterprise LAN to connect with the femtonode subsystem with additional features as higher transmission power and the support of handover between the enterprise femtonodes belonging to the same group, using the enterprise LAN to convey some intergroup handover and synchronization signalling.

The investigation were not focused on femtonode as an independent entities features, but rather were focused on the behaviour of a group of networked femtonodes, inserted in the communications infrastructure (LAN) of an enterprise, where the femtocells backhaul traffic will have to co-exist with the enterprise traffic and cannot be expected that the femtocells traffic can be prioritized over the LAN traffic.

It is important to highlight that as in the case of standalone femtocells the networked femtocells, the connection to the operator core network is carry out through an intermediate femtocell subsystem, as it is shown in Figure 3.37.



**Figure 3.37: Networked femtocell deployment setup.**

Specifically it was studied:

The LAN configuration changes that are needed to support the femtocell group, keeping in mind that ideally these changes will have to be minimum and easy to implement by the enterprise network administrator.

The logical connectivity of the femtonodes to their corresponding femtonode subsystems.

Initial femtonode group radio planning, and radio coverage

Network and radio femtonodes self-configuration.

Mobility, synchronization, and performance

### 3.7.2 Deployment description

The pilot was deployed in the ground floor of a Telefónica I+D building in C/Don Ramón de la Cruz 82-84, in Madrid. Femtonodes from two manufacturers that will be referred as manufacturer A and manufacturer B were installed in the same locations in order to allow a consistent performance comparison. Also two femtonodes subsystems (A and B) were used in this pilot to provide service to the corresponding femtonodes

The femtonodes were connected to TID's LAN, and connected through a IPSEC tunnel to each respective femtocell subsystem located in Telefónica's radio labs in Pozuelo (Madrid). It must be noted that the enterprise femtonodes pilot was not connected to the general mobile core network that provides service to Spain's customers, but to an evaluation core network used for testing purposes and used a different frequency and PLMN ID that the commercial mobile network, to avoid uncontrolled interactions with real customers.

Two femtocell subsystems were used; one for each manufacturer, as it is not possible to disclose some concrete details, next figure shows the generic architecture of a femtocell subsystem and the network architecture used in this pilot. The generic elements that compose a femtonode subsystem are:

Security GW (IPSec Gateway).

IP Clock Server, that provides a reference clock for femto synchronization.

Configuration Server (Femtocell Manager), which configures the femtonodes radio parameters.

Femto Manager (Femtocell Home Register).

Femto subnetwork manager, controller of the femtonode subsystem network.

AAA Server (Authentication, Accounting and Authorization Server).

Femto Gateway (Access Gateway).

Figure 3.38 depicts the network architecture of the enterprise femtocell pilot. In this picture the Femto Manager includes the AAA Server



**Figure 3.38: Network architecture of the enterprise femtocell pilot**

Three femtonodes from every vendor were used during the trial. The characteristic of each femtonode are presented in Table 3-7.

**Table 3-7 Femtonode Characteristics**

| Characteristics | Manufacturer | |
|---|---|---|
| | **VENDOR A** | **VENDOR B** |
| Simultaneous services | 1 CS & 3 PS per user | 1 CS & 3 PS per user |
| Number of users | 32 CS and 16 PS | 8 CS and 8 PS |
| HSDPA | 7,2 Mbps | 7,2 Mbps |
| HSUPA | 1,44 Mbps | 1,44 Mbps |
| Power | 20 dBm | 20 dBm |
| Sensitivity | -110 dBm | -113 dBm |

Figure 3.39 shows the basic scheme of the deployment scenario in TID premises, where three femtonodes per vendor were installed and the coverage overlap area was important.

The total area to cover has been 1,900 m2, and the femtonodes were installed in the upper side of three pillars, close to the ceiling and beside the enterprise Wi-Fi access points. The installed femtonodes picture is a low resolution one, in order to not unveil the femtonode manufacturers deployed in the trial.

**Figure 3.39: Femtonodes locations in TID premises**



2 femtonodes + Wi-Fi AP

**Figure 3.40: Two femtonodes and one WiFI access point installed in one of the office's pillars in TID premises.**

In order to avoid any interference with the macro layer, it was decided to make use of the UMTS carrier assigned in Telefónica Spain for indoor coverage, i.e. UARFCN 10788. The next table shows the available radio carriers, and grey highlighted the carrier used in the trial.

| UARFCN DL | DL Centre Frequency (MHz) | UARFCN UL | UL Centre Frequency (MHz) |
|---|---|---|---|
| **10788** | **2157.6** | **9838** | **1967.6** |
| 10813 | 2162.6 | 9863 | 1972.6 |
| 10838 | 2167.6 | 9888 | 1977.6 |

Regarding the PSC (Primary Scrambling Code) that were used by the femtonodes, they were:

Vendor A; the femtonodes used PSC's in the range from 50 to 55. The femtonodes automatically scan the PSC's in use and select the one with best signal to noise ratio (Ec/N0).

Vendor B, the femtonodes used PSC's in the range from 134 to 136. In this case, the femtonodes select the first PSC from a preloaded list, after detecting if it is not in use by other femtonode of the group.

### 3.7.3 Femtonodes IP network configuration

In order to connect the femtonodes to the Enterprise LAN, and to ensure its connectivity to the IPSec and Access gateways of the femtonode subsystems, in Pozuelo labs, through the public IP network, the following steps were followed:

Request to the network administrator of static local IP addresses for the femtonodes and other femtocell subnetwork parameters (network mask, broadcast address, subnetwork name, address range, etc)

Configuration of DHCP server, with the mapping of femtonodes MAC addresses with the femtonodes local IP addresses and the DNS server IP address.

Configuration in the Network Address Translation (NAT) server of the mapping between the femtonodes local and public IP addresses (static NAT).

Configuration of the IPSec tunnel for the inbound and outbound connection of the networked femtonodes.

Physical installation of the femtonodes including Ethernet cabling needed to connect the femtonodes to the switch that connect the different building subnetworks.

The network administrator provided the following femtonode subnetwork configuration details::

- Subnetwork name: Pilot Femto.

- Network: 10.95.87.80 / 28

- Range: 10.95.87.81 - 10.95.87.94

- Netmask: 255.255.255.240

- Broadcast: 10.95.87.95

Once the femtonode has a local IP address, the femotonode is connects with the Configuration and AAA Servers to be authorized to operate, after a security process that includes the exchanging of security keys. If the process is successful it is established a femtonode - SeGW IPSeC tunnel. Then the network administrator must configure access to the public IP network of the femtonodes, through the firewall/edge router. For this purpose each manufacturer provided the IP address assigned to its SeGW and list of the types of protocols and ports used in the communication between the femtonodes and the femtonode subsystems as:

- IKE (Internet Key Exchange) the protocol used to key exchange between the femto and femto subsystem. IKE use UDP as its transport protocol. Initially IKE goes over UDP port 500 (IKE_SA_INIT), and in a second phase (IKE_AUTH) goes over UDP port 4500, because it is used NAT-T (NAT traversal in the IKE). NAT-T is a method of enabling IPsec-protected IP datagrams to pass through Network Address Translation (NAT).

- IPSeC. One time the femtocell is authenticated and authorized, the IPSeC tunnel is established and all the traffic, except the IPClock traffic, is encapsulated.

- NTP, the protocol to transmit the IPClock client-server signals from/to the femto subsystem. Depending of the manufacturer this traffic was in IPSeC tunnel or outside

### 3.7.4 Femtonodes authentication and authorization

It was verified that the femtonodes authentication and authorization processes worked properly, when the femtonodes and the enterprise firewall was configured to share the enterprise IP address and also when they were using and individual IP address.

The process of authentication and authorization was successful. Figure 3.41 presents the authentication and IPSeC tunnel creation process.

IKE (Internet Key Exchange) is the protocol used to set up a security association in IPSeC protocol suite. IKE consists of two phases: phase 1 (IKE_SA_INIT) in which is established a secure authenticated communication channel by using a key exchange algorithm to generate a shared secret key to encrypt further IKE communications and phase 2 (IKE_AUTH) in which the IKE peers use the secure channel established in Phase 1 to negotiate Security Associations for IPSeC. IPSeC tunnel payload is ESP (Encapsulating Security Payload)

As we can see the femtonodes, with IP address 10.95.87.85 initiates the tunnel creation with an exchange of security keys process (IKE_SA_INIT). The authentication and authorization process finalized correctly and the IPSEC tunnel was created.



**Figure 3.41: Key exchange and IPSeC tunnel establishment between the femtocell and the femtocell subsystem**

Figure 3.42 depicts the certificate based mutual authentication between the femtonode and the femtonode subsystem. Initially the femtonode contact with the AAA server, through the SeGW, which in turn accesses the HSS that is where the authentication data is stored. Once the AAA server finds that the femto-node is authorized for operation, it completes the process and establishes the IPSec tunnel. The HSS, in other parts of this document is called Configuration Server of femtonodes.



**Figure 3.42: Authentication and authorization based on EAP-AKA, before the creation of the IPSec tunnel (3GPP TS 33.320)**

### 3.7.5  IP connectivity between the femtonodes and the femtonode subsystem

Two options were considered to connect the group of femtonodes to Internet, **femtonodes sharing the enterprise public IP address (dynamic NAT)** and **femtonodes using an individual IP public address (static NAT)**. The first option was the preferred, because is the usual enterprise policy for management and cost reason. Initially it were analysed the connectivity of one femtonode of each manufacturer and it was validated (only the individual IP addressed method worked properly at the end) it were carried out the group of femtonodes test (e.g. handover).

#### 3.7.5.1  Femtonodes Sharing the Enterprise Public IP Address

As it was said before, this was the initial preferred option to connect the group of femtonodes to the femtonode subsystem. This approach was implemented in the enterprise firewall using *Hide NAT*.

With *Hide NAT*, multiple femtonodes can access the Internet sharing the same public IP address, so that the firewall uses dynamically assigned port numbers to different connections, corresponding to each internal host. When a packet from the internal network try to reach Internet, the firewall intercepts it and modifies the source IP and port. The internal source IP is replaced by the public IP (usually shared for all enterprise traffic) and source port is changed to a dynamically allocated port to uniquely identify the connection. The port allocation is dynamic and deleted where a timer expires. The relationship between the dynamically assigned port and internal IP is stored in the *dynamic* firewall state tables, so that when a response packet arrives, the firewall uses the destination port to determine which connection is for the package and adjust the port and the destination IP properly.

Establishing a many to one relationship between the private IP and public IP behind which "hide", which allows internal hosts to establish connections with external servers but cannot make connections to the contrary. On the other hand, the dynamic NAT cannot be used with protocols that require that does not change the port number or required to distinguish between internal hosts based on IP address. Hide NAT present the advantage that only one public address is needed.

To verify the connectivity femtonode – femtonode subsystem, once the femtonodes were authenticated, the femtonode subsystem was configured to send PINGs (ICMP packets) to one femtonode of each manufacturer, using static NAT. The results of these initial tests were:

For the femtonode_A. It was observed that packet loss exceeds 10%.

For the femtonode_B. Also, there were packet losses, but in this case it was not possible to maintain communication for more than 5 minutes, by IPSeC tunnel drops. IPSeC tunnel automatically restarted process after this period.

After analyse traffic traces under the supervision of the network administrator, responsible of the enterprise network configuration, it were reached the following conclusion:

***Femtonode subsystem A***

By the default the externally originated traffic is blocked by the firewall, but when communication is established between an internal source and an external destination, it is created a temporal response channel that allows accept answer / acknowledgment packets automatically, without having to configure a firewall rule to accept external data sources. This feature is configured with the functionality "Accept UDP replies", with the parameter "UDP Virtual Session Timeout" that determine the amount of time a UDP response channel can remain open, which by default is set to 40 seconds. Then the response channel (direction SeGW-femtonode) is closed every time expires "UDP Virtual Session Timeout" (in our case every 40 seconds). Conversely, in the case of TCP packets, the response channel (direction SeGW-femtonode), is keep open to allow bidirectional traffic.

When an outbound packet was sent by a femtonode (e.g. a keep alive packet used by the IPSeC connection), the firewall accepted external packets from the destination direction during 40s, but after the external packets blocked (e.g. from the femtonode subsystem), until another internal packet was issued by the femtonode to the femtonode subsystem. This supposes that some femtonodes incoming communication can be blocked. This did not affect to the process of key exchange and authentication because is a highly bidirectional traffic and initiated by the femtonode.

It was enabled the external access for traffic from the SeGW (UDP/4500), but it was found that was not working properly, because after the "UDP Virtual Session Timeout" the firewall state table is deleted, being the firewall unable to identify the internal host to which it should delivered the incoming packet, returning to the source an ICMP message type "Time-to-live-expired". This solution supposes some security concerns, because it is opened continuously a port to enter in the company LAN.

*Femonode subsystem B*

In addition to the 10% of packet losses, like the femtonodes system A, every 5min the IPSeC tunnel that connected the femtonodes and femtonodes subsystem B dropped. This was motivated because the SeGW belonging to femtonode subsystem B, always sent the packets to the port (and public IP address) where the initial femtonode initiated the communications with the SeGW. When the NAT dynamically change the output port of femtonodes output traffic, the SeGW do not recognized the change and continued to send traffic to the original port, as the femtonodes did not answered, it occasioned the IPSeC tunnel dropping. After a time-out, the femtonode B restarted the process to establish a new IPSeC tunnel with the SeGW, recovering the communications with the femtonodes subsystem. The SeGW of the femtonodes subsystem A did not experienced this problem, because it was able to follow the changes of the femtonode ports in their outbound messages using the IPSeC consecutive sequence numbers.

To solve this issue, there are two suboptimal possibilities:

- Ensure traffic generation outbound direction (femtonode-> SeGW) more frequently than the "UDP Virtual Session Timeout". It supposes to increase traffic load.

- TCP instead of UDP in the connection between femtonodes and its SEGW. But this supposes more protocol surcharge.

### 3.7.5.2 Femtonodes with dedicated Public IP Address

Due to the problems found to implement dynamic NAT it was implemented static NAT, even when this solution could be difficult to be implemented in enterprise LAN, because each femtonode needs a public IP address.

*Static NAT*. Static NAT establishes a direct relationship between the femtonode internal IP address and its external and public IP address. Because of this relationship one to one, a public IP is required for each node that accesses the Internet with static NAT. This translation allows for connections to internal nodes through their public IP, when the necessary permissions are enabled. The correspondence between privates and publics IP addresses are stored in the firewall *static* state tables.

In the case of the femtonodes IPSeC connections, it was observed that IKE and ESP, is encapsulated over UDP with source and destination port 4500. When applying static NAT the packets internal IPs addresses are changed by the external ones and vice versa in response packets. The ports are not changed, as the case of Hide NAT.

The protocol ports were the same that in the case of sharing IP address.

Conversely to the case of sharing IP address there were not packet losses problems because the communication channel never expire.

### 3.7.6 Coverage simulation

It was estimated the femtonodes coverage by simulations in the deployment scenario, using an internal simulation tool, that takes into account the materials used in the simulation scenarios (walls thickness and composition, crystal windows, etc ). In deployment it was used a fix pilot transmission power (CPICH) of 10 dBm (10 mW) and a maximum power of 20 dBm (100 mW), these data also were used in the simulations.

Figure 3.43 presents the simulation results. The coverage is quite good with few zones in blue (poor coverage) and those zones correspond to a courtyard and auxiliary rooms, separated by thick walls from the working zone.

**Figure 3.43: Estimation of femtonodes coverage**

The power radiated by the femtonodes can be adjusted by the operator through a change in the femtonodes profile, in the Femtonodes Configuration Server.

### 3.7.7 Mobility between networked femtonodes

Mobility tests were conducted between the three test femtonodes, isolated from the 3G macrocell commercial network, for this purpose it was used a PLMID and frequency (DL UARFCN 10788) different of the macrocells. It was tested voice call handovers (CS handover), because only one of the femtonodes supported packet data handover (PS handover).

The tests for checking handovers between femtocells were done using an engineering phone equipped with a SIM card for the PLMID used in trials, calling a voice server which delivered a continuous sound and going through the femtocell group (F_3, F_2 and F_1). The engineering phone provided information on the serving, the monitored, and the detected cells and also on possible call interruptions and call drops during the handovers. The test was repeated by initiating the call in each of the positions shown in Figure 3.43, where we can observe the overlapping areas between femtonodes.

The transfer between femtonodes is based on measurements performed by the mobile on both, the serving femtocell and as the neighbours. The handover is triggered when the signal quality of the femtocell neighbour is better than the current service with a certain hysteresis value.

The handover between femtonodes is not performed by the core mobile network, it is carried out by the subsystem of femtonodes. The control and user plane is switched from the source to the target femtonode by the Femto Gateway (see Figure 3.38). In both cases, the PDP context of the user equipment is maintained in the SGSN and is restored when the terminal is camping on another cell when the transfer of voice (CS) is completed

The two femtonodes systems were configured in a different way. Femtonodes A were configured to use hard handover and femtonodes B used soft handover.

#### 3.7.7.1 Configuration of networked femtocell groups

The two femtonodes systems used a different way of creating the group of networked femtonodes, among those will took place handovers. The femtonodes A system approach was based on the realization of initial and periodical scanners (typically every 12 - 24 hours) to detect neighbours femtonodes using the same frequency and PLMID. In the system of femtonodes B, femtonodes which constituted the group,

should be registered in the management system, and this information is communicated to each femtonodes of the group, in order to be included automatically in the list of neighbouring femtonodes.

In both cases (femtonodes B and A) the individual femtonodes carry out an initial scanning to choose the best scrambling code, not in use by neighbouring femtonodes, the difference is on the way in which each femtonode creates the list the neighbouring femtonodes. The approach of the femtonode group A is totally based on scanning results and option B is based on the O&M configuration data.

### 3.7.7.2  Handover between femtonodes belonging to the femtonodes subsystem A

The first step was to configure the femtonode cluster, using the femtocell configuration server, configured with the frequency and the list of the possible scrambling codes that the femtonodes choose in the initial scanning.

In the test conducted with the femtocell system A, it was observed that the handover of calls from the femtocell F_2 to the femtocell F_1 consistently failed. Figure 3.44 shows the scrambling codes (PSCs) transmitted in the BCH channel and PSCs detected by engineering the telephone in the vicinity of each femto node. Femtonode the F_2, did not detected F_1 as neighbour, not being broadcasted F_1 by F_2 and therefore not been taken as femtonode target by the engineering phone when handover, losing the call when was approaching to F1. However, the call initiated in the F_1 femtonode is transferred to F_2, because F_1 was transmitting correctly all neighbouring femtonodes. By repeating the tests sometimes F_2 detected correctly all the other femtonodes as neighbours, depending on radio channel conditions.

It should be highlight that femtocells F_3 and F_1 were nearly in line of sight, and F_2 was not in line of sight respect F_3 and F_1.



**Figure 3.44: Scrambling codes broadcasted by the femtonodes belonging to subsystem A**

Figure 3.45 shows the Wireshark traces corresponding to handover failure when the mobile was moving between F_2, F_3 and F_1. When the mobile was moving from F_2 to F_1 the call dropped. We can see how the femtocell F_3 to transferred briefly the call to F_2, because it was the only one in its list of neighbours, but because the mobile moves out of radio coverage of both femtonodes and did not try to connect with F_1, the call dropped. After the call drop, it was initiated a new call, that was processed by F1 and it was carry out a handover from F1 to F_2 without problem.

The femtonodes were restarted to force the creation of a new the list of neighbours in each femtonode, based on the scanning results, and sometimes F1 was included in F_2 list of neighbours, and in this case there were not handover problems with F1.

**Figure 3.45: Call drop during a handover among F_2 and F_1 using self-scanning of eibur femtonode to create the femtonode neighbouring list and hard handover**

We can extract the conclusion that the creation of group of networked femtonodes, based on scanning of neighbours is a highly adaptive technique that can automatically add or remove neighbour femtocells, but present the inconvenience that in the moment to proceed the scanning of a neighbour femtocell, its signal is temporally hidden or interfered, it will be not included in the neighbour list, and therefore handover to this femtocell will be not possible.

Regarding to the hard handover used by femtonodes, it was noticed a small communication interruption just in the moment of call transfer between femtonodes.

### 3.7.7.3 Handover between femtonodes belonging to the femtonodes subsystem B

Like the femotonodes subsystem A, the first step was to configure the femtonode cluster, using the femtocell configuration server, configured radio parameters as; frequency the list of the possible scrambling codes that the femtonodes choose in the initial scanning, Cell ID and RNC ID

In all the femotonode subsystem B handover tests were performed without problems during handover when a voice call was transferred from a femtonode to another, even without appreciating the small interruption in the continuous beep in the phone used in this test. This was due to in the femtonode subsystem B was configured to use soft handover that translate in a seamless handover. However, in some cases it was observed that communication, as a result of the transfer, was degraded to the point that was left to hear the beep and all I heard was a continuous noise.

Figure 3.46 shows schematically the soft handover procedure in the femtonode subsystem B. The user terminal sends regularly measurements reports to serving femtonode of both the serving femtonode (Femto-1) and also of the neighbouring femtonodes (Femto-2) included in the list of neighbouring femtonodes. When the RSCP measurement of the serving femtocell is below a certain threshold, the serving femtocells communicates to the terminal that it will also receive data from the neighbour femtocell with a better signal to noise ratio (Femto-2).

At the time of soft handover, the terminal receives data simultaneously Femto-1 and Femto-2. Data for Femto-2 is provided by Femto-1 via the LAN to which are connected the group femtonodos (similar to the Iur interface, between the RNCs in UMTS). It is important to highlight that, the subsystem sends to each femtonode a table with the local IP address of each femtonode of the same group.

Once the target femtocell definitely takes over the call, the data flow from the Femto Gateway (similar to an Iuh interface) is switched from the Femto-1 to Femto2.

**Figure 3.46: Sof handover procedure in femtocells**

In all the femotonode subsystem B handover tests were performed without problems, it was not appreciated any interruption in the continuous beep during handover, when a voice call is transferred from one femtonode to another. This was due to in the femtonode subsystem B was configured to use soft handover, which translate in a seamless handover.

Figure 3.47 shows the correct operation of the smooth transfer reflected in the trace captured by Wireshark for a voice call initiated in the femto cell F_2 (central position of the pilot deployment).



**Figure 3.47: Handover between femtocells using neighbouring list creation using O&M subsystem and soft handover**

### 3.7.7.4 Conclusions

**Networked femtonodes and femtonode subsystem IP connectivity**

Two options were considered to connect the group of femtonodes to Internet, femtonodes sharing the enterprise public IP address (dynamic NAT) and femtonodes using an individual IP public address (static NAT). The first option was the preferred, because is the usual enterprise policy for management and cost reason

- *Sharing the Public IP Address*. It was not possible to connect the femtocell subnetwork to the femtocell subnetwork sharing the public IP address of the enterprise LAN (dynamic NAT with port translation). With this approach it was created the IPSeC tunnel without problems, but there were packet losses when the IPSeC tunnel closed in direction (SeGW → femtonode) when the firewall timer "UDP Virtual Session Timeout" expired and there were not traffic from the femtonode to the SeGW.

- *Dedicated Public IP Address*. With femtonodes with dedicated Public IP Address there were not packet losses between the femtonodes and the SeGW.

**Authentication, authorization and IPSeC tunnel creation**

- It was configured the enterprise firewall/NAT for the connection with the femtonodes subsystem. There were no problems on the femtonodes authentication and authorization and creation of IPSeC tunnel processes, because the associated signalization is highly bidirectional, and therefore it is not affected by the expirations of timers.

**Creation of networked femtonodes group and handover**

- Two approaches have been analysed on the creation of the networked femtonodes group. The first one based on the individual scanning of networked femtocells and the second one is based is based on the O&M subsystem configuration. The result has been that the option based on self scanning can fail due to shadowing, interference or noise in the moment in which it is carried out the scanning and the second one in which the femtonode neighbouring list is communicate by the O&M subsystem provide more stable and predictable results, without possible handover failures, but present the drawback of lack of flexibility due to manual configuration. An hybrid solution in which new femtonodes could be automatically added looks to be an optimal solution.

# 4. Network Management

## 4.1 Distributed Fault Diagnosis

### 4.1.1 Introduction

Fault Diagnosis in BeFEMTO focuses on the enterprise networked femtocells scenario and is designed as a distributed framework that allows local management capabilities inside the femtocell network. In order to implement this distributed framework, a multi-agent approach is followed. Fault diagnosis is therefore cooperatively conducted by a set of cooperation agents distributed in different domains which share their knowledge about network status.

In particular, Fault Diagnosis will target the diagnosis of problems related to the use of video services while being under the coverage of an enterprise femtocell network.

Diagnosis knowledge in BeFEMTO is modelled as a Bayesian Network (BN) which relates causes and symptoms by means of probabilities. This approach is well suited to deal with uncertainty and lack of full status information. Furthermore, diagnosis knowledge can be split up and appropriately distributed to the different domains involved.

D5.1 described this approach and the work progress made until end of 2011. In this document a more detailed description of the diagnosis framework is provided.

### 4.1.2 Architecture overview

In BeFEMTO there will be different types of agents which are briefly described below:

Interface agents: These agents serve as the interface with other BeFEMTO components and external applications (such as a video client). Upon receipt of Service Error events, they request a diagnosis agent to start the corresponding diagnostic process. They are also in charge of propagating the resulting report to other BeFEMTO components, external applications and to a storage agent located in the LFGW.

Diagnosis agents: They receive diagnosis requests, together with observations made for the diagnosis. A diagnosis agent will create its own BN and fill it with the evidences it has for a given request and the related evidences it may have in its cache. It will then infer the new probabilities and try to come up with the cause of the problem. If it needs further evidences to make a conclusion, it may request further observations/evidences to Observation Agents. In addition to observations, they may also request for beliefs on a particular Bayesian node state to a Belief Agent.

Observation agents: They provide observations by performing specific tests upon request. Each Observation Agent will be specialised in a certain type of tests.

Belief agents: They provide a belief on a certain node state. Like the diagnosis agents, they also have their own BN and make Bayesian inference to obtain the requested belief. The main difference with diagnosis agents is that they just provide the belief of the state of a particular node, rather than providing a whole diagnosis report.

Knowledge agent: The knowledge agent is in charge of distributing diagnosis knowledge to all interested agents and performing Self Learning by processing the results of past diagnoses. Note Self Learning capabilities will be implemented in future iterations.

Storage agent: This agent is in charge of storing diagnosis reports in a central repository to allow for self-learning.

Besides the agents previously mention, the BeFEMTO Fault Diagnosis functionality provides:

A web interface to view all diagnosis reports performed in BeFEMTO, as well as to manually validate them.

A common repository database where all diagnosis reports will be stored, as well as relevant information about their validation.

In order to implement the multi-agent diagnosis system introduced in the previous sections, the WADE/JADE multi-agent architecture will be used. This agent platform can be distributed across machines (which not even need to share the same OS). JADE is FIPA compliant and it is completely implemented in Java language. The minimal system requirement is the version 1.4 of JAVA (the run time environment or the JDK). WADE environment runs on top of JADE and provides important mechanisms to manage the deployment of the system agents.

In BeFEMTO, most JADE agents will be located in the LFGW, and may run on top of OSGi to ease deployment (this is no still finally decided). To do so, the JADE-OSGi bundle provided by JADE, which creates a JADE container on top of OSGi, should be installed. Then, other bundles can be used to start/kill agents within the JADE container. Another option would be to deploy agents as a WADE application to take advantage of the advance features provided by WADE but, in this case, there would be no support for OSGI.

### 4.1.3 Fault Diagnosis processes

#### 4.1.3.1 Discovery process

The JADE platform offers a Directory Facilitator (DF) agent where agents register the services that they offer. Then, any agent wishing to discover other agents can interrogate the DF agent. Therefore, diagnosis, observation, belief and knowledge agents in BeFEMTO have to publish their services in the DF. Upon start-up, these agents register to the DF indicating the particular Enterprise Network they belong to.

Diagnosis agents publish the diagnosis service offered; observation agents publish the observation service provided together with all the possible observations and their cost. Finally, belief agents publish the belief service provided together with the belief node name they can obtain and its cost.

#### 4.1.3.2 Diagnosis process

Fault Diagnosis is designed as a recursive algorithm. BEFEMTO will iteratively perform Bayesian inference with available information until it reaches a given predefined confidence level. Otherwise, it tries to gather additional information by performing tests. From all the set of possible tests that can be performed, the one with the largest difference between value and cost is chosen first. The value and cost parameters are described below:

Value of the test: A test has a large value if its result implies a significant increase in the confidence of the diagnosis. The value of a given test may change depend on the scenario at hand.

Cost of the test: It represents the cost in terms of resources used, time consumed to perform such a test, price, etc.

The steps of the diagnosis process are as follows:
1) It adds to the BN all stored observations related to it.

2) It performs Bayesian inference

3) It checks if the confidence in the diagnosis is high enough. If so, it stops the diagnosis. Otherwise it goes to 4.

4) It selects the most appropriate action (test, request observation or request belief). Depending on the action to be done:

    If the given action is to perform a test, it performs it and returns to 1.

    If the given action is to request an observation, it sends the request, and waits a maximum for an answer. When receiving a proper answer, it goes to 5

    If the given action is to request a belief, it sends the request, and waits a maximum time for an answer. When receiving a proper answer, it goes to 5

5) It waits for an answer. When receiving one, it evaluates the response to continue the diagnosis, going to 1 again.

6) If no more actions are found, it finishes the diagnosis procedure. This can happen because it cannot get more evidences or beliefs, or because the time to perform the diagnosis has been exceeded.

#### 4.1.3.3 Self-learning process

Since service operators usually have deep knowledge of the problem domain, they should initially define the structure and Conditional Probability Tables of the BN used to diagnose each service. However, an important challenge for fault diagnosis is being able to improve its diagnostic intelligence over time and to automatically adapt it to the particularities of each Enterprise Femtocell Network deployment.

Therefore, self-learning capabilities have been incorporated to the Fault Diagnosis functionality. For that purpose, diagnosis information is stored in a repository and is further processed in order to train the BN. This training requires diagnosis results to be validated, indicating whether they were right or wrong.

Although the ideal solution would be to somehow automatically validate diagnosis results, manual validation is assumed for the time being.

### 4.1.4  Fault diagnosis database design

The following figure depicts the relational diagram for the design of the database used to store and validate diagnosis operations.



**Figure 4.1: Fault diagnosis database tables.**

Following is an explanation of each of the tables used:

**kow_operation:** Allows storing diagnosis operations. Only general information is stored here.

operation_id: The diagnosis operation identifier.

scene_id: The identifier of the scene associated to the operation. In general, a scene is the same as a BN.

scene_version: The version of the associated scene.

timestamp_start: Timestamp that marks the start of the operation.

timestamp_end: Timestamp that marks the end of the operation.

operation_status: The status of the diagnosis ('terminated', 'terminated_with_critical_error', 'terminated_with_warning', 'in progress')

**kow_operation_belief:** Allows storing the beliefs involved in one operation.

operation_id: The diagnosis operation identifier

belief_type: The type/name of the belief. In a BN, the name of a belief node (e.g. ProviderHANCongestion).

belief_value: The logical value of the belief.. This attribute is part of the PK because all possible belief values are to be stored, each one with its associated probability.

belief_probability: The probability associated to the belief value (e.g. 0.8596575). Min value = 0.0; Max value = 1.0

**kow_operation_observation:** Allows storing the observations involved in one operation.

operation_id: The diagnosis operation identifier

observation_type: The type/name of the observation. In a BN, the name of an observation node.

observation_value: The value of the observation..

timestamp: Timestamp that marks when the observation was taken.

additional_info: Any additional information that we may want to store.

**kow_operation_validation:** Allows storing if a diagnosis operation has been validated. If it has been validated, the diagnosis result (correct or incorrect) is stored. If the diagnosis result was incorrect, the real cause of error can be stored.

    operation_id: The diagnosis operation identifier

    diagnosis_result: The result of the diagnosis operation (null, 'correct' or 'incorrect'). A null value means that the operation has not been validated yet.

    real_cause: If the diagnosis result is 'incorrect', the real cause of the problem can be stored here.

**kow_operation_parameter:** Allows storing operation parameters. A possible parameter type is "AlternativeOperationID".

    operation_id: The diagnosis operation identifier

    parameter_type: The type/name of the parameter (e.g. AlternativeOperationID).

    parameter_value: The value of the parameter.

**kow_operation_error:** Allows storing information about errors occurred when performing a diagnosis operation.

    id_operation: The diagnosis operation identifier

    error_source: The source of the error. Maybe the name of the agent where the error occurred.

    error_details: The details of the error.

    level: The level of criticism of the error (critical and warning).

    timestamp: Timestamp that marks when the error occurred.

### 4.1.5  Information model: BEFEMTO Fault Diagnosis Ontology

BEFEMTO Fault Diagnosis Ontology (BEFEMTO_FDO) is the basis for agent communications. Agents communicate between them using ACL Messages. The content of these messages can be free text or any kind of serialization of structured data as text. For simple communications, it is possible to define an easy convention so that agents can understand each other. For complex messages, it is more convenient to set up structured data formats and a mechanism to marshall and unmarshall these messages from Java code: FIPA-SL, XSD/XML, RDF/N3, etc.

The Jade platform provides a mechanism to define this structured data as Java classes and an interface to serialize/de-serialize these objects into FIPA-SL language. It is possible to code Jade ontologies in Java. However, it is a messy task so BEFEMTO_FDO has been defined using the Protègè and Ontology Bean Generator tools.

BEFEMTO_FDO defines three groups of classes in order to represent bayesian knowledge, information about current operations and information about the actions that agents perform. These classes are described in the next subsections.

#### 4.1.5.1  Bayesian Knowledge Classes

BEFEMTO agents must have knowledge about what evidences can be acquired, which hypothesis can be achieved, and their dependencies in terms of conditional probabilities. This knowledge is, in fact, equivalent to the definition of a BN. Since different scenarios must define different BNs, fault diagnosis ontology only defines general terms about the structure of the Bayesian knowledge and the specific hypothesis, evidences, and probabilities are defined as instances in this ontology.

The following classes are defined:

BNKnowledge: a container to store information related a specific BN that will be used by an agent with diagnosis capabilities. Besides the different components of a BN (hypothesis nodes, observation nodes, Conditional Probability Tables, etc), it includes information about the BN name (or scenario), and version. This name and version will be used within the knowledge update mechanism. Also, if the agent using this container is a diagnosis agent, it will publish its capacity to diagnose that scenario in the DF.

ObservationNode: the class of all the types of evidences that can be obtained by an agent. It includes information about all possible values, the time during which an observation of this type is still valid, the cost of obtaining this kind of observations and the information that it is required in order to carry out this observation.

HypothesisNode: the class of all the possible diagnostics for a given scenario. It includes information about their types and their possible values. It may be related to several Thresholds.

Thresholds: the class of all the thresholds that can apply to a particular value of a given hypothesis.

CPTNode: instances are the conditional probabilities of each node (hypothesis or observer).

Bridgenode: a class that defines auxiliary observation nodes just are used to exchange beliefs between agents as opposed to change concrete evidences.

Agent: is the class of all the types of agents that exist, and include information about what kind of actions they can perform.



**Figure 4.2: Bayesian Knowledge Classes Diagram.**

All agents with diagnosis capabilities (diagnosis and belief agents) create upon start up an instance of the BNKnowledge with the BN that models their view of the problem to be diagnosed.

### 4.1.5.2 Operational Classes

The operational classes represent all the information that agents use and share regarding a specific diagnosis process or operation. All the information that is gathered or deduced is aggregated with these classes. The following classes are defined:

Operation: aggregates all the information related to an operation. It includes basic information for the operation (identifier, timestamp, etc.) together with specific parameters, observations, found errors, etc.

Observation: the class of all the evidences that can be created by the agents. Their type and value corresponds to one of the possible values defined by an observer. They also include the timestamp of the evidence and a free text details field.

BeliefProbability: represents the certainty of an hypothesis value to be true.

Belief: aggregates the probabilities of all possible values of an hypothesis.

OperationParameter: its instances are pairs of parameter name & value gathered for this operation.

OperationError: its instances are created whenever an error occurs during the diagnosis process.

**Figure 4.3: Operational Classes Diagram.**

### 4.1.5.3 Action Classes

In Jade ontologies mechanisms, only classes that inherit from ContentElement can be serialized into FIPA-SL: AgentAction, IRE and Predicate. In BEFEMTO_FDO, some actions have been defined so that agents can request and inform other agents about different things.

The main actions, related to the diagnosis procedure, are:

RequestObservation. This action allows an agent to request an observation to another agent.

RequestDiagnosis. This action allows an agent to request a diagnosis to another agent.

RequestBelief: This action allows an agent to request a belief to another agent.

InformObservation. This action allows an agent to send the results of one or more observations.

InformDiagnosis. This action allows an agent to send the result of a diagnosis to another one.

InformBelief. This action allows an agent to send the result of a belief it has about the state of a particular node.

Next figure shows the class diagram for the previous actions.

**Figure 4.4: Class diagram for actions related to the diagnosis procedure.**

Also, the BeFEMTO Fault Diagnosis ontology defines actions related to the distribution of the BN knowledge:

StoreDiagnosis. This action allows an agent to request storing in a shared repository all the information related to a given diagnostic operation.

DistributeKnowledge. This action is used to send the Bayesian knowledge for a given scenario to another agent.



**Figure 4.5: Class diagram for actions related to knowledge distribution.**

### 4.1.6 Communication between agents

BeFEMTO agents communicate using the ontology previously described. Each message is composed of two parameters:

- Action: It describes the type of action carried in the message. It can be any of the Actions defined previously (RequestDiagnosis, RequestBelief, RequestObservation, InformDiagnosis, …)

- Operation: All operations contain an OperationId used to correlate answers with requests. The extra information carried within the operation depends on the type of Action.

The following figure shows a generic message flow triggered by an Interface Agent upon receipt of a ServiceErrorEvent. The agent delegates the diagnosis to a diagnosis agent via a RequestDiagnosis action. The diagnosis agent may request observations to observations agents, or even beliefs to a belief agent. A belief agent could also request observations from other observation agents. The flow diagram just shows the main parameters contained within the operation part of the message.

**Figure 4.6: Message flow for Fault Diagnosis.**

### 4.1.7 Bayesian Network modelling

The diagnosis functionality uses Bayesian inference to come up with the most probable cause of error. Therefore, a BN has to be defined, modelling the relationship between the evidence that can be obtained and the different hypothesis.

The main source of information for the diagnostic process identified so far as are follows:

Status information provided by the femtonodes via management interfaces.

Status information provided by the Iuh-tap component

Test agents developed to provide specific information such as connectivity status between different network domains.

### 4.1.8 Interfaces with other components

The Fault Diagnosis will mainly interface with the actual femtonodes in the Enterprise Network and with the Iuh-tap component sitting in the LFGW. Although the actual specification of these interfaces is not finished, the initial thinking is to base it on the use of asynchronous notifications. Thus Fault Diagnosis would subscribe to a set of relevant status events and would in turn publish the result of each diagnostic process so all subscribed components will be properly notified.

In particular Fault Diagnosis will subscribe to Service Error Problems detected by BeFEMTO that will act as trigger of the diagnostic process. Service Error Events must convey the type of error and other specific information about the detected problem.

The Diagnosis Report Event published by Fault Diagnosis must contain all information received in the ServiceErrorEvent that triggered the diagnosis process plus all observations gathered by fault diagnosis and a list of beliefs showing the root causes that exceeds a configured threshold, together with their probabilities. For each belief, both the beliefType (and name of the hypothesis of the root cause of the error), the beliefValue (or name of the state of that belief) and the associated probability will be provided.

### 4.1.9 User interface

The results of the diagnosis process are stored in a Data Base installed in the LFGW server. To allow LFGW operators accessing information about past diagnoses, a graphical user interface has been implemented which provides the following capabilities:

Query diagnosis results. It lists the diagnosis information stored, filtered according to a set of parameters introduced by the user. Only the main parameters of each diagnosis result are shown.

Show diagnosis details. It provides all relevant information related to a given diagnosis result. For instance, all the evidences used in the diagnosis process, the value of each of the hypothesis, etc.

Validation of diagnosis results. By means of this option, an operator can mark diagnosis results are correct or incorrect. In case of incorrect diagnosis, the operator can also introduce the real cause by means of the user interface. Note this functionality is needed to enable self-learning in the fault diagnosis application.

## 4.2  Enhanced Power Management in Femtocell Networks

Powering a femtocell without a single user attached to it is a waste of energy – even more if it is a femtocell network with many femtocells consuming energy. A smart power management procedure can shutdown one or more femtocells in times when they are not used and reactivate it when required. Depending on the capabilities of the femtocell, various power save modes can be applied. The procedure can be divided in two sub-tasks: 1) activation of femtocells required to serve UEs coming into communication range, and 2) putting femtocells into a power sleep mode as soon as they are not required anymore.

In the enterprise network scenario, such power management is provided by functions on the LFGW. In year one of the project we have devised a mechanism to coordinate power management of femtocells when all users have moved to the macro cell network. As soon as the last UE has left the communication range of the femtocell network, the femtocells are put into power sleep mode. A power management application hosted in the mobile core or on the LFGW (if such an entity exists), wakes up the femtocells when a UE which is authorised to attach to the femtocell network is getting close to the femtocell communication area again. In year two, we have extended this scheme to in-house power saving mechanisms. These enhancements are described in the next section together with the original scheme, followed by a simulative evaluation of this scheme.

### 4.2.1  Enhanced power management in femtocell networks

A trivial power management system can be aligned to office hours of an enterprise switching off the femtocell network during off-office hours. A more advanced power management scheme could utilize output power adaptation based on the scheme described in [91]. This scheme adapts the output power of a femtocell in the centre of a femtocell network to cover the area of surrounding femtocells during very low utilization. These surrounding femtocells are put into a power save mode until higher demand arises again.

Our power management scheme based on UE tracking on the **macro network** works as follows: Femtocells have to be active only if at least one UE, authorized to attach to the femtocell, is within the communication range of the femtocell. As soon as all these UEs have left the communication range of the femtocell network and are not likely to move back into this area soon, the femtocells can be put into a power-saving state. To be more precise, tracking of all UEs is not required but only the UE which is closest to the femtocell. One way to realize this is to have a specialized function on the UEs which can signal to the power management application to indicate that it has left or entered a specific area. The management application directs the femtocells in this network to power on when the first *enter* signal has been received. For each additional *enter* signal a counter is incremented to keep track of the number of UEs within this area; likewise the counter is decremented for each *leave* signal received. This modelling provides enough information to the power management application while preserving the privacy of the users – the identity of the UE is not required in this scheme at all.

The same observation can be exploited to implement a **femtocell network internal** power management scheme. In an office building, employees usually have a fixed routine which halls to walk down and which rooms to enter; one reason can be that the employee has only access to certain areas within the building or there is no good reason to divert from this routine during normal business operations. Thus, a power saving scheme which has powered down all femtocells during off office hours could exploit the above mentioned procedure to detect the first employee coming to work or it simply switches on a single femtocell in the building's entrance area. Each time an employee comes to work his UE will perform a regular hand-over procedure from the macro network to the femtocell in the entrance area and authenticate itself in this process. Based on this UE identity the LFGW can determine the areas in the building the employee can access and the most likely path the employee is going to take. Depending on the current status of the UE, active or idle, the LFGW will activate a subset of the femtocells within the building. If the UE is idle when the employee enters the building, the LFGW will only activate the femtocell in the employee's office (the most likely place any employee will go to in the morning); if the

way to this place is rather long, some more femtocells on this path could be activated in case the employees starts using his phone while walking. Is the UE active when the employee enters the building, the LFGW will instruct all femtocells on the most likely path to get fully operational to provide continuous connectivity for the employee; femtocells in all other areas where the employee has access to but are not on the most likely path are put into a power state which allows for very quick activation (depending on the femtocells' power saving capabilities).



**Figure 4.7: Activation of femtocells when employee A comes to work**

A further refinement of the management system takes additional external data into consideration to decrease the power consumption further. For example, if employee A and B have a joint meeting in the near future, indicated in their respective calendars, the likelihood that employee B arrives shortly after employee A is high. Similar correlations can be found by analysing historic information, i.e., which users come to work together, or at nearly similar times. This allows for configurations as depicted in Figure 4.8 where a femtocell at the end of the path is selected to provide connectivity to more than a single user.

Correlating the arrival time of the user with his calendar provides the additional benefit of predicting deviations of the user's daily routine. If the user enters the building briefly before his first scheduled meeting in a conference room, the path prediction function can map the way to this conference room in addition to the path to the user's office. Thereby, continuous connectivity can be provided on either path the user takes.



**Figure 4.8: Predictably select a network node which will serve more users simultaneously.**

Combining the path prediction with output power adaptation (see [91]) provides for even higher energy savings. Figure 4.9 shows an example where the configuration of the femtocell network requires the capability of the femtocell to change its output power to increase its transmission area. The management system is expected to have prior information regarding the power consumption of a femtocell using different output power settings. Based on such information the management system can compute the total energy consumption for different configurations and apply the configuration that is more cost efficient in terms of power consumption.

**Figure 4.9: Utilizing output power adaptation capabilities of the network node**

**Power saving management algorthim**

The proposed algorithm is centralized with prior access to information regarding

- enterprise network: location and coverage compensation of each network node
- monitoring methods of user activity
- information regarding arriving users
- additional user specific data sources (electronic calendars)
- history information

It is worth noting that the history information is referred to record that are kept by the management system regarding the arrival or leaving times of each user, his load activity, mobility inside the building.



**Figure 4.10: Overview of the energy management algorithm for powering on cells.**

In the following description we consider the most complex embodiment of our scheme, where certain locations are covered by several wireless access points at the same time and the management system is equipped with history information about user activity and arriving/leaving patterns. The management system monitors the user's arriving or leaving patterns and performs accordingly the energy management

of the femtocell enterprise network. The management system is supposed to maintain a sorted list of user arriving in the network, which is based on history and electronic calendar information. An overview of the energy management process, for arriving users that enter the enterprise network is illustrated in Figure 4.10.

When a user enters the building, he is removed from the arriving list, which is also updated. The management system then identifies based on given information the set of cells that may serve the user in terms of coverage and examines whether any of these cells are powered-on already. In case none of the identified cells is powered-on, the management system examines the user load activity to determine the optimal cell to power-on. If more than one cell can support that particular user, the next forthcoming user within the same set of cells is also considered, until only one cell is specified that may support the user that arrived and the maximum set of forthcoming users.

The reason is the support of forthcoming users limits the need of operating additional cells, while at the same time it ensures stability limiting the handover of users from a certain cell to another. In case a cell from the selected set is already powered-on, the system checks whether it is feasible to add the new user and in a positive case it performs the corresponding updates. Otherwise, it needs to examine whether powering on another cell instead may serve all current users including the entering user. An overview of such search process is illustrated in Figure 4.11.



**Figure 4.11: Overview of the process that specifies whether another set of cell may be able to accommodate all current users including the new incoming one.**

To do so, the current users that camp at cells within the coverage range of the new user are examined to identify the potential of being shifted to another common cell to all users that could accommodate their load demand. If a common cell does not exist, the system abandons the process. Otherwise, the maximum load user is first selected with the management system identifying the set of cells that may provide coverage to such user and at the same time may provide coverage to all other affected users. From that set of cells the management system selects the least loaded one.

If the selected cell can accommodate the first user, the next user is considered until all users are examined. In case a user cannot be accommodated by the selected cell, the process also is terminated. Otherwise, at the end of the process the system switches on the selected cell and instructs the affected users to perform a handover, instructing the previous cell to power-off. In case none other cell is

identified, the management systems need to select a new cell to power-on using the same process as in the case of an incoming user with none cells being powered-on.

When a user leaves the network, the system adds the user to the list and performs the list history related updates, before identifying the cell affected from the user absence. An overview of the energy saving management process for the case of a user leaving the network is illustrated in Figure 4.12. In particular, when a user leaves the enterprise network, the management system checks whether more users are camping in that particular cell, and in a negative case it simply switches off such cell and performs the appropriate updates.



**Figure 4.12: Overview of the energy management algorithm for powering off cells.**

Otherwise, it examines whether an alternative already powered-on set of cells can accommodate the remaining users in order to hand them over and power-off the previously operating cell. In case of a positive case the management system instructs a handover procedure before adapting the power configuration of the affected cell and performing the appropriate updates else, it just performs the appropriate load related updates without powering off any cells and continues to monitor the network. The process of identifying the appropriate set of powered-on cells that may accommodate the existing users is similar to the process described in Figure 4.11.

It should be noted that the power configuration of a femtocell is adapted, the management system checks whether such cell is within a geographical range where users may be need it back in operation again in case they deviate from their common directions and behavior. Based on such assessment, the management system either switches off the selected cells completely or maintains them within an easily powered back state.

The proposed method blends both arriving and leaving cases with the ultimate objective to provide a stable solution that minimizes energy consumption for enterprise networks.

### 4.2.2 Simulative Evaluation

We evaluated our proposed energy saving scheme in simulations. Figure 4.13 shows the office layout used. The femtocell at the reception is considered always on. All other femtocells are turned off when all employees in the area around a femtocell have left.



**Figure 4.13: Scenario layout for simulative evaluation of power management scheme**

Based on the proposed scheme we have developed three algorithms to activate femtocells when employees are entering the office. The first and most simple algorithm turns on the office femtocell of the employee detected at the reception and all femtocells in the hall leading to this office. The second algorithm refines this behaviour by checking whether a neighbouring femtocell which has already been activated can host the just arriving user. An office femtocell is only turned on when the communication requirements of the user exceed the capacity neighbouring femtocells can provide. The third and last algorithm is further refining this algorithm by predicting the employees' arrival order and based on that turning on femtocells. The prediction of the employee arrival order is based on historic data – when do the individual employees usually come to work in the morning. The arrival prediction could be improved by taking additional information sources, e.g., the employees' digital calendars, into account. Based on the predicted arrival order the algorithm optimizes the femtocell activation by empowering the femtocells which can host most of the arriving employees first. In order to simulate this prediction we have defined a fixed order of arrival which is fed into the algorithm. The actual arrival of employees is varied using a Gaussian distribution around the defined arrival time plus some small probability that the employee does not come to work this day.

Figure 4.14 shows simulation results for the three algorithms just introduced. On the X-axis the office femtocells are depicted and on the Y-axis the power saved in comparison to a simple time based, i.e., turning on all femtocells at begin of business, power saving scheme. Hence, the higher the bars the more power is saved by this implementation of the power management scheme.

**Figure 4.14: Simulation results for prediction based power management scheme**

# 5. Security

## 5.1 Secure, Loose-Coupled Authentication of the Femtocell Subscriber

### 5.1.1 Work during Year 1

The work carried out within BeFEMTO during year 1 has presented for a standalone scenario a new model for Fixed Access Network Authentication over a FTTH backhaul, in a way that is decoupled from the physical elements of the access network infrastructure, secure, compatible with current standards, and near the market.

Current fixed access networks, regardless of the access technology, are natural heirs of traditional PSTN lines. This is probably the reason why the client identification (or subscriber, or the person who holds the contract with the Telco operator) is based, generally, on parameters directly linked with physical elements of the access network. An example is the BRAS (Broadband Remote Access Server) port at which the virtual ATM circuit (or Ethernet VLAN) of the DSLAM is terminated (in xDSL deployments). In this case, such port depends on physical parameters (BRAS identifier, card number and connector) as well as on logical ones (e.g., ATM virtual circuit identifier). Identification in FTTH deployments works similar to that.

The usage of these identification schemes for Femtocell deployments leads to highly complex systems, since they require a system able to control each network element between the Femtocell node and the IP terminator (BRAS).

This complexity, in turn, translates into highly complex provision systems, with long and prone-to-failure workflows, to run all the actions needed to configure the network to deliver a new service. The operative costs of this model are enormous, and in a short/mid-term environment where users will want to use their services the very first moment they buy them, it will prove unable to cope with this prerequisite.

Our target is to find a mechanism to identify the subscriber of the Femtocell to the fixed access line that acts as a backhaul for the Femtocell node, in a way that is decoupled from the physical elements of the access network infrastructure, secure, compatible with current standards, and near the market (due to the high costs that a complete change of the access network infrastructure would imply).

**UICCs supporting EAP-AKA** authentication has been selected as the storage point for the subscriber credentials. A pluggable UICC offers enough security guarantees to the end user and provides other advantages related with mobility scenarios (a subscriber could extract the UICC card from a Femtocell node and insert it in another, triggering an authentication process in which the backhaul would be configured identically than if he were at home).

The algorithm developed covers the following situations:

Initial attachment: The user inserts the UICC card for first time.

Full Re-authentication: IMSI is used and fresh authentication vectors are generated to avoid security threads.

Fast Re-authentication: TIMSI is used to avoid overloading the network and detect whether or not the UICC card is removed.

### 5.1.2 Solution extension for Networked Femtocell Scenario

The architecture already studied defines the Authentication of the Femtocell Subscriber in a residential environment when one subscription line and one Femtocell is involved. However, this model is not appropriate for an enterprise environment where the Femtocell has to provide service to a larger number of users. In this case, a networked Femtocell approach fits better. Regarding the scalability issues, two approaches are under study:

The first one implies having a single LFGW in charge of management issues and several Femtocells providing service to end users.

The second one implies having multiple LFGWs in charge of management issues and several Femtocells providing service to end users.

In both approaches, the Femtocells in charge of the management of the network of Femtocell will hold the connection to the wired backhaul and will thus run the authentication protocol for the subscription in a similar way to the procedure described in section 5.1.5 of deliverable D5.1 However, since we are in an enterprise environment, an additional requisite appears: backup line

**5.1.2.1 Network of Femtocell with centralized management**

As already mentioned, this solution consists in one central Femtocell (LFGW) in charge of management issues and several Femtocells providing service to end users. This LFGW would be the one holding the connection to the wired backhaul and should perform the activities related with subscriber authentication.

The next figure depicts the architecture of a network of Femtocell with centralized management.



**Figure 5.1: Network of Femtocell with Centralized Management Architecture.**

**5.1.2.2 Network of Femtocell with distributed management**

This solution consists in several LFGW in charge of management issues and several Femtocells pending from each LFGW providing service to end users. This solution is intended to support a large number of end users in an Enterprise in which, due to the equipment constraints, having a single LFGW is not feasible.

The authentication procedure in a first approach, where there is one subscription line, would be similar to the residential scenario. The LFGW would have identical functionalities to the Network of Femtocell with Centralized management scenario above described. However, some additional functionalities regarding to the management of the network of LFGW (traffic routing, handovers management, etc) should also be addressed.

The next figure depicts the architecture of a network of Femtocell with .distributed management.

**Figure 5.2: Network of Femtocell with .Distributed Management Architecture**

### 5.1.2.3 Backup line. Technologies involved

Enterprises often require an additional SLA for communication services. In case there are problems in the Broadband Access Router or power outage, a fall-back solution can be offered at least to provide some critical services.

The Backup line purpose is resilience rather than more ambitious targets such as load balancing or voice session continuity. Not all services may be available, or able to be given an appropriate QoS, when the backup interface is in use. It must be possible to specify which services will be maintained (for example, security alarms). Both the user and the remote management system need to be informed when there is a change of WAN interface in use. It is needed to be some hysteresis so that transient conditions do not cause frequent and unnecessary changes.

The function of a backup WAN interface is that it takes care of the most important services (SOS calls, banking etc.) when the primary interface is not working. Depending on operator's service agreement with the customer, two alternatives can be seen:

The backup interface is turned off (cold) during normal operation (when the primary interface is working). It is turned on when the primary interface stops working.

The backup interface is always up and running (warm) and the services can be moved to the secondary interface almost directly, when the Broadband Access Router decides to switch to the secondary WAN interface.

The IP connectivity procedure, basically means that L1, L2 and L3 are up and running (possibly without traffic, possibly with keep-alive), and that the Broadband Access Router has been assigned an IP address (either via DHCP or PPP) and IP-configuration by the BSP. To periodically check the IP connectivity, Ping can be used. Ping is a basic method that allows verifying that a particular IP address is assigned and allows communication, by sending Internet Control Message Protocol (ICMP) Echo Requests, and checking a reply to it.

The following flowchart shows the Dual WAN handling as indicated above.

**Figure 5.3: Backup WAN interface process flow**

The model presented in this document for the authentication of the Fixed Access Network over a FTTH backhaul must be compliant with several technologies and architectures besides the FTTH one in order to be able to switch to the backup line and perform the same authentication procedure.

The decision of the technology involved in the backup line will depend on several factors like the bandwidth that the critical services need to be supported or the compliance with the NGN access network infrastructure already in use for the FTTH technology. In the following section four technologies are under study: FTTH, xDSL, HSPA and LTE.

### 5.1.2.3.1 Backup technologies

#### 5.1.2.3.1.1 FTTH

A backup WAN interface using a FTTH backhaul is a solution that can not be considered only for critical service support. A secondary line in the ONT is feasible and only the line failure detection should be implemented because the authentication mechanism towards the backhaul would be the same.

Figure 5.4 shows the deployment for the backup line. This secondary line would be connected to a different port of the OLT. Therefore in case of failure in the primary line, the ONT would be able to switch to the secondary WAN interface and support all the services until the primary line is restored. The secondary line would be connected to another port of the OLT.

**Figure 5.4: FTTH as Backup Technology**

Unfortunately this architecture is not realistic in a real deployment of FTTH based on GPON. A real deployment is not point-to-point based but includes splitters to deliver the signal from the OLT to several ONT's. Figure 5.5 depicts the real deployment of GPON FTTH backhaul.



**Figure 5.5: FTTH Access Network Architecture**

Deploying a FTTH GPON backup WAN line implies that both lines would be likely connected to the same $2^{nd}$ level splitter or at least to the $1^{st}$ level splitter. As a consequence the primary and the backup line

would be connected to the same GPON port in the OLT and therefore, in case of failure in the access network, both lines would be affected.

### 5.1.2.3.1.2   LTE

A backup WAN interface using LTE technology is a very interesting choice because none only critical services could be supported, allowing to overcome seamlessly failures in the primary line subscription.

From a corporate perspective, LTE promises much and has several key characteristics. It has very high theoretical data rates (100Mbps downlink and 50Mbps uplink), if spectrum is available. This makes it a viable alternative to large corporations Wi-Fi deployments.

It has very low latency (5msec for small packets) and will be an all-internet protocol (IP) network for voice and data.

The all-IP nature of LTE, combined with low latency, may make LTE particularly suitable for applications that involve multiple types of media streams or multimodal interactions, such as communication-enabled business processes.

Long Term Evolution (LTE) will be a widely adopted next-generation cellular technology and it is envisaged that it will reach 50% of subscribers in the US and Western Europe by 2015.

However, despite the promises, in most countries, LTE coverage will be incomplete until 2020, and actual bandwidth will likely be around just 10-20% of the theoretical peak.

Figure 5.6 depicts the architecture proposed for a LFGW/Broadband Access Router holding two subscription lines. The primary subscription throughout the fixed fibre backhaul, and the backup subscription throughout the LTE backhaul. The architecture proposed for the enterprise environment includes the Optical Network Terminator (ONT) and the Local Femtocell Gateway + Broadband Access Router both integrated as a single box model.



**Figure 5.6: LTE as Backup Technology**

#### 5.1.2.3.1.2.1   Functional Description

**eNodeB**: The only node in the E-UTRAN is the E-UTRAN Node B (eNodeB). The eNodeB is a radio base station that is in control of all radio related functions in the fixed part of the system. Base stations such as eNodeB are typically distributed throughout the networks coverage area, each eNodeB residing near the actual radio antennas.

Functionally eNodeB acts as a layer 2 bridge between UE and the EPC, by being the termination point of all the radio protocols towards the UE, and relaying data between the radio connection and the corresponding IP based connectivity towards the EPC.

The eNodeB is also responsible for many Control Plane (CP) functions. The eNodeB is responsible for the Radio Resource Management (RRM), i.e. controlling the usage of the radio interface, has an important role in Mobility Management (MM), it controls and analyses radio signal level measurements carried out by the UE, make similar measurements itself, and based on those makes decisions to handover UEs between cells.

**MME**: Mobility Management Entity (MME) is the main control element in the EPC. Typically the MME would be a server in a secure location in the operator's premises. It operates only in the CP, and is not involved in the path of UP data.

In addition to interfaces that terminate to MME in the architecture as shown in Figure 5.6, the MME also has a logically direct CP connection to the UE, and this connection is used as the primary control channel between the UE and the network. The following lists the main MME functions in the basic System Architecture Configuration:

Authentication and Security

Mobility Management

Managing Subscription Profile and Service Connectivity

**SeGW**: In the Basic System Architecture configuration, the high level function of S-GW is UP tunnel management and switching. The S-GW is part of the network infrastructure maintained centrally in operation premises.

The SGW routes and forwards user data packets, while also acting as the mobility anchor for the user plane during inter-eNodeB handovers and as the anchor for mobility between LTE and other 3GPP technologies (terminating S4 interface and relaying the traffic between 2G/3G systems and PGW). For idle state UEs, the SGW terminates the DL data path and triggers paging when DL data arrives for the UE. It manages and stores UE contexts, e.g. parameters of the IP bearer service, network internal routing information. It also performs replication of the user traffic in case of lawful interception.

**P-GW**: Packet Data Network Gateway (P-GW, also often abbreviated as PDN-GW) is the edge router between the EPS and external packet data networks. It is the highest level mobility anchor in the system, and usually it acts as the IP point of attachment for the UE. It performs traffic gating and filtering functions as required by the service in question. Similarly to the S-GW, the P-GWs are maintained in operator premises in a centralized location.

Typically the P-GW allocates the IP address to the UE, and the UE uses that to communicate with other IP hosts in external networks, e.g. the internet. It is also possible that the external PDN to which the UE is connected allocates the address that is to be used by the UE, and the P-GW tunnels all traffic to that network. The IP address is always allocated when the UE requests a PDN connection, which happens at least when the UE attaches to the network, and it may happen subsequently when a new PDN connectivity is needed.

5.1.2.3.1.3   High Speed Packet Access (HSPA)

A backup WAN interface using HSPA technology is a possibility of special interest due to the support of critical services can be ensured when using a mobile secondary interface that acts as a complement to the primary (that is considered to be a high speed interface, i.e. Fibre).

A mobile backup WAN interface using HSPA (High Speed Packet Access) can be integrated into the Broadband Access Router, attached to the USB or to one of the Ethernet ports as a module (dongle). When attaching the mobile WAN module to the Ethernet or USB interface, these ports on the Broadband Access Router need to be configured as WAN port.

This section tries some aspects of HSPA (3G) mobile backup WAN interface. Table below summarizes typical and theoretical data rates that can be expected for a HSPA (3G) mobile backup WAN interface. Note that a HSPA module normally has fallback to WCDMA (3G), EDGE (GSM) and GPRS (GSM).

| Packet data service | | Theoretical max data rate | Typical data rate |
|---|---|---|---|
| HSPA | Upload | 2.0* Mbps | 300-800 kbps |
| | Download | 7.2** Mbps | 600-5000 kbps |
| WCDMA | Upload | 384 kbps | Over 300 kbps |
| | Download | 384 kbps | Over 300 kbps |
| EDGE | Upload | 118 kbps | 50-60 kbps |
| | Download | 236 kbps | 100-130 kbps (with bursts over 200 kbps) |
| GPRS | Upload | 43 kbps | 20 kbps |
| | Download | 86 kbps | 40 kbps |

*) Later releases will provide 5,67 Mbps
**) Later releases will provide 21 Mbps

The coverage is very much depending on the frequency band used in different parts of the world. In Europe the 3G (UMTS) is using 2100/1900 MHz and GSM is using 900 MHz. The data rates in the table are for HSPA module using one frequency carrier that is most commonly used today. But the standard allows (and some vendors have implemented it) the use of more than one frequency carrier, and then the data rates in the table need to be multiplied with the number of carriers.

Figure 5.7 depicts the architecture proposed for a LFGW/Broadband Access Router holding two subscription lines when HSPA is employed as backup technology.



**Figure 5.7: HSPA as Backup Technology**

### 5.1.2.3.1.3.1 Functional Description

**NodeB:**

Node B is the term used within UMTS to denote the base station transceiver. It contains the transmitter and receiver to communicate with the UEs within the cell.

**Radio Network Controller (RNC):**

This element of the radio network subsystem controls the Node Bs that are connected to it. The RNC undertakes the radio resource management and some of the mobility management functions, although not all. It is also the point at which the data encryption / decryption is performed to protect the user data from eavesdropping. In order to facilitate effective handover between Node Bs under the control of different RNCs, the RNC not only communicates with the Core Network, but also with neighbouring RNCs.

**MSC:**

The MSC provides the interface between the radio system and fixed network, performing all necessary functions to handle CS services to and from mobile terminals. As such, an MSC will interface with several base stations.

In effect it is an exchange which performs switching and signalling functions for mobiles within its designated area of control. It needs to take into account the allocation of radio resources and the mobile nature of users, which impact the location registration & handover between cells.

**Serving GPRS Support Node (SGSN):**

As the name implies, this entity was first developed when GPRS was introduced, and its use has been carried over into the UMTS network architecture. The SGSN provides a number of functions within the UMTS network architecture.

Mobility management When a UE attaches to the Packet Switched domain of the UMTS Core Network, the SGSN generates MM information based on the mobile's current location.

Session management: The SGSN manages the data sessions providing the required quality of service and also managing what are termed the PDP (Packet data Protocol) contexts, i.e. the pipes over which the data is sent.

Interaction with other areas of the network: The SGSN is able to manage its elements within the network only by communicating with other areas of the network, e.g. MSC and other circuit switched areas.

Billing: The SGSN is also responsible billing. It achieves this by monitoring the flow of user data across the GPRS network. CDRs (Call Detail Records) are generated by the SGSN before being transferred to the charging entities (Charging Gateway Function, CGF).

**Gateway GPRS Support Node (GGSN):**

Like the SGSN, this entity was also first introduced into the GPRS network. The Gateway GPRS Support Node (GGSN) is the central element within the UMTS packet switched network. It handles inter-working between the UMTS packet switched network and external packet switched networks, and can be considered as a very sophisticated router. In operation, when the GGSN receives data addressed to a specific user, it checks if the user is active and then forwards the data to the SGSN serving the particular UE.

**Home location register (HLR):**

This database contains all the administrative information about each subscriber along with their last known location. In this way, the UMTS network is able to route calls to the relevant RNC / Node B. When a user switches on their UE, it registers with the network and from this it is possible to determine which Node B it communicates with so that incoming calls can be routed appropriately. Even when the UE is not active (but switched on) it re-registers periodically to ensure that the network (HLR) is aware of its latest position with their current or last known location on the network.

#### 5.1.2.3.1.3.2   PDP Context and QoS

For a given frequency carrier, one or multiple primary PDP Contexts can be created. To each primary PDP Contexts a secondary PDP Context also can be configured. A PDP Contexts connection can be seen as a layer 3 connection between the Broadband Access Router and the GGSN over the mobile network. To set up the connection, APN (Access Point Name), APN username (optional) and password (optional) and PDP authentication need to be configured. Normally these parameters are taken from the UICC card. Depending on the number of PDP Contexts used and the maximum and guarantied bit rates, the Broadband Access Router needs to adapt its internal bandwidth and QoS algorithm accordingly.

#### 5.1.2.3.1.3.3   VPN Termination and L2TP Tunnelling

A HSPA backup WAN interface can be configured to use an IPsec/VPN tunnel per PDP Context. An IPsec VPN serves as a point-to-point tunnel interface allowing the Broadband Access Router to send some or all of its WAN traffic across an encrypted tunnel rather than in clear text.

A HSPA backup WAN interface supports L2TP tunnelling, providing mechanisms for tunnelling Ethernet frames between two peers over an existing network (usually the Internet). The LAN/WLAN is bridged into the tunnel using BCP over PPP over L2TP. The BCP/PPP/L2TP tunnel can be used to connect a branch office LAN to a corporate office LAN over a HSPA network connection.

### 5.1.2.3.1.4  xDSL

A backup WAN line using xDSL technology is a solution that may go beyond of supporting only critical services when the primary subscription falls down. Besides the very critical services such SOS calls or banking, some services considered important to the business development may need more bandwidth in order to operate properly or meet some QoS requirements.

The end-to-end ADSL network reference model is depicted in Figure 5.8.



**Figure 5.8: xDSL Architectural model reference**

#### 5.1.2.3.1.4.1  Functional Description

**DSLAM:**

The DSLAM is the concentration point for broadband and narrowband data. The DSLAM may be located at a central office or a remote site. The DSLAM serves as an ATM layer multiplexer/concentrator between the ATM Core Network and the Access Network. In the downstream direction it may perform routing/demultiplexing, while in the upstream direction it may perform multiplexing/concentration and higher layer functions, e.g., co-location with Core Network functions

The DSLAM contains a Core Network Interface Element that performs the ATM and PHY layer functions to interface to the ATM Core Network. Non-ATM core networks are not precluded. The VPI/VCI translation and higher-layer function performs the multiplexing/demultiplexing of the VCs between the Access Network interfaces (ATU-Cs) and the Core Network interface on a VPI and/or VCI basis. This block may also perform other higher layer protocol functions. The Access Network side ATM layer functions, if present, support the ATU-Cs, which terminate the Access Network lines in the DSLAM. If an ATU-C supports both 'Fast' and 'Interleave' channels two ATM TC sub layer functions may be needed.

**Broadband Remote Access Server:**

The BRAS interconnects the access nodes in a geographical area. This text does not assume ATM is the only transport technology, although ATM has been increasingly deployed in this infrastructure to provide broadband connectivity among COs. A service interworking function is required if the Regional Broadband Network includes multiple transport technologies such as ATM, Frame Relay or IP.

#### 5.1.2.3.1.4.2   L2TP Tunnelling

Two approaches are recommended for access to corporate networks. The first uses tunnelling (i.e., IPSec [69] or L2TP [70]) across an IP network, possibly the Internet to the corporate network. This approach replaces dial-up modems by using a virtual private network across the Internet. In the case of L2TP, which is referred to as a "compulsory" tunnelling model, the tunnel is created without any action from the user, and without allowing the user any choice in the matter. The L2TP tunnelling model can be on top of IPSec if the security is required.

The second approach is to use the Regional Broadband Network to provide direct high-speed connectivity to the corporate network. This has the advantages of being able to offer higher speed, QoS guarantees and greater security. Nevertheless, a "path" or virtual circuit between the user and the ingress service point MUST be established. The provisioning data SHOULD include the relationship between the service and the virtual circuit and the security policy in case of IPSec or tunnelling end-to-end.

#### 5.1.2.3.1.4.3   xDSL and NGN Architecture coexistence

An important aspect if we choose a wired based access technology such as xDSL for the backup line is the possibility that it could coexist with the current TISPAN NGN core network deployment for the FTTH technology.

Figure 5.9: Example architecture with xDSL access illustrates a possible realization of the TISPAN NGN functional architecture, with an xDSL-based access network.



**Figure 5.9: Example architecture with xDSL access**

This configuration assumes the following:

A Core-Border Gateway Function (C-BGF) is implemented in a Core Border Node sitting at the boundary between the access network and a core network, at the core network side.

A Resource Control and Enforcement Function (RCEF) is implemented in an IP Edge node sitting at the boundary between core networks, at the access network side. In this example, this node also implements the ARF functional entity.

A Interconnection Border Gateway Function (I-BGF) is implemented in a Border GateWay (BGW) sitting at the boundary with other IP networks.

A Trunking Media Gateway Function (T-MGF) is implemented in a Trunking GateWay (TGW) at the boundary between the core network and the PSTN/ISDN.

An Access Media Gateway Function (A-MGF) is implemented in an Access Node (AN) which also implements a DSLAM.

A Residential Media Gateway Function (R-MGF) is implemented in a Broadband Access Router located in the customer premises.

### 5.1.2.3.2 Backup line Requirements

5.1.2.3.2.1 General Requirements

| Number | Requirement | |
|--------|-------------|---|
| G.1 | The Broadband Access Router MUST support a second WAN interface for resilience which may be of the same type as the primary interface, or it may be of a different type. | |
| G.2 | The Broadband Access Router MUST be able to detect when the main WAN interface is down. | |
| G.3 | The Broadband Access Router must have an internal function that handles the Dual WAN behavior, i.e the switching between the primary WAN interface and the backup WAN interface. | |
| G.4 | The Broadband Access Router MUST switch over to the backup WAN interface after a configurable period (default 60 seconds) when the primary WAN interface is down. | |
| G.5 | When using the backup interface, the Broadband Access Router must periodically check for IP connectivity on the primary | |
| G.6 | The Broadband Access Router MUST automatically revert to the main WAN interface when it comes back up for a configurable period (default 60 seconds) | |
| G.7 | The Broadband Access Router MUST be able to send an alarm to the Broadband Service Provider (BSP) to indicate that the main interface has failed. | |
| G.8 | The Broadband Access Router MUST be able to provide an indication to the user that the main interface has failed e.g. by means of an LED | |
| G.9 | The Broadband Access Router MUST be able to power down the secondary WAN interface when it is not in use | |
| G.10 | The switch to the backup line MUST trigger the authentication procedure for the backup subscription line. | |

5.1.2.3.2.2 xDSL Requirements

| Number | Requirement | |
|--------|-------------|---|
| X.1 | xDSL (in case of Fiber primary access) backup WAN interface MUST be available on the Broadband Access Router. One of the Ethernet ports MUST be able to be reconfigured to act as the back-up WAN interface<br><br>The Backup WAN interface are built into the Broadband Access Router | |
| X.2 | The Broadband Access Router MUST contain one of the following WAN interfaces:<br>ADSL2+<br><br>ADSL2+ Annex M<br><br>G.SHDSL<br><br>VDSL2 | |
| X.3 | The Broadband Access Router MUST be able to route WAN-LAN traffic at a minimum rate of 100 Mbps (pps tbd). | |
| X.4 | The Broadband Access Router SHOULD support routing both downstream and upstream traffic at the full PHY rate of the access technology. | |

| Number | Requirement | |
|--------|-------------|---|
| **X.5** | The Broadband Access Router MUST be able to bridge LAN-LAN traffic at a minimum rate of 200 Mbps (pps tbd). | |
| **X.6** | The Broadband Access Router MUST support multiple, simultaneous L2 connections per WAN interface (VLANs, ATM VCs) | |
| **X.7** | The Broadband Access Router MUST be able to simultaneously support PPPOA and IPOA, or PPPOE and IPOE on different L2 interfaces. | |
| **X.8** | The Broadband Access Router Should be able support a different QoS profile for the backup interface | |

### 5.1.2.3.2.3  LTE Requirements

| Number | Requirement | |
|--------|-------------|---|
| **L.1** | The Broadband Access Router MUST be able to support a wireless back-up interface | |
| **L.2** | LTE (in case of Fibre primary access) backup WAN interface MUST be available on the Broadband Access Router. If not, one of the USB ports MUST be able to be reconfigured to act as a back-up WAN interface (e.g. to support a broadband wireless USB dongle) | |

### 5.1.2.3.2.4  HSPA Requirements

| Number | Requirement | |
|--------|-------------|---|
| **H.1** | The Broadband Access Router MUST be able to support a wireless back-up interface | |
| **H.2** | HSPA (in case of Fibre primary access) backup WAN interface MUST be available on the Broadband Access Router. If not, one of the USB ports MUST be able to be reconfigured to act as a back-up WAN interface (e.g. to support a broadband wireless USB dongle) | |

### 5.1.2.3.2.5  FTTH Requirements

| Number | Requirement | |
|--------|-------------|---|
| **F.1** | The ONT MUST support two WAN fibre interfaces | |
| **F.2** | In case two ONT are present, backup WAN interface MUST be available on the Broadband Access Router. Otherwise, one of the Ethernet ports MUST be able to be reconfigured to act as the back-up WAN interface | |

#### *5.1.2.3.3  Subscriber Authentication Procedure for Backup line*

##### 5.1.2.3.3.1  LTE Subscriber Authentication

EPS AKA is the authentication and key agreement procedure that shall be used over E-UTRAN.

A Rel-99 or later USIM application on a UICC shall be sufficient for accessing E-UTRAN, provided the USIM application does not make use of the separation bit of the AMF in a way described in TS 33.102 [71] Annex F. Access to E-UTRAN with a 2G SIM or a SIM application on a UICC shall not be granted.

In this section it is described the procedures that enable the CPE to perform the subscriber authentication over LTE lines. These procedures are: Initial Subscriber Attachment, Fast Re-Authentication and Full Re-Authentication.

##### **5.1.2.3.3.1.1  Initial Subscriber Attachment**

Figure 5.10 depicts the Initial Subscriber Attachment for a CPE when a UICC card is introduced.

**Figure 5.10: Initial Attachment Procedure**

A connection is established between the CPE and the LTE network using a procedure that is out of scope for the present document.

[1]- The MME invoke this mechanism once the connection is established. The MME sends a subscriber Identity Request asking for the International Mobile Subscriber Identity (IMSI) that is the permanent identity of the subscriber.

[2]- The CPE answer the MME sending the Subscriber Identity response message which includes the IMSI in clear text.

[3]- The MME request to the HSS the EPS authentication vectors (RAND, AUTN, XRES,$K_{ASME}$) to authenticate the subscriber. Each EPS authentication vector can be used for authenticate the subscriber. If the Network Type equals E-UTRAN, the "separation bit" in the AMF field of AUTN shall be set to 1 to indicate to the CPE that the authentication vector is only usable for AKA in an EPS context.

[4]- The HSS sends an authentication response back to the MME that contains the requested information. If multiple EPS authentication vectors had been requested, then they are ordered based on their sequence numbers. The MME shall be aware of the order of the EPS authentication vectors and shall use it in order.

[5]- The MME sends to the USIM via the CPE the random challenge RAND and an authentication token AUTN for network authentication from the selected authentication vector. It also includes a $KSI_{ASME}$ for the CPE which will be used to identify the $K_{ASME}$ (and further keys derived from the $K_{ASME}$) that results from the EPS AKA procedure. The MME also includes in this message a Global Unique Temporary Identifier (GUTI) to be used in future Re-Authentications.

At receipt of this message, the USIM shall verify the freshness of the authentication vectors by checking whether the AUTN can be accepted as described in TS.33.102 [71]. If so, the USIM computes a response RES. USIM shall compute CK and IK which are sent to the CPE. A CPE accessing E-UTRAN check during authentication that the "separation bit" in the AMF field of AUTN is set to 1. The "separation bit" is bit 0 of the AMF of field AUTN.

[6]- The CPE answers with subscriber authentication response including RES in case of successful AUTN verification and successful AMF verification. In this case the CPE computes $K_{ASME}$ from CK, IK and the serving network's identity (SN id).

[7]- The MME checks that the RES equals XRES (expected Response). If so the authentication is successful. If not the MME may initiate further identity request or send an authentication reject message towards the CPE.

#### 5.1.2.3.3.1.2 Subscriber Re-Authentication

As described in section 5.1.5.1.2.2 of deliverable D5.1, the Subscriber Re-Authentication is needed to detect the removal of the UICC card from the CPE. For the sake of not overloading the HSS with successive request of authentication data, and for preventing transmitting the IMSI on the radio path, the fast authentication procedure is performed employing a Global Unique Temporary ID (GUTI). The MME retrieve the IMSI of the user throughout the GUTI received.



**Figure 5.11: Fast Re-Authentication Procedure**

[1]- The MME sends a subscriber Identity Request to identify the user.

[2]- The CPE answer the MME sending the Subscriber Identity response message which includes the GUTI.

[3]- At receipt of this message, the MME obtains the IMSI throughout the GUTI. It also can check that it is the MME which sent this GUTI to the CPE. The MME retrieves the authentication data derived from the HSS in a previous Authentication (Full or Initial). The MME sends to the USIM via the CPE the random challenge RAND and an authentication token AUTN for network authentication from the selected authentication vector. It also includes a $KSI_{ASME}$ for the CPE which will be used to identify the $K_{ASME}$ (and further keys derived from the $K_{ASME}$) that results from the EPS AKA procedure.

At receipt of this message, the USIM shall verify the freshness of the authentication vectors by checking whether the AUTN can be accepted as described in TS.33.102 [71]. If so, the USIM computes a response RES. USIM shall compute CK and IK which are sent to the CPE. A CPE accessing E-UTRAN check during authentication that the "separation bit" in the AMF field of AUTN is set to 1. The "separation bit" is bit 0 of the AMF of field AUTN.

[4]- The CPE answers with subscriber authentication response including RES in case of successful AUTN verification and successful AMF verification. In this case the CPE computes $K_{ASME}$ from CK, IK and the serving network's identity (SN id).

[5]- The MME checks that the RES equals XRES (expected Response). If so the authentication is successful. If not the MME may initiate further identity request or send an authentication reject message towards the CPE.

Figure 5.12 describes the procedure for a Full re-Authentication. Full Re-Authentication procedure is performed when the GUTI is not valid anymore and for increase the security requesting fresh authentication vectors to the HSS.



**Figure 5.12: Full Re-Authentication Procedure**

[1]- The MME sends a subscriber Identity Request to identify the user.

[2]- The CPE answer the MME sending the Subscriber Identity response message which includes the GUTI.

[3]- The MME is either not able to map the GUTI to the IMSI or it is time to refresh the EPS authentication vectors. The MME sends a Subscriber Identity Request requesting the IMSI identity.
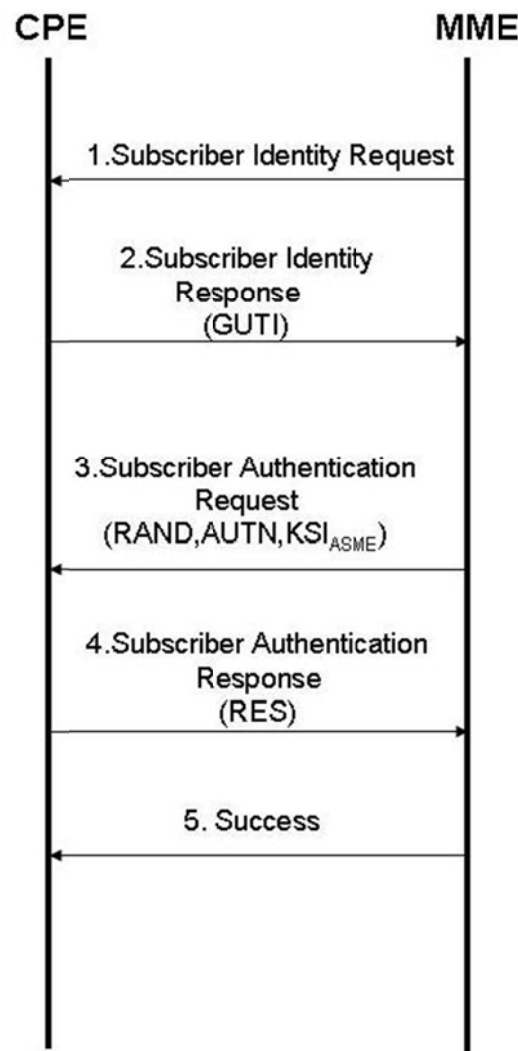
[4]- [9] The procedure is similar to the Initial Attach Procedure.

### 5.1.2.3.3.2 xDSL Subscriber Authentication

As described in section 5.1.2.3.1.4, one of the key concepts for choosing xDSL as Backup technology is the possibility of compliance with the NGN (NASS + RACS) subsystems. The EAP server located within the Network Attachment Subsystem (NASS) needs no adjustments to operate either with a fibre or a xDSL backhaul. The 802.1x Authenticator that filters 802.1x packets and maps the authentication messages into RADIUS packets would be integrated within the DSLAM instead of the OLT.

Summarizing the procedures for the Subscriber authentication are unaware of the underlying fixed technology even though the authenticator is integrated within another server.

### 5.1.2.3.3.3 HSPA Subscriber Authentication

The protocol used for authentication within the EPC is EPS AKA, which is based on UMTS AKA. Therefore if HSPA is chosen as backup technology, the procedures are very similar to the ones shown in section 5.1.2.3.3.1 for LTE technology. As occurred in section 5.1.2.3.3.2 for xDSL, the main changes are related to the actors involved in the authentication process. In UMTS, the VLR/SGSN will be the node in charge of handling the backend authentication and the repository where subscriber profile related information, and the authentication vectors are kept will be the HLR instead of the HSS node.

## 5.1.3 Conclusions

This section has described the subscription authentication solution for a scenario different from the standalone. In this new scenario several Femtocells give service to a larger amount of users than in a residential scenario. The subscription authentication in this scenario is the same than in the residential scenario because the subscription is also fibre based. Nonetheless, there are some features that an enterprise may require that are not necessary in the standalone scenario. This extra feature is for instance the necessity of a backup line to support critical services in case of failures in the main subscription. As a consequence some logic would need to be developed to detect failures in the main subscription and trigger the activation of the secondary subscription.

Within the networked Scenario section it is studied different possibilities for the backup technology which will support the authentication for the secondary subscription. The choice for this secondary subscription is based on throughput, security, and compliance with other subsystems already operating for the primary subscription.

Regarding having a secondary subscription based on fibre using a FTTH/GPON technology, may be a good choice because it could be supported all services already operational for the primary subscription.

Although it is not common to have an ONT with two fibre interfaces, the main constraint is that in current fibre deployments, neighbouring subscriptions are attached to the same PON port on the OLT. This leads that in case of failure in the primary subscription, depending where the failure is located, it will likely affect both subscriptions.

Regarding having a secondary subscription based on xDSL technology results in a very feasible choice due to the wide spread of this technology nowadays. Moreover, xDSL provides enough bandwidth and throughput to support most services, at least the critical ones, already deployed within the primary subscription. Furthermore there would be no problems for compliance with the TISPAN NGN functional architecture so coexistence among services for both technologies is ensured. The main drawback is that the choice for backup technology should last in time. With the imminent deployment of fibre, it may occur that the cooper access will likely yield to the fibre infrastructure even in the last mile.

Regarding having a secondary subscription based on HSPA is a good choice for a short term because as well as xDSL technology, HSPA is wide spread and already deployed by service providers. It supports reasonable data rates and fits with the idea of supporting critical services rather than more ambitious targets such as load balancing or voice session continuity. On the other hand, as well as the xDSL technology, new mobile technology releases such as LTE are in a mature state and extensive deployments will take place in a medium term.

Regarding having a secondary subscription based on LTE is the best choice in a medium term because it will largely support critical services in case of failures in the primary subscription and will fit with the NGN subsystems leading to a convergence for both fixed and mobile paradigms.

## 5.2 Access Control to Local Network and Services

The problem of access control in the BeFEMTO network architecture can be divided into two different aspects: the access control to services provided through the mobile network operator, and the access control to services of the network of femtocells.

The LFGW must have the ability to control which user equipments (UE) can have access to the local network, and the femtocell network operator must have the ability to decide what resources wants to share and which UE can make use of these resources. This access restriction is needed because the network of femtocell owner may not be happy to share neither the backhaul link nor the own services with all the users.

The first part of this section is devoted to access control to the services provided through the MNO. In this case, the situation does not introduce any changes especially relevant with respect to scenario of a single femtocell. From the MNO point of view, the LFGW plays the same role as the H(e)NB in the case of single femtocell. Recall that in the BeFEMTO architecture, the LFGW is seen as a HeNB by the MNO and as the EPC by the HeNBs inside the network of femtocells. The LFGW must establish secure tunnels with the different H(e)NB and aggregates all these tunnels into a new one that provides secure communications between the LFGW and SeGW. The elements involved in this scenario are shown in the next figure:



**Figure 5.13: Entities involved when accessing services provided through the MNO**

The user equipment is authenticated by means of the HSS/AAA server by using EPS-AKA. Moreover, the tunnels between the H(e)NB and LFGW, and between LFGW and SeGW are IKEv2 based [97].

The document explains how to set these tunnels and how to authenticate the involved entities. The access to the resources provided by the mobile operator depends on the CSG subscription information that is stored in the HSS and retrieved by the MME. Therefore, the document also describes how to manage such CSG data subscription.

In terms of access control to services of the local network of femtocell, it seems clear that it is important to minimize the processing load and communication with entities of the MNO. The management of permissions for different UE to access the different services provided by the femtocell network should be only dependent on the network of femtocell operator. It is important to point out that the LFGW acts as a proxy MME and proxy S-GW, hence enabling local mobility management and data forwarding (see D2.2). The second part of this section analyses this problem.

### 5.2.1 Access control to services provided through the MNO

The access control to services provided through the MNO is based on the concept of closed subscriber group (CSG) membership [98]. The closed subscriber group (CSG) identifies a group of UEs allowed to access to one specific femtocell, group of femtocells or network of femtocells. One way to link the network of femtocells to the UEs allowed to use them is to define a CSG identity that is unique. This unique CSG identity is broadcast by all femtocells or networks of femtocells that support access to UE belonging to the related CSG. The UE also needs to store the CSG identity of the CSG of which it is a member in order to verify that it can access the network of femtocells. It is expected that a given UE may be subscribed to one or more CSGs; hence, each UE should store a list of the allowed CSG identities. The mobile network operator performs access control based on the CSG ID advertised by the CSG cell and the CSG subscription data of the UE stored in the network. The mobile network verifies whether the UE is a CSG member or not.

The CSG is addressed in 3GPP TS 25.367 [99], TS 24.285 [100] and 22.220 [101]. TS 25.367 [99] describes CSG Identification, CSG Selection, CSG Cell Reselection and CSG and Hybrid Cell Handover. TS 24.285 [100] presents allowed CSG list management object, which is used to manage the list of allowed CSG IDs, the Operator CSG IDs and the related restricted access information at UE. TS 22.220 [101] describes common requirements including access control requirements and display requirements for CSG. TS 22.011 [102] denotes Access Control including Access Classes description and Emergency calls.

Once a UE has camped on a network of femtocell, it does not mean that the UE is automatically allowed to make or receive a call via the H(e)NB and LFGW especially in a case where the owner expects the LFGW to be accessed only by an allowed set of users. Providing secure network access service requires access control based on the mutual authentication to prevent from illegal access and authorization of the clients and the access networks to ensure users privilege for the access, and supply security functions for data exchange.

The access control functions ensure a User Equipment (UE) has a valid subscription at a CSG cell where it performs an access or a handover and indicates whether the UE is a member or non-member with the associated CSG ID.

### 5.2.1.1 CSG functions and management

CSG provisioning includes adding or deleting a subscriber from the CSG list of subscribers as well as viewing the list of subscribers. There are four aspects to be considered to support CSG provisioning: a) How the subscriber is added or removed from the list of subscribers for a CSG; b) How the CSG subscription data for a subscriber is stored in the network; c) How the Allowed CSG list and the Operator CSG list are updated at the UE, and d) How to manage any delays between when the CSG subscription data is updated on the network and when it is updated at the UE.

The CSG provisioning network elements include:

- The CSG List Server hosts functions used by the subscriber to manage membership to different CSGs. For example, the CSG List Server includes the UE CSG provisioning functions which manage the Allowed CSG List and the Operator CSG list stored on the UE.
- The CSG Administration Server hosts functions used by the CSG manager to manage the CSG. For example, the CSG Administration Server includes the CSG administration function which manages the list of subscribers for a CSG, i.e., the access control list for the CSG.

The management of the CSG subscription data at the UE must fulfil the following requirements

- The UE shall contain a list of allowed CSG identities (Allowed CSG List). It shall be possible to store the Allowed CSG List in the USIM. When available, the list on the USIM shall be used. It shall be possible for both, the operator and the UE, to modify the Allowed CSG List.
- The UE shall allow the user to introduce new CSGs to the Allowed CSG List by means of manual CSG selection only.
- The UE shall maintain an operator controlled list of allowed CSG identities (Operator CSG list). It shall be possible to store the Operator CSG list in the USIM. When available, the list on the USIM shall be used. It shall be possible for the operator to modify the Operator CSG List.
- The two lists are maintained independently from each other. A change in the Operator CSG list shall not trigger the UE to modify the Allowed CSG list to reflect such change automatically.
- All CSG cells belonging to a CSG identity not included in the Allowed CSG List or Operator CSG list shall be considered not suitable by the UE

The CSG subscription data is permanently stored in the HSS as part of the subscriber data as specified in TS 23.060 [103] and TS 23.401 [104] and retrieved by the MME for access control during the attach procedure, service request procedure or tracking area updating procedure as part of the UE's subscription profile. The CSG subscription information is transferred to the MME using the MAP (TS 29.002 [105]) or Diameter (TS 29.272 [106]) protocols. The structure of the CSG related parameters is defined in TS 23.003 [107]; for each CSG-Id there is an optional associated expiration date which indicates the time when the subscription to the CSG-Id expires (an absent expiration date indicates unlimited subscription).

When a UE accesses a CSG cell, the MME shall check that the CSG ID of the CSG cell corresponds to a CSG ID in the CSG subscription data, and that the expiration time (if present) is still valid [108]. The following figure shows an example.

**Figure 5.14: Access control at a CSG cell or a hybrid cell**

1.- The UE initiates the NAS procedure by sending, to the H(e)NB, the appropriate NAS Request.

2. The H(e)NB forwards the NAS Request message together with the Cell Access Mode, CSG ID, and other information of the cell from where it received the message to the MME through the LFGW that contains the Proxy-MME. The CSG ID is provided by the H(e)NB if the UE sends the NAS Request message via a CSG cell or a hybrid cell. The Cell Access Modes is provided by the H(e)NB if the UE sends the NAS Request message via a hybrid cell. For example, the HeNB selects the MME from the GUTI and from the indicated Selected Network and forwards the TAU Request message to the MME along with the CSG ID of the HeNB. In the case where the HeNB or the LFGW is connected to a HeNB GW, the HeNB forwards the TAU Request message to the HeNB GW and the HeNB GW performs the function of selecting the MME from the GUTI and from the indicated Selected Network and forwards the TAU Request message to the MME along with the CSG ID of the HeNB.

3.- The LFGW decrypts the tunneled information, and the proxy MME store information corresponding to the UE device (identification)

4.- The MME verifies whether it holds subscription data for the UE. If there is no subscription data in the MME for this UE then the MME sends an Update Location Request message to the HSS.

4. The HSS acknowledges the Update Location Request message by sending an Update Location Ack (IMSI, Subscription Data) message to the MME.

The MME sends this information to the H(e)NB through the LFGW. The LFGW decrypts the tunneled information, and the proxy MME updates the subscription UE data that may contain the CSG subscription data for the MNO.

### 5.2.1.2 Communications between H(e)NB and Security Gateway

Current deployments of H(e)NBs always require IPSec tunnels between every H(e)NB and a Security Gateway, in order to prevent external attacks or even a possible misuse by the customers. All signalling, user, and management plane traffic over the interface between H(e)NB and SeGW shall be sent through an IPsec ESP tunnel that is established as a result of the authentication procedure.

Because encryption and decryption processes are performed at the ends of the communication link, it is not possible to divert some of the IP packets in any intermediate point to provide a LIPA or SIPTO service. The adopted solution is to ensure that the LFGW is a trusted environment like the HeNBs, such that it can terminate the IPsec tunnels from HeNBs and concentrate them into a single tunnel towards the Security Gateway. It is important to note that LFGW is a very critical point in the architecture, and must contain a trusted environment to store critical data. This Trusted Environment (TrE) shall be a logical entity which provides a trustworthy environment for the execution of sensitive functions and the storage of sensitive data. All data produced through execution of functions within the TrE shall be unknowable to unauthorized external entities. The TrE shall be built from an irremovable, HW-based root of trust by way of a secure boot process, which shall occur whenever a H(e)NB is turned on or goes through a hard reset.

In addition, the LFGW is responsible of the authentication of the Security Gateway and H(e)NBs. The LFGW authentication function resides in the Security Gateway. After successful mutual authentication between the LFGW and the SeGW, the SeGW connects the LFGW to the operator's core network. Any connection between the LFGW and the core network is tunnelled through the SeGW. The SeGW provides the LFGW with access to the HeMS and HeNB-GW.

The HSS stores the subscription data and authentication information of the H(e)NBs. When hosting party authentication is required, AAA server authenticates the hosting party based on the authentication information retrieved from HSS.

#### 5.2.1.2.1 SeGW, LFGW and H(e)NB. Authentication Procedure

As hundreds of thousands of femtocells are deployed, the scalability issue imposes costly re-configuration and operation in MME/S-GW. Because femtocells use enterprise broadband as the backhaul to connect to the mobile core network (CN), the security issue needs be considered in order to protect the integrity of the network from malicious operations. Therefore, the femtocell network needs to consider both problems.

The H(e)NB shall use IKEv2 protocol [109] to set up at least one IPsec tunnel to protect the traffic with the LFGW due to the insecurity of the corresponding link, i.e. a pair of unidirectional SAs between H(e)NB and LFGW. All signalling, user, and management plane traffic over the interface between H(e)NB and LFGW shall be sent through an IKEv2 tunnel that is established as a result of the authentication procedure.

Mutual authentication shall be performed using IKEv2 with public key signature based authentication with digital certificates. Therefore, the LFGW shall authenticate itself to the different H(e)NB connected to it with a certificate based on the globally unique and permanent LFGW identity, signed by an operator authorized entity. The H(e)NB sees the LFGW as a SeGW, so the name used in the subjectAltName field of the LFGW digital certificate must fulfil the rules specified in the 3GPP TS 33.320, i.e, if DNS is available, the LFGW's name is the FQDN used to resolve its IP address; otherwise it is the IP address of the LFGW

The same procedure must be used by the HeNB to be authenticated to the LFGW. The H(e)NB's identity and LFGW identity shall be stored in the TrE (Trusted Environment) and shall not be modifiable. The H(e)NB's private key and the LFGW's private key shall be stored in the corresponding TrE and shall not be exposed outside of these TrE. The root certificate used to verify the signatures shall be stored in the corresponding TrE and shall be writable by authorized access only. The verification process for signatures shall be performed by the H(e)NB's TrE, and by LFGW's TrE.

Moreover, the SeGW shall authenticate itself to the LFGW using a certificate signed by an operator trusted CA. In this authentication procedure, the SeGW must see the LFGW as a H(e)NB entity. As the , subjectAltName of the LFGW digital certificate is its own FQDN, the same digital certificate can be used in the two authentication procedures of the LFGW. This digital certificate and the intermediate CA certificates shall be obtained from the IKEv2 CERT payload. The trusted root CA shall be obtained from a SeGW local store of trusted CA certificates. The LFGW shall include its identity in the IDi payload of the first IKE_AUTH request, and the LFGW identity in the IDi payload may be used for policy checks. Initiating/responding end entities are required to send certificate requests in the IKE_INIT_SA exchange for the responder and in the IKE_AUTH exchange for the initiator. The messages for the IKE_AUTH exchanges shall include a certificate or certificate chain providing evidence that the key used to compute a digital signature belongs to the identity in the ID payload.

The LFGW may check the revocation status of the SeGW using OCSP. For security reasons, the use of SHA-1 is not recommended for newly created OCSP responses. The OCSP communication between LFGW and OCSP server may use the in-band signalling of certificate revocation status in IKEv2 according to RFC 4806 "Online Certificate Status Protocol (OCSP) Extensions to IKEv2" [110], through which the SeGW can include an OCSP response within IKEv2.

The SeGW may check the revocation status of the LFGW certificate using CRLs or OCSP as with the exception that the SHA-1 and SHA-256 hash functions shall be mandatory to support. For security reasons, the use of SHA-1 is not recommended for newly created CRLs and OCSP responses.

The SeGW shall implement support for either CRL checking or OCSP or both. The locations of the CRL Server and OCSP Responder may be in the operator's network or provided by the manufacturer/vendor. If the H(e)NB certificate contains CRL or OCSP server information, then the SeGW may contact this server for revocation information. If the CRL or OCSP server is located at manufacturer of H(e)NB, the distribution of revocation information is provided by the manufacturer directly. To use such revocation information, normally the SeGW needs a CRL or OCSP client capable to reach the public Internet to contact these servers.

Validity check of LFGW certificates in SeGW shall be configurable by the operator, i.e. whether to use CRLs, OCSP or both and whether to use operator CRL or OCSP server, manufacturer CRL or OCSP server, or more than one of them.

Autonomous validation is performed during secure start-up and performs validation of the LFGW. As IKEv2 allows the inclusion of information data into Notify Payload, information regarding the trustworthy state of the LFGW may be carried in the Notify Payload during IKEv2 procedures from the

LFGW to the SeGW. Notify Payload within IKEv2's IKE_AUTH message is protected by IKEv2 SK and AUTH. In addition, the Notify Payload, as constructed by the TrE, should include a nonce and should be cryptographically signed by the TrE.

If the SeGW and the H(e)NB-GW are not integrated, then the interface between them may be protected using NDS/IP. In the absence of a HeNB-GW, the HeNB is directly connected to the MME via the SeGW.

## 5.2.2 Access control to Local Network and Services

The list of subscribers that are members of a given CSG and the list of HeNBs that allow access only to members of that CSG are stored on mobile operator servers. As it has been previously mentioned, typically, there are web interfaces over which femtocell users can modify the CSG member lists of their femtocells. Nevertheless, this process is not very flexible, introduces some delay and can introduce privacy problems. Basically, the main problem is that this solution does not allow a fine-granular control over which subscribers can access which global and local networks and services. The proposed solution is based on the use of an access control list located in the LFGW managed by the network of femtocells operator, that indicates the different offered services and for each one of them, which UE is able to access. This LFACL (Local Femtocell Access Control List) is complementary to the CSG mechanism, and it has the advantage of being independent of the MNO. The functions of the LFACL are similar to the ones provisioned by the UAAF (User Access Authentication Function) (see the transport architecture of BeFEMTO in D2.2), but are limited to the local services.

When a UE connected to a H(e)NB wants to use some local network or services, a breakout is performed at a L-GW placed in the LFGW (see D5.1). In this case, LIPA services should be provided, but taking into account the functional architecture of the BeFEMTO LFGW node. Once LIPA or traffic is allowed, it is important the use of a LFACL (Local Femtocell Access Control List) to control this traffic. The LFACL contain the rules that are applied to the different services provided by the network, each with a list of UE permitted to use the service.

The network of femtocell manager shall be able, under the operator supervision, to add, remove and view UE membership in the LFACL (Local Femtocell Access Control List). For temporary members, it shall be possible to limit the period of time during which the subscriber is considered a member of a LFACL (granted access rights). It shall be possible to configure a time period for each temporary member. Moreover, unlimited membership to the LFACL must be allowed. When the time period expires, the network of femtocell shall no longer be considered to be available to provide its own services.

For adding a UE at a LFACL, the UE needs to be provisioned at the LFACL Administration Server using a permanent and unique identifier that is preferably easily accessible to the subscriber such as the MS international ISDN number (MSISDN). The structure of this identifier is provided at 3GPP TS 23.003.

The following figure shows the process to include new user equipment in the LFACL.



**Figure 5.15: Process to include a new user equipment in the LFACL**

The LFACL manager sends a request to the LFACL Administration Server to add or remove a subscriber to the LFACL. For example, the LFACL manager logs into a web page with a secure access (i.e., authenticated by using digital certificate or login and password); clicks a Tab on the web page for their LFACL; and selects a subscriber to add (including services an optional time limit for membership) or remove. Adding or removing a subscriber at a LFACL is subject to approval by the local network of femtocell operator. Once approved, the LFACL Administration Server communicates with the AAA Relay to update the subscriber's LFACL subscription data stored in the AAA Relay. It is important to point out that the AAA Relay component includes a subset of the information contained in the HSS, but there is not any critical information regarding the MNO point of view.

The expiration time may also be set if the LFACL manager has set a time limit for membership.

To offer access to local network services, a PDN connection can be required after the UE has been attached to the MNO. Moverover, the access to the local services will be offered after the inclusion of the UE identity in the LFACL List Server. The general procedure for a PDN connection is the following:



**Figure 5.16: general procedure for establishing a PDN connection**

The steps 1 to 9 are analogous to those presented in figure 2.26 of document D5.1 (in this figure it has been shown the case when the IP address is provided by the L-GW). When a local service is required, it is necessary to analyse if this service can be provided to this UE. This information is contained in the LFACL List Server, but the address contained in this list is the MSISDN, so the corresponding mapping between the IP address and the MSISDN identification is required. Therefore, the L-GW requests this mapping to the AAA Relay, and that identification (MSISDN) is introduced in the query sent to the LFACL Server. In the case that the services are permitted, the packet data will be exchanged between the UE and the corresponding local server.

## 5.3 Architecture and IP Security

Femtocell security requires careful consideration primarily due to the fact that equipment participating in operator's network communication is installed outside of its controlled environment. While a brief review of typical femto network architecture or a femto equipment vendor presentation may try to convince us about a strong security measures applied in a given case the reality may not be so bright. Recent (July 2011) announcement made by one of the security groups provides details of a hack that could lead to calls being intercepted and listened to. Some analysts are also suggesting that femto cell security will be broken and it is only a question of time. Their suggestion is to protect networks from hacked femtos.

Following section provides bases for security considerations and outlines possible initial security analysis of a femto network. The architecture and IP security section was written at a generic level to highlight the fact that security is a key component for any solution. Due to the fact that innovative solutions proposed by BeFemto will be tested primarily in a simulated environment security testes as such are only highlighted but will not be conducted. Steps described below have been used in a real life test environment. These steps would be used as an initial assessment of the infrastructure. Any issues found during this phase of security analysis would further be explored with use of custom tools and testing methods.

### 5.3.1 Architecture and security overview

Security of a system must be considered at the design time - related aspects must be taken into consideration within system architecture. Femto solutions are no different.

While BeFemto project does not have a strong security focus it is essential to have such topic included and considered along with innovative efforts in other aspects of femto technology. Work carried out in security is not expected to bring to life significant innovations; however, its goal is to ensure that BeFemto project provides up to date security guidelines aligned with proposed innovative solutions. This approach underlines importance of tackling all architecture components at the same time.

Various publications of femto related standards seems to provide sufficient amount of information and references for a proper architecture design and for implementation of secure environment. BeFemto may in certain areas change the approach, solution, or configuration of components. While these changes proposed by BeFemto consortium should not augment security aspect of system components a thorough approach must take under consideration the bases provided by published standards, changes proposed within BeFemto, and security of the final product.

### 5.3.2 Approach

In order to asses security of the femto model proposed by BeFemto consortium a number of aspects will be considered:

- Verification of BeFemto architecture against latest femto standards

- Analysis of modifications introduced by BeFemto

- Review of architecture from IP security perspective (feeding testing methodology)

- Definition of testing methodology and development of test procedure

- Analysis of test results and conclusions

Analysis, reviews, and verifications will be conducted 'on paper' and where applicable and justifiable, within a demo environment.

### 5.3.3 Reference model

A generic reference model will be used, as published in Femtocell Security Framework document [63]



**Figure 5.17: Security Reference Model.**

Other references have been introduced earlier in the document, chapters 3 and 5. All of these references will be considered during the analysis and verification process.

### 5.3.4  Security verification processes

The following verification processes will be further developed and customized to the architecture of BeFemto proposed deployment scenarios.

#### 5.3.4.1  Architecture review

**TEST CASE Description:**

The point is to verify security of the system. Scope: RBS, SGW, infrastructure available via IPSec tunnel.

**Pre-Requisites:**

Configuration and architecture of system components is available for review.

Documentation of system components is available.

**Procedure:**

Review of the architecture:

is the architecture compliant with PTC requirements (D-NLB-00-02, I-NLB-00-01),

is defense-in-depth principle applied?

is the architecture compliant with best practices?

#### 5.3.4.2  Configuration review

**TEST CASE Description:**

The point is to verify security of the system. Scope: RBS, SGW, infrastructure available via IPSec tunnel.

**Pre-Requisites:**

Configuration and architecture of system components is available for review.

Documentation of system components is available.

**Procedure:**

Review of configuration taking into account:

manufacturer documentation,

guidelines and requirements (D-NLB-00-02, I-NLB-00-01),

international standards,

legal requirements,

best practices.

#### 5.3.4.3  RBS network scan (applicable if demo is available)

**TEST CASE Description:**

The point is to verify the list of open ports.

**Pre-Requisites:**

RBS is up and available.

RBS is operational (connected to BSC).

**Procedure:**

The level of security means can be assessed by using network scanners and standard network tools. The scan is performed from the RBS local network. Only the following protocol/ports on RBS should be in the listening state:

IP proto 50 (ESP)

IP proto 51 (AH)

IP ICMP

Any other traffic must be dropped or icmp unreachable should be sent.

#### 5.3.4.4 SGW network scan (applicable if demo is available)

**TEST CASE Description:**

The point is to verify the list of open ports.

**Pre-Requisites:**

SGW is up and available.

System is operational (there are some RBSes connected to BSC).

**Procedure:**

The level of security means can be assessed by using network scanners and standard network tools. The scan is performed from the SGW local network. Only the following protocol/ports on SGW should be in the listening state:

IP proto 50 (ESP)

IP proto 51 (AH)

IP icmp

Any other traffic must be dropped or icmp unreachable should be sent.

#### 5.3.4.5 Internal network scan (via IPSec tunnel) (applicable if demo is available)

**TEST CASE Description:**

The point is to verify internal hosts accessible via IPSec tunnel. This case requires establishing IPSec tunnel from test workstation with RBS credentials.

**Pre-Requisites:**

SGW is up and available.

System is operational (there are some RBSes connected to BSC).

Valid RBS credentials are configured on the test workstation.

**Procedure:**

The level of security means can be assessed by using network scanners and standard network tools. The scan is performed from behind IPSec tunnel.

Only the following hosts/protocols/ports should be accessible:

Abis_Synch

AbisIP

AbisIP

OSS

OSS

NEDSS

O&M client

Abis_Synch IP icmp

AbisIP IP icmp

OSS IP icmp

NEDSS IP icmp

O&M client IP icmp

Any other traffic must be dropped or icmp unreachable should be sent.

#### 5.3.4.6 Traffic outside of IPSec tunnel (applicable if demo is available)

**TEST CASE Description:**

The point is to verify if traffic is generated outside of the tunnel. Scope: RBS, SGW, infrastructure available via IPSec tunnel.

**Pre-Requisites:**

System is operational.

Physical access to RBS network (network hub).

**Procedure:**

Run network sniffer, power up RBS, wait until operational, perform voice call, data connection. Verify that no traffic is generated outside of IPSec tunnel.

#### 5.3.4.7 RBS secure boot

**TEST CASE Description:**

The point is to verify that RBS boot process is secured.

**Pre-Requisites:**

System is operational.

Physical access to RBS network (network hub).

**Procedure:**

Power up RBS. During the startup process verify periodically (in a loop) that no well known ports are in the listening state (except from ports open during RBS operational state).

#### 5.3.4.8 System design guidelines development

**TEST CASE Description:**

The point is to verify security of the system. Scope: RBS, SGW, infrastructure available via IPSec tunnel.

**Pre-Requisites:**

The architecture and configuration is defined and well documented

Assumption of physical access to a working RBS device taken into account.

**Procedure/scope** includes but is not limited to:

Discussion of OS security (RBS, SGW, systems accessible via IPSec tunnel).

Enhancements (if any) to application security (RBS, SGW, systems accessible via IPSec tunnel).

General security.

### 5.3.5 Conclusions

Project BeFemto relies on general security standards applicable to femto technology. With this in mind a solution specific analysis of modifications proposed by the consortium must be performed in order to verify whether BeFemto introduces any changes significant to security, and if so, what measures should be applied in order to maintain high security level of the final architecture.

# 6. Revenue Sharing in Multi-Stakeholder Scenarios

## 6.1 Introduction

In 2010 the concept of Revenue sharing mechanisms related to complex H(e)NB scenarios has been defined both on a business and on technical level. General description of potential actors involved in this process and key technical functionalities have been prepared as well as main business and technical challenges have been identified. Also one of the most complex scenarios – Stuck at the Airport – has been in details described from the revenue sharing perspective. The results were included in D.1.

This document focuses primarily on further important aspects related revenue sharing for complex business and technical scenarios:

- Legal and privacy aspects analysis pertaining to mobile services

- Further description of Shopping Mall scenario and services.

- Discussion on FemtoCells National Roaming.

## 6.2 Revenue sharing FemtoCells service deployment obstacles

In general mobile services could be roughly classified into two different groups:

- The first group includes applications typically downloaded from UpStores which use mostly only local information collected by handsets and local data stored by device owners. During the application installation process the user granted the given application an access to certain resources and data and authorizes the application do carry out some actions on his behalf. In many cases, the applications interact with third parties, but these companies typically are not ruled by the telecommunication law. The mobile operator usually doesn't directly control or even doesn't know what their subscribers install and use. The number of such applications available on various upstores grows really rapidly.

- The second group includes applications and services which take advantage of mobile operator managed network data and features like: location service, SIM authorization and information about subscriber behaviour. Currently, these services are offered mostly only by mobile operators (in many cases an operator cooperates with a third party). The operator is obliged to comply with telecommunication law which gives more confidence and security to the subscribers, but also limits a way the operator data is used. The services are sold under the operator brand and then centrally maintenance by the operator, usually a big company thus the costs of the service is usually high. In addition a large company is typically not too well prepared (and has high overhead costs) to address needs of small user groups. In the results the niche markets are not too well covered.

To speed up the creation of new applications related to the second group, is seems, a new business approach is needed: instead of centralized application creating, marketing and maintenance processes, new smaller, more agile, not too much burden by administrative cost structures are necessary. Small branches of a large companies or small or medium enterprises would be closer to clients, would know better clients' needs and would be able to cover new service areas traditionally not enough profitable for big players.

This idea, though promising, is heavily impacted via law limitation deeper described in the next paragraph.

## 6.3 Legal and privacy aspects

In EU member countries there are three major groups of legal acts impact design and usage of mobile services (including FemtoCells specific)

- Personal data protection law

- The law related to the process of service provisioning via electronic media

- Telecommunication law

In addition to these the most important there are also a number of other regulations which may be applicable in some situations and for certain services (e.g. pertaining to payment procedures, permission ect).

The specific implementation of these acts may slightly vary from country to country, but generally they should based on a number of appropriate European Commission directives:

- Directive 2002/21/EC of the European Parliament and of the Council of 7 March 2002 on a common regulatory framework for electronic communications networks and services (Framework Directive)

- Directive 2002/20/EC of the European Parliament and of the Council of 7 March 2002 on the authorization of electronic communications networks and services (Authorization Directive)

- Directive 2002/19/EC of the European Parliament and of the Council of 7 March 2002 on access to, and interconnection of, electronic communications networks and associated facilities (Access Directive)

- Directive 2002/22/EC of the European Parliament and of the Council of 7 March 2002 on universal service and users' rights relating to electronic communications networks and services (Universal Service Directive)

- Directive 2002/58/EC of the European Parliament and of the Council of 12 July 2002 concerning the processing of personal data and the protection of privacy in the electronic communications sector (Directive on privacy and electronic communications)

- Commission Directive 2002/77/EC of 16 September 2002 on competition in the markets for electronic communications networks and services (Text with EEA relevance)

- Directive 1999/5/EC of the European Parliament and of the Council of 9 March 1999 on radio equipment and telecommunications terminal equipment and the mutual recognition of their conformity

- Council Directive 89/336/EEC of 3 May 1989 on the approximation of the laws of the Member States relating to electromagnetic compatibility.

For purposes of this document it is enough to concentrate only on certain provisions.

- The personal protection law describe a process and details how the personal data may be collected, stored and processed. The key point is that the personal data processing should be fully transparent and the person the data concerns should explicate express assent on his data storing and processing and he should be informed for which purposed his data may be used and by whom. The person may at any time revoke his permission. In such a case the data administrator has to delete the data pertained to this person. By definition personal data is all the data which could be used to identify a person the data refers to. According to this definition many information produced by telecommunication systems should be treated as personal data.

- The law related to service provisioning via electronic media is an essential act because in particular it regulates to whom a company may send its ads and other marketing information in an electronic way, so the act automatically concerns almost all Internet and mobile services which directly or indirectly generate revenue. The key provision of this act is that before communication with a client is started, first this person should explicitly express his/her interest on receiving marketing information. This declaration may be done in different forms (SMS, e-mail, paper docs, phone conversation, finally via checking in an appropriate option in a given application). It is also acceptable to ask the client first if he agrees and then to send him a proper offer message.

- The telecommunication law significantly limits the storage and usage of subscribers' related data (including traffics, location, financial and personal related). In general, this data can't be processed, disclosure nor passed to a third party by a mobile operator (unless the party is directly involved in the service offering process and carries out some tasks on the operator behalf). Also an operator is not allowed to use its subscriber database for advertisement purposes (unless explicitly agreed by a subscriber).The only justified subscriber data usage is improvement, security, clearing and documentation of the operator services which in practice includes only billing, fraud detection and network optimization. However, these limitations can be suspended provided that a subscriber explicate allows an operator to process his data (telco and other) /or / and to receive commercial offers. These declarations are typically collected during signing new / or updating existing contracts.

There are basically two classes describing clients (in this case operator subscribers)

OPTIN   who only agree on processing their telecommunication data

OPTIN +  who agree on processing their telco (personal) data  and receiving commercial offers (including the third party advertisements) send by on operator.

For a mobile operator case currently a quite significant percentage of subscribers are members of OPTIN and OPTIN+ groups and these statistic are steadily growing. It is expected that soon these groups may count about 50%. If so it would significantly enlarge an operator opportunities to offer mobile services to a large number of its clients .

## 6.4 Legal and privacy aspects and H(e)NB specific revenue sharing services

By definition revenue sharing process takes place when more than one party participates in a service delivery process, provides some value the delivery chain, receives some revenue or participates in losses. The party contribution may include certain information delivery (like marketing offer), access to a client group (via client database) or technical means to execute an application or to establish communication session with a customer.

There are a few observation resulted from the legal analysis:

- A trivial observation is that there must basically be always one party who offers a service and acts as a front-end to a client. The party may (often case) in background cooperate with other third parties. However, the less obvious remark is that the transfer of clients related information between parties involved in service delivery process is significantly limited. An example of such mechanisms is an operator cooperation with an insurance company. The operator may select a sub-group of his subscriber (from OPTIN) group according to the criteria  defined in advance by an insurance company and then offer a certain insurance product on behalf of its partner. However, the insurance company has never a direct access to the operator subscribers database and don't directly participate in subscriber segmentation process.

- An operator can't locate its subscribers without their explicate expressed permission. If a subscriber who is not  a member of OPTIN the operator is not allowed to carry out any actions when f.e. this subscribers enters the H(e)NB range. In this case sending any information to the subscribers is not legal (even a question if he/is would agree to receive ads), because this would require the earlier subscriber location.

- The personal and telecommunication data can be processed provided that the data is anonymized before this process

- If a subscriber gives his assent to process his personal data (including telecommunication data) it should be clearly stated which party his assent concerns. Referring to a complex scenario of a Shopping Mall, the client actually should give his assent to all shops located in the Shopping Mall. If location information is provided by another party (e.g. a mobile /network infrastructure operator) then the situation is getting more complicated because an addition client agreement is necessary to pass the location information to the shops. The process of obtaining many assents may be pretty annoying and discouraging for clients and in the result may decrease usability and potential of the mobile services for Shopping Mall scenario.

 Conclusions:

- The current legal regulations protect well the clients/ subscribers' privacy, but they also significantly limit and complicate the way the customer data may be utilized.

- It seems the legal logic and provisions related to privacy protection do not efficiently address many scenarios which have appeared in the last years as the result of technological advancement (sophisticated mobile services, advanced service and product personalization based on client related data, deep data analysis ect.). A few new amendments would  be necessary to speed up deployment of services based on client data. The point is not to decrease the client privacy and security, but to simplify some processes and to introduce a more coherent legal interpretation in several situations.

- Despite some limitations , the currently legal situations allow practical implementation of advanced services with revenue sharing between parties. (A separate subchapter describes proposal how a vision related to Shopping Mall scenario may be implemented as well as list some classes of services)

### 6.4.1  Other factors to be taken into account

Even in an anonymised form (usage of which is legally permissible) telecommunication and other people' activity related data contains a lot of information on individual and group behaviour. Recently many companies try to use this data for commercial purposes. Apart from strictly legal aspects (in many

particular cases it is not clear what is legally acceptable and what not) the customers and public opinion in many countries complain on illegitimate usage of their data and on increasing invigilation and interference in their life by commercial companies. These opinions may result in adoption of more restricted legal regulations or at least may destroy reputation of some companies.

Simultaneously many people agree on processing their behaviour data expecting that this information may help the companies and service providers to prepare tariffs and commercial offers which will more accurately address their specific needs and interests.

Regardless the result of these discussions it seems that in the future there will be significant social acceptance for usage of activity data unless this process will keep some ethical standards and won't break the legal rules.

## 6.5 Shopping Mall scenario and specific services.

Thanks to latest technical achievements (especially availability of powerful handsets and other sophisticated mobile devices with network connectivity) the processes of virtual and physical shopping and entertainment could be closer integrated. It could create new opportunities both for service providers and their clients.

Based on currently available technology and legal regulations this vision may be realized for example in the following way.

- A person enters a Shopping Mall with an intention to spend some time doing shopping, relaxing and having fun and entertainment.

- The person has his smartphone, tablet or could borrow such devices at the Shopping Mall (if doen't have his own)

- The Shopping Mall has implemented an integrated shop assistance and entertainment services platform. The platform includes information system, advanced advertisement, location services, user authorization (based on SIM), interactive games, mall navigation system and other dedicated applications. The system is connected with clients profiles collected by shops, restaurants and pubs and client loyalty programs databases. Depending on a current situation shops / restaurants may dynamically compose general, segmented, or personal offers.

- As the person enters the shopping mall, he is encouraged to install on his smartphone a special application (the details /web addresses could be provided via QR code placed on mall billboards) or to open a dedicated WebPages. The person could login to the system using his SIM credentials (the mobile operator plays a role of Identify Provider) or to work with the application in a guest (an anonymous) mode. Then the person is asked to set up several options – in particular to decide if he agrees to be localised by the network and to pass his coordinates to specific companies connected also to this system. Because the decision to start an application comes from the client - the legal requirements are satisfied.

- The person visits various areas of shopping mall. Even in the anonymous mode he can still have access to many services: his smartphone can determine his location thanks to the mall location services (without being located by mobile operator) and the person may check the newest adds. and current general offers published by shops.

- Visiting a given shop (if he wants) a person may disclosure his identity and take advantage of special offers (e.g customer loyalty program, discount for credit card owners, subscriber of a mobile operator the shop has an agreement with)

- Being not anonymously login (if he agrees) the person could actively receive profiled advs. and information. Two legally consistent options are possible:

  o The systems sends selected authorised by the person and his Identify Provider information to a shop. The shop composes a special offer taking also into account own information about this person (if it already has his profile) and send it to the client. The system role is limited to the transmission media.

  o The systems (effectively the mobile / or infrastructure operator) using own info about an user composes an offer based on criteria provided by the shop.

- Due to Femtocells technology Shopping Mall may deliver information and ads. in a rich multimedia form (interactive, special audio/ video clips, dedicated applications)

- Thank to shopping mall location systems the clients may enjoy reality games with other shoppers. These games my include specific places discovery, answering some questions (f.e. related to shops offers) and interaction with other players.

- Among the other services the Shopping Mall platform may offer it is worth to mention:
  - o Friends discovery (integrated with the popular social services). In this case legal requirement can also be easily met.
  - o A restricted access to Voice / SMS /data services in some public places (cinema, exhibitions).
  - o A customer location sensitive and interactive billboard which would display customized information tailored to a customer profile who approaches the billboard. Additional functionalities may allow the customer to navigate the billboard via a dedicated mobile application (plus the standard navigation method provided by "traditional" interactive billboard).
  - o Place specific services and applications (for cinemas, application offered by shops
  - o Push personalized public information services (eg. Info about flight delays)
  - o Downloading multimedia content for free (agreeing to receive ads)
  - o Location of selected subscribers (e.g. where is my child case)
  - o Access to price comparison system and opinions about given products (shops)

### 6.5.1 PicoCell trial

While PTC has no yet experience in Femtocells based commercial advertisement, a few observations and forecasts may be done analysing the results of a PicoCell based location advertisement trial in one of the Warsaw shopping malls.

#### 6.5.1.1 Trial description

This trial was executed in the end of 2010. PTC subscribers who belonged to the OPTIN+ group (see. 1.2.2) and who were localized at this shopping mall received SMSs / or MMSs adv. encouraging him to visit a selected shop / or shops located in this mall. The shop or / and a mobile operator offered some discount provided that a client carried out a certain transaction. 2-3 days after receiving an offer the results were evaluated via phone inquiry. The trial trail confirmed high potential of mobile advertisement (10% of those who receive SMS/MMS visited a given shop) and also showed that the subscribers in general positively reacted on receiving SMSs/ MMSs advs. provided that their number is not greater than 2-3 adv. a week.

#### 6.5.1.2 Conclusions from the PicoCells trial.

The PicoCell based location system used in the trial described above was far less accurate than location offered by Femtocells systems.

The system based on Cell ID informed only that a subscriber is inside / or near the shopping mall

The system used a passive location method (analysing network signalization data) triggered by a user / network actions (periodic location update, voice / data sessions ). In the result only a part of subscribers who did visit this shopping mall finally received adv. SMSs / MMSs.

The PicoCell trial showed that SMS/MMS services are good communication channels to distribute targeted advertisements. However, these services have a significant limitation – people accept and memorize at the most only up 10-15 SMS/MMS communicates monthly, these messages are quite short, not interactive and it is difficult to send in this way (if the customer is interested) additional information.

## 6.6 FemtoCells National Roaming Analysis

Although the Femtocells idea was introduced a few year ago, the Femtocells deployment process don't progress rapidly so far. Many various factors contribute to this situations, but one of them is doubtfully related to national roaming issues, which significantly slows down especially the stand alone Femtocells

deployment. It is pretty easy to understand why subscribers are not too eager to install H(e)NB at their homes: if household members use more than one mobile operator services (a typical situation) it is difficult for a subscriber to imagine installation of separated H(e)NBs for each operator. A quick solution for this issue seems to be crucial for Femtocells commercial success. As Femtocells roaming we will understand in this document – national roaming between mobile operators. International Femtocells roaming is also a very interesting topic with a lot of new opportunities and a potential market impact, but due to its complex nature we won't discuss it in this document.

### 6.6.1  Aspects impacting FemtoCells roaming

Despite the fact that implementation of Femtocells roaming procedures seems to be necessary, it has not been done / or even proposed in a satisfactory way yet. There are definitely a lot of reasons why Femtocells roaming is very difficult to address:

- Not obvious interest of the mobile operators – from an operator perspective one of the main motivation for Femtocells implementation is churn reduction. As it directly impacts total customer acquisition cost, churn reduction is typically a key factor responsible for positive Femtocells Business Case
- Not sufficient roaming support on the network architecture standardization level: Femtocells idea was introduced many years after the basis network standards were proposed and implemented. A standard mobile network roaming procedures and legal regulation can't be directly reused for Femtocells roaming. Also the standard roaming was originally focused more on international roaming than on national roaming
- Currently by default some network services may not be available for roamers– typically it is not a big problem for international roaming, but f.e. a lack of location service may drastically limit the Femtocells functionality
- Legal and confidentially problems (subscriber related activity data) – telecommunication law doesn't explicate address the national roaming case. So far only one operator has access to its subscribers data and is sole responsible for any data infringement.
- Not too popular infrastructure sharing practices in many countries yet – Femtocells sharing is a special case of infrastructure sharing.
- Not clear signals from national regulators how to approach this issue – national regulator main objective is to stimulate telco services development and to protect and to well address the client needs. Often used regulators strategy is to foster operator competition, while sharing operator infrastructure is not necessarily a step in this direction (as is implies closer cooperation between operators)
- Security issues – H(e)NB installed at the subscriber premise is more vulnerable to potential attacks
- Billing issues – it is also related to the legal issues - mechanisms for verification and access to selected information of subscribers activities via more than one operator is necessary
- Image and responsibility for Quality of Service – in a Femtocells roaming case other network subscribers would use only a limited part of their home operator infrastructure. The home operator wouldn't have direct influence not even control of QoS of this sessions. But from the subscriber perspective it still be responsible for service availability and quality.

### 6.6.2  Key drivers speeding up the Femto roaming process

- Strong pressure to lower network transmission costs / higher network capacity and new mobile services
- Femto is considered as one of the key approaches for LTE deployment - the national regulators are currently discussing various options to support Femtocells: dedicated Femtocells operators with separated bandwidth, Femtocells / LTE national roaming is also considered as a conditions to frequency tenders
- LTE as new technology allows easier implementation of roaming procedures
- some changes of telecommunication law could possibly be done to better address the Femto / Small Cell challenges
- Evolution of Femtocells  in the direction of closer unification with other small cell technologies

### 6.6.3 Remarks on Femtocells National Roaming options

A few potential solutions are possible. Each of them should marry many parties interests (operators, regulators, subscribers, service providers), satisfy legal and privacy requirements as well as be technically feasible. Unless LTE forces deeper standardization, there won't probably be one standards Femtocells roaming regulations, however every options will have to some extent address the following topics:

- competitiveness preservation - the H(e)NB is deployed and controlled only by one mobile operator / or dedicated Femtocells operator
- Symmetrical conditions - some (or all in the country) mobile operators sign agreements regulating access of other operator subscribers to its H(e)NB and vice versa (such agreement will probably include QoS monitoring, billing and payment rules). One of the possible provisions may state that the agreement does not change the other operator subscribers tariff (or give other subscribers a small discount to encourage them to user Femtocells
- Regulatory requirements – national regulators is a national stimulator of Femtocells roaming implementations. Regulator means may include: frequency managements, interconnect rates, network quality monitoring
- Rules regulating subscribers access to N(e)NB and Macro - an N(e)NB owner may have some freedom to decide if and with whom to share her/his N(e)NB (e.g. member of his family, neighbours, everybody ect.). The specific conditions to be proposed by an operator – a subscriber level will create a new field for competitiveness. One the other hand the subscriber should have some freedom regarding the decision if he/she want to use Marco or Femtocells

### 6.6.4 Summary

The Femto / Small Cells roaming is a crucial step toward cheaper and more capacious mobile network. The latest tendencies observed in telecommunication world (popularization of RAN and IP network infrastructure sharing, explosion of new mobile services taking advantage of various network features, requirement on more data transmission, approaching LTE massive deployment) challenge all telco stakeholder to come up with new technical and organizational solutions. In a few year perspective it will probably results also in implementation of Femtocells roaming. As the concept is very complex – it needs still deeper legal business and technical analysis.

# 7. References

[1] H. Zhang, X. Wen, B. Wang, W. Zheng, and Y. Sun, "A Novel Handover Mechanism Between Femtocell and Macrocell for LTE Based Networks," in Proc. of IEEE ICCSN '10, 2010, pp. 228–231.

[2] A. Ulvan, R. Bestak, and M. Ulvan, "Handover Scenario and Procedure in LTE-based Femtocell Networks," in Proc. of UBICOMM '10, Florence, Italy, 2010.

[3] L. Wang, Y. Zhang, and Z. Wei, "Mobility Management Schemes at Radio Network Layer for LTE Femtocells," in Proc. of IEEE VTC Spring '09, Barcelona, Spain, Apr. 2009, pp. 1–5.

[4] F. A. Zdarsky, A. Maeder, S. Al-Sabea, and S. Schmid, "Localization of Data and Control Plane Traffic in Enterprise Femtocell Networks," in Proc. of IEEE VTC-Spring 2011 Workshop on Broadband Femtocell Technologies, 2011.

[5] H. Soliman, C. Castelluccia, K. ElMalki, and L. Bellier, Hierarchical Mobile IPv6 (HMIPv6) Mobility Management, IETF RFC 5380, Oct. 2008.

[6] Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access Network (E-UTRAN); S1 Application Protocol (S1AP), 3GPP Std. TS 36.413, Rev. 9.4.0, Sep. 2010.

[7] Technical Specification Group Services and System Aspects; General Packet Radio Service (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access, 3GPP Std. TS 23.401, Rev. 10.2.1, Jan. 2011.

[8] J. Roberts, "A survey on statistical bandwidth sharing," Computer Networks, vol. 45, no. 3, pp. 319–332, Jun. 2004.

[9] T. Elteto, P. Vaderna, and S. Moln´ar, "Performance analysis of tcp networks loaded by web traffic," in Proc. of the 18th International Teletraffic Congress (ITC), Berlin, Germany, Aug. 2003.

[10] P. Lassila, H. van den Berg, M. Mandjes, and R. Kooij, "An integrated packet/flow model for tcp performance analysis," in Proc. of the 18th International Teletraffic Congress (ITC), Berlin, Germany, Aug. 2003.

[11] J. W. Roberts, "Internet Traffic, QoS and Pricing," Proceedings of the IEEE, vol. 92, no. 9, pp. 1389–1399, Aug. 2004.

[12] J. P. Curtis, J. G. Cleary, A. J. McGregor, and M. W. Pearson, "Measurement of Voice Over IP Traffic," in Proc. of Passive & Active Maesurement Workshop, Hamilton, New Zealand, Apr. 2000.

[13] A. Veres, Z. Kenesi, S. Molnar, and G. Vattay, "On the Propagation of Long-Range Dependence in the Internet," in Proc. of ACM SIGCOMM, Stockholm, Sweden, Sep. 2000, pp. 243–254.

[14] F. Delcoigne, A. Prouti`ere, and G. R´egni´e, "Modeling integration of streaming and data traffic," Performance Evaluation, vol. 55, no. 3–4, pp. 185–209, Feb. 2004.

[15] M. Mathis, J. Semke, and J. Mahdavi, "The macroscopic behavior of the tcp congestion avoidance algorithm," ACM Computer Communication Review, vol. 27, no. 3, Jul. 1997.

[16] T. V. Lakshman and U. Madhow, "The performance of tcp/ip for networks with high bandwidth-delay products and random loss," IEEE/ACM Transactions on Networking, vol. 5, no. 3, Jun. 1997.

[17] R. W. Wolff, "Poisson arrivals see time averages," Operations Research, vol. 30, no. 2, pp. 223–231, Mar. 1982.

[18] H. Fathi, R. Prasad, and S. Chakraborty, "Mobility Management for VoIP in 3G Systems: Evaluation of Low-Latency Handoff Schemes," IEEE Wireless Communications Magazine, vol. 12, no. 2, pp. 96–104, Apr. 2005.

[19] ITU-R, "FINAL EVALUATION REPORT FROM WINNER+ ON THE IMT ADVANCED PROPOSAL IN DOCUMENTS IMT-ADV/6, IMT-ADV/8 AND IMT-ADV/9," Tech. Rep., Jun. 2010.

[20] J.-H. Yun, M. Lee, and S. Choi, "Comparison of handover schemes for 3GPP Long Term Evolution and 3GPP2 Ultra Mobile Broadband," in Proc. of IEEE PiMRC '08, Cannes, France, Sep. 2008, pp. 1–5.

[21] 3GPP, "TR 25.912 V9.0.0, Feasibility study for evolved Universal Terrestrial Radio Access (UTRA) and Universal Terrestrial Radio Access Network (UTRAN)," Tech. Rep., Sep. 2009.

[22] Lundberg, T.; de Bruin, P.; Bruhn, S.; Hakansson, S.; Craig, S.; , "Adaptive thresholds for AMR codec mode selection," Vehicular Technology Conference, 2005. VTC 2005-Spring. 2005 IEEE 61st , vol.4, no., pp. 2325- 2329 Vol. 4, 30 May-1 June 2005

[23] Juejia Zhou; Xiaoming She; Lan Chen; , "Source and Channel Coding Adaptation for Optimizing VoIP Quality of Experience in Cellular Systems," Wireless Communications and Networking Conference (WCNC), 2010 IEEE , vol., no., pp.1-6, 18-21 April 2010

[24] Mkwawa, I.-H.; Jammeh, E.; Lingfen Sun; Khan, A.; Ifeachor, E.; , "Open IMS Core with VoIP Quality Adaptation," Autonomic and Autonomous Systems, 2009. ICAS '09. Fifth International Conference on , vol., no., pp.295-300, 20-25 April 2009

[25] Johnny Matta, Christine Papin, Khosrow Lashkari, and Ravi Jain. 2003. A source and channel rate adaptation algorithm for AMR in VoIP using the Emodel. In Proceedings of the 13th international workshop on Network and operating systems support for digital audio and video (NOSSDAV '03). ACM, New York, NY, USA, 92-99.

[26] A. Krendzel, J. Mangues, M. Requena, J. Núñez, "VIMLOC: Virtual Home Region Multi-Hash Location Service in Wireless Mesh Networks", in Proceedings of the IFIP Wireless Days Conference. 24-27 November 2008, Dubai (United Arab Emirates)

[27] 3GPP. TS 23.401 – v8.6.0, General Packet Radio Services (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access, 3G Release 8, 2009. http://www.3gpp.org/ftp/specs/archive/23 series/23.401/

[28] M.J. Neely, "Stochastic Network Optimization with Application to Communication and Queueing Systems", Morgan&Claypool Publishers, 2010.

[29] B. Jennings, et al., "Towards autonomic management of communications networks," Communications Magazine, IEEE, vol. 45, no. 10, pp. 112-121, 2007.

[30] S. Davy, et al., "Policy-based architecture to enable autonomic communications - a position paper," in Consumer Communications and Networking Conference, 2006. CCNC 2006. 3rd IEEE, vol. 1, 2006, pp. 590-594.

[31] K. Mase and Y. Toyama, "End-to-end measurement based admission control for VoIP networks," in Communications, 2002. ICC 2002. IEEE International Conference on, vol. 2, 2002, pp. 1194-1198vol2.

[32] J. Fitzpatrick, S. Murphy, and J. Murphy, "SCTP based Handover Mechanism for VoIP over IEEE 802.11b Wireless LAN with Heterogeneous Transmission Rates," in Communications, 2006. ICC '06. IEEE International Conference on, vol. 5, 2006, pp. 2054-2059.

[33] M. A. Hoque and F. Afroz, "Call admission control: QoS issue for VoIP," in 3rd International Conference on Communication Systems Software and Middleware and Workshops, 2008. COMSWARE 2008., 2008, pp. 757-761.

[34] V. Joseph and B. Chapman, Deploying QoS for Cisco IP and Next Generation Networks: The Definitive Guide. Morgan Kaufmann Publishers Inc., 2009.

[35] ITU-T. G.107, "The E-Model, A Computational Model For Use In Transmission Planning," 2005.

[36] Femto Forum, "Femtocell Synchronization and Location," June, 2010.

[37] D. Mills, J. Martin, J. Burbank, and W. Kasch, "RFC 5905: Network Time Protocol Version 4: Protocol and Algorithms Specification," Internet Engineering Task Force, 2010.

[38] J. C. Eidson, "IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems," IEEE Std 1588-2002, pp. i-144, 2002.

[39] 3GPP and T. B. Forum, "PCRF-BPCF Functional Split and Information Exchange," France Telecom / Orange, NEC, Alcatel-Lucent, 2010.

[40]   H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson, "RFC3550: RTP: A Transport Protocol for Real-Time Applications," *RFC Editor United States*, 2003.

[41]   S. Blake, et al., "An architecture for differentiated services," Citeseer, 1998.

[42]   E. Weingartner, H. v. Lehn, and K. Wehrle, "A Performance Comparison of Recent Network Simulators," in *IEEE International Conference on Communications, 2009. ICC 2009.*, 2009, pp. 1-5.

[43]   O. Tipmongkolsilp, S. Zaghloul, and A. Jukan, "The evolution of cellular backhaul technologies: Current issues and future trends,"Communications Surveys Tutorials, IEEE, vol. 13, no. 1, pp. 97 –113, 2011.

[44]   P. Bhagwat, B. Raman, and D. Sanghi, "Turning 802.11 inside-out," SIGCOMM Comput. Commun. Rev., vol. 34, pp. 33–38, January 2004. [Online]. Available: http://doi.acm.org/10.1145/972374.972381

[45]   B. Raman and K. Chebrolu, "Experiences in using wifi for rural internet in india," IEEE Communications Magazine, vol. 45, no. 1, pp.104–110, 2007.

[46]   S. Sen, S. Kole, and B. Raman, "Rural telephony: A socio-economic case study," in Information and Communication Technologies and Development, 2006. ICTD '06. International Conference on, May 2006, pp. 301 –309.

[47]   United Nations International Telecommunications Union, "Measuring the Information Society", 2010

[48]   E. S. Roman, "Bringing Broadband Access to Rural Areas: A step by step approach for regulators, policy makers and unviersal program

[49]   administrators. The Experience in the Dominican Republic," 9th Global Symposium for Regulators (GSR), 2009.

[50]   R. Flickenger, S. Okay, E. Pietrosemoli, M. Zennaro, and C. Fonda, "Very long distance wi-fi networks," in Proceedings of the second ACM SIGCOMM workshop on Networked systems for developing regions, ser. NSDR '08. New York, NY, USA: ACM, 2008, pp. 1–6. [Online]. Available: http://doi.acm.org/10.1145/1397705.1397707

[51]   "The ns-3 network simulator." [Online]. Available: http://www.nsnam.org/

[52]   A. Uvliden, S. Bruhn, and R. Hagen, "Adaptive multi-rate. a speech service adapted to cellular radio network quality," in Signals, Systems Computers, 1998. Conference Record of the Thirty-Second Asilomar Conference on, vol. 1, Nov. 1998, pp. 343 –347 vol.1.

[53]   R. F. H. Schulzrinne, S. Casner and V. Jacobson, "RTP: A transport protocol for real-time applications," July 2003. [Online]. Available: http://tools.ietf.org/html/rfc3550

[54]   "Wireshark - packet capture utility." [Online]. Available: http://http://www.wireshark.org/

[55]   A. L. J. Sjoberg, M. Westerlund and Q. Xie, "Real-time transport protocol (RTP) payload format and file storage format for the adaptive multi-rate (AMR) and adaptive multi-rate wideband (AMR-WB) audio codecs," June 2002. [Online]. Available: http://www.ietf.org/rfc/rfc3267.txt

[56]   ITU-T, The E-model, a computational model for use in transmission planning, International Telecommunication Union Recommendation G.107, May 2000.

[57]   J. Fitzpatrick, S. Murphy, M. Atiquzzaman, and J. Murphy, "Using cross-layer metrics to improve the performance of end-to-end handover mechanisms," Computer Communications, vol. 32, no. 15, pp. 1600 – 1612, 2009. [Online]. Available: http://www.sciencedirect.com/science/article/B6TYP-4WHFD6N-2/2/d8bfd3dcb414da3b25fec6e818479292

[58]   L. Carvalho, E. Mota, R. Aguiar, A. Lima, and J. de Souza, "An e-model implementation for speech quality evaluation in voip systems," in Computers and Communications, 2005. ISCC 2005. Proceedings. 10th IEEE Symposium on, 2005, pp. 933 – 938.

[59]   ITU-T, Transmission Impairments, International Telecommunication Union Recommendation G.113, Feb 1996.

[60]   L. Sun and E. Ifeachor, "Voice quality prediction models and their application in voip networks," Multimedia, IEEE Transactions on, vol. 8, no. 4, pp. 809 –820, 2006.

[61]    J. Matta, C. P´epin, K. Lashkari, and R. Jain, "A source and channel rate adaptation algorithm for amr in voip using the emodel," in Proceedings of the 13th international workshop on Network and operating systems support for digital audio and video, ser. NOSSDAV'03. New York, NY, USA: ACM, 2003, pp. 92–99. [Online]. Available: http://doi.acm.org/10.1145/776322.776338

[62]    F. Mertz and P. Vary, "Efficient voice communication inwireless packet networks," in Proceedings of Sprachkommunikation 2008 - 8. ITG-Fachtagung, 2008, pp. 92–99. [Online]. Available: http://www.vde-verlag.de/proceedings-de/453120011.html

[63]    ITU-T, Methodology for the derivation of equipment impairment factors from instrumental models, International Telecommunication Union Recommendation P.834, July 2002.

[64]    UTRAN Iuh Interface RANAP User Adaption (RUA) signalling, 3GPP Recommendation TS 25.468 R9, Dec 2009.

[65]    K. Medepalli, P. Gopalakrishnan, D. Famolari, and T. Kodama, "Voice capacity of ieee 802.11b, 802.11a and 802.11g wireless lans," in Global Telecommunications Conference, 2004. GLOBECOM '04. IEEE, vol. 3, Nov. 2004, pp. 1549 – 1553 Vol.3.

[66]    N. Baldo, M. Requena-Esteso, J. N´u˜nez Mart´ınez, M. Portol`es-Comeras, J. Nin-Guerrero, P. Dini, and J. Mangues-Bafalluy, "Validation of the ieee 802.11 mac model in the ns3 simulator using the extreme testbed," in Proceedings of the 3rd International ICST Conference on Simulation Tools and Techniques, ser. SIMUTools '10. ICST, Brussels, Belgium, Belgium: ICST(Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2010, pp. 64:1–64:9. [Online]. Available: http://dx.doi.org/10.4108/ICST.SIMUTOOLS2010.8705

[67]    Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (EUTRAN); Overall description; Stage 2 (Release 10), 3GPP Std. TS 36.300 v10.2.0, Dec. 2010.

[68]    L. Georgiadis, M. Neely, and L. Tassiulas. Resource allocation and cross layer control in wireless networks. Now Publishers. 2006.

[69]    IETF RFC 2401, Security Architecture for the Internet Protocol, November 1998.

[70]    IETF-PPPEX-L2TP-15.TXT, Layer Two Tunneling Protocol "L2TP", March 1999.

[71]    3GPP TS 33.102: "3G security; Security architecture".

[72]    A. Jamalipour, The Wireless Mobile Internet. Architectures, Protocols and Services, 1st ed. John Wiley, 2003, pp. 249-323.

[73]    3GPP TS 24.008: "Mobile radio interface Layer 3 specification; Core network protocols; Stage 3 v10.3.0 (Release 10)", June 2011.

[74]    A. Chandra and K. Mal, "Genetic algorithm based optimization for location update and paging in mobile networks", in Proceedings of 2004 Asian Applied Computing Conference (AACC), pp.222-231.

[75]    B. Liang and Z.J. Haas, "Predictive distance-based mobility management for PCS networks", in 18th Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM), 1999, vol. 3, pp.1377-1384.

[76]    Y. Tseng, L. Chen, M. Yang, and J. Wu, "A stop-or-move mobility model for PCS networks and its location-tracking strategies", Computer Communications, Volume 26, Issue 12, 21 July 2003, pp. 1288-1301.

[77]    NETGEAR Femtocell Voice Gateway DVG834GH, datasheet [Online]. Available: http://www.netgear.com

[78]    R. Langar, N. Bouabdallah, and R. Boutaba, "A comprehensive analysis of mobility management in MPLS-based wireless access networks," IEEE/ACM Trans. Netw., vol. 16, no. 4, pp. 918–931, 2008.

[79]    Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2010-2015, Whitepaper, February 2011

[80]    K. Lee, Y. Yi, J. Lee, I. Rhee, S. Chong, "Mobile Data Offloading: How Much Can WiFi Deliver?", in proc. of ACM CoNEXT 2010, USA, December, 2010.

[81] M. Taqqu, W. Willinger, and R. Sherman, "Proof of a fundamental result in self-similar traffic modeling", Computer Communications Review 26:5-23, 1997.

[82] Stoev, G. Michailidis, and J. Vaughan, "On Global Modeling of Network Traffic", INFOCOM 2010, The 29th Conference on Computer Communications, San Diego, California, March 2010.

[83] 3GPP TR 36.806, Evolved Universal Terrestrial Radio Access (E-UTRA), Relay architecture for E-UTRA (LTE-Advanced), Release 9, V9.0.0 (2010-03).

[84] 3GPP TR 36.912, Feasibility study for Further Advancements for E-UTRA (LTE-Advanced), Release 10, V10.0.0 (2011-03).

[85] 3GPP TS 36.300, Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN), Overall description V10.4.0 (2011-06).

[86] 3GPP TR 25.912, Feasibility study for evolved Universal Terrestrial Radio Access (UTRA) and Universal Terrestrial Radio Access Network (UTRAN), Release 10, V10.0.0 (2011-03)

[87] 3GPP TSG-RAN WG3 #66, R3-092738, "Architecture Options Comparison: UE Mobility Support", Fujitsu, November, 2009.

[88] ARTIST4G (Advanced Radio InTerface TechnologIes for 4G SysTems), "D3.4 - Relay configurations", July, 2011.

[89] S. Mohan, R. Kapoor, B. Mohanty, "Latency in HSPA Data Networks", white paper of Qualcomm, February, 2011.

[90] ARTIST4G (Advanced Radio InTerface TechnologIes for 4G SysTems), "D3.2 - Advanced Relay Technical Proposals", February, 2011.

[91] K. Samdanis, D. Kutscher, M. Brunner: Self-Organized Energy Efficient Cellular Networks. 21st Annual IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC 2010); September 2010.

[92] 3rd Generation Partnership Project; Technical Specification Group Radio Access Network; UTRAN architecture for 3G Home Node B (HNB); Stage 2 (Release 10), 3GPP Recommendation 3GPP TS 25.467, Sep 2011.

[93] R. Stewart, Q. Xie, M. Tuexen, S. Maruyama, and M. Kozuka, "Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration," RFC 5061 (Proposed Standard), Internet Engineering Task Force, September 2007. [Online]. Available: http://www.ietf.org/rfc/rfc5061.txt.

[94] R. Stewart, M. Ramalho, Q. Xie, M. Tuexen, and P. Conrad, "Stream Control Transmission Protocol (SCTP) Partial Reliability Extension," RFC 3758 (Proposed Standard), Internet Engineering Task Force, May 2004. [Online]. Available: http://www.ietf.org/rfc/rfc3758.txt.

[95] J. Fitzpatrick, S. Murphy, M. Atiquzzaman, and J. Murphy, "ECHO: A Quality of Service Based Endpoint Centric Handover Scheme for VoIP," in Wireless Communications and Networking Conference, 2008. WCNC 2008. IEEE, 31 2008-april 3 2008, pp. 2777 –2782.

[96] R. Stewart, "Stream Control Transmission Protocol," RFC 4960 (Proposed Standard), Internet Engineering Task Force, September 2007. [Online]. Available: http://www.ietf.org/rfc/rfc4960.txt.

[97] M. Olsson, S. Sultana, S. Rommer, L. Frid, C. Mulligan. "SAE and the Evolved Packet Core: Driving The Mobile Broadband Revolution," Academic Press, 2009.

[98] Golaup, A.; Mustapha, M.; Patanapongpibul, L.B.; "Femtocell access control strategy in UMTS and LTE" IEEE "Communications Magazine,: September 2009, vol: 47 Issue:9, pp 117 – 123.

[99] 3GPP TS 25.367, Mobility procedures for Home Node B (HNB); Overall description; Stage 2 V10.0.0 (2011-03).

[100] 3GPP TS 24.285Allowed Closed Subscriber Group (CSG) list; Management Object (MO) V10.2.0 (2010-12).

[101] 3GPP TS 22.220, Service requirements for Home Node B (HNB) and Home eNode B (HeNB) V11.3.0 (2011-10).

[102] 3GPP TS 22.011, Service accessibility V11.1.0 (2011-10).

[103] 3GPP TS 23.060, General Packet Radio Service (GPRS); Service description; Stage 2 V11.0.0 (2011-12).

[104] 3GPP TS 23.401, General Packet Radio Service (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access V11.0.0 (2011-12).

[105] 3GPP TS 29.002Mobile Application Part (MAP) specification V11.1.0 (2011-12).

[106] 3GPP TS 29.272Evolved Packet System (EPS); Mobility Management Entity (MME) and Serving GPRS Support Node (SGSN) related interfaces based on Diameter protocol V11.1.0 (2011-12).

[107] 3GPP TS 23.003, Numbering, addressing and identification V11.0.0 (2011-12).

[108] Gavin Horn "3GPP Femtocells:Architecture and Protocols," QUALCOMM Incorporated White paper, Sept. 2010.

[109] C. Kaufman, P. Hoffman, Y. Nir, P. Eronen "Internet Key Exchange Protocol Version 2 (IKEv2)," IETF RFC 5996, Sept. 2010.

[110] M. Myers, H. Tschofenig. "Online Certificate Status Protocol (OCSP) Extensions to IKEv2," IETF RFC 4806, Feb. 2007.