

**SELF-HEALING SOLUTIONS FOR LTE EVOLVED PACKET
CORE**

by

Md. Mustafizur Rahman

A thesis submitted to

the Faculty of Business and Information Technology

In conformity with the requirements for the degree of MSc in Computer Science

University of Ontario Institute of Technology

Oshawa, Ontario, Canada

(July, 2012)

Copyright ©Md. Mustafizur Rahman, 2012

Abstract

The 3GPP Long Term Evolution (LTE) is considered as a dominant future cellular wireless technology in terms of performance and user experience. With technological advancement of the wireless networks, dependencies and business impact of the mobile network services have increased phenomenally. It is, therefore, crucial to address the issues regarding network infrastructure or service failure. In this thesis, a self-healing solution is presented for the LTE Evolved Packet Core (EPC) with a view to maintaining service continuity in the event of core network elements - the MME and S-GW failures. The core network element failures have significant impact on a larger number of subscribers in comparison to the access network element failures. In the proposed self-healing scheme, the restoration mechanisms and associated failover recovery procedures with regards to service survivability are described in details from the LTE network and protocol perspective.

This thesis studies two different self-healing approaches - the centralized active-backup and distributed active-active and conducts simulation for each approach in various failure scenarios. The performances of each of these scenarios are evaluated in terms of service restoration time, throughput, EPS (Evolved Packet System) bearer delay etc. The results show that the proposed self-healing system can ensure service continuity at a certain level if resources are properly provisioned. And in terms of restoration delay, in general, the active-backup configuration performs better than the active-active configuration.

The thesis presents analytical and simulation methods to estimate signaling message overhead at the LTE EPC that arises due to the recovery process. It also analyzes the bandwidth requirements of the signaling traffic that is incurred by the other operational procedures of the self-healing scheme and their ramification to the LTE core network.

Acknowledgements

I would like to take this opportunity to express my sincere appreciation to my thesis supervisor Dr. Shahram S. Heydari for his expert guidance and valuable feedback throughout the research. I gratefully acknowledge the role of the Natural Sciences and Engineering Research Council of Canada (NSERC) in funding this project. I deeply thank my parents, sister-in-law and brothers for their continual encouragement and support. Finally, many thanks go to all of those who helped me in any respect during the completion of this thesis.

Table of Contents

Abstract.....	ii
Acknowledgements.....	iii
Chapter 1.....	1
1.1 Overview.....	1
1.2 Background.....	3
1.2.1 Evolution of Cellular Wireless Networks.....	3
1.2.2 Global System for Mobile Communications (GSM).....	4
1.2.3 General Packet Radio Service (GPRS).....	6
1.2.4 Universal Mobile Telecommunications System (UMTS) and High Speed Packet Access (HSPA).....	7
1.2.5 Long Term Evolution (LTE).....	9
1.2.6 Long Term Evolution – Advanced (LTE-A).....	11
1.3 Motivation.....	12
1.4 Literature Survey and Related Works.....	13
1.5 Research Objectives.....	18
1.6 Approach.....	19
1.7 Thesis Contribution.....	20
1.8 Thesis Outline.....	21
Chapter 2 LTE EPC Architecture and Protocols.....	22
2.1 The EPS Architecture.....	22
2.1.1 The Evolved UMTS Terrestrial Radio Access Network (E-UTRAN).....	23
2.1.2 The Evolved Packet Core (EPC) Network.....	24
2.1.3 LTE Interfaces.....	27
2.2 LTE Network Protocol Architecture.....	28
2.2.1 Description of Protocols Used in Control Plane and User Plane.....	29
2.3 LTE EPS Bearer and QoS.....	34
2.4 UE Context Information.....	37
2.5 S1-Flexibility (S1-Flex).....	38
2.6 LTE Self-Organizing Network (SON).....	39
2.6.1 Main functionalities of SON in LTE.....	40
2.6.2 LTE SON Architecture.....	42
Chapter 3 LTE Self-healing Core Proposal.....	45

3.1 Self-healing System Overview.....	45
3.1.1 Fault-tolerant System Architecture	48
3.1.1.1 N:M Active-Backup Configuration	49
3.1.1.2 1:1 Active-Active Configuration	51
3.1.2 Failure Detection, Notification and Fault Isolation.....	53
3.1.3 Failover Recovery Coordination.....	54
3.1.4 Recovery Mechanism and Service Continuity Procedures	56
Chapter 4 Performance Evaluation of the Self-healing Scheme.....	61
4.1 Simulation Setup.....	61
4.2 Performance Metrics	66
4.2.1 Service Restoration Time.....	66
4.2.2 Signaling Message Overhead.....	68
4.3 Results and Analysis	70
4.3.1 N:1 Active-Backup Configuration	70
4.3.2 1:1 Active-Active Configuration	77
4.3.3 Service Restoration Time Calculation	79
4.3.4 Signaling Message Overhead Estimation	81
Chapter 5 Bandwidth Analysis	83
5.1 Descriptions of the Logical Connections	83
5.2 Description of the Signaling Messages	84
5.2.1 Definitions of Information Elements (IEs)	88
5.3 Analytical Model for Bandwidth Calculation	89
5.4 Numerical Results for Bandwidth Calculation	92
Chapter 6 Conclusion and Future Works	94
References.....	96

List of Figures

Figure 1-1: Evolution of GSM based cellular wireless network.....	3
Figure 1-2: GSM network architecture.	5
Figure 1-3: GSM-GPRS network architecture.....	6
Figure 1-4: UMTS/HSPA network architecture.	8
Figure 1-5: LTE Network Architecture.	11
Figure 2-1: Evolved Packet System (EPS) architecture.....	23
Figure 2-2: LTE E-UTRAN architecture [2].	24
Figure 2-3: Functional split between the E-UTRAN and the EPC [2].	26
Figure 2-4: Control plane protocols stack.....	28
Figure 2-5: User plane protocols stack [38].....	29
Figure 2-6: Protocol architecture between the eNodeB and the MME.....	31
Figure 2-7: Inter-connected eNodeB protocol architecture.	32
Figure 2-8: GTP-C over S11 and S5/S8 interfaces.	33
Figure 2-9: GTP-U over S1-U and S5/S8 interfaces.....	33
Figure 2-10: LTE EPS bearer architecture [2].	34
Figure 2-11: LTE EPS bearers and corresponding service data flows (SDFs).	35
Figure 2-12: Two unicast EPS bearers (GTP based S5/S8) [38].	36
Figure 2-13: LTE S1-flex configuration.	39
Figure 2-14: LTE centralized SON architecture [51].	43
Figure 2-15: LTE distributed SON architecture [51].	43
Figure 2-16: LTE hybrid SON architecture [51].	44
Figure 3-1: Logical architecture of the Self-healing system in the centralized approach.....	47
Figure 3-2: Logical architecture of the Self-healing system in the distributed approach.	47
Figure 3-3: N:M Active-Backup configuration.	51
Figure 3-4: 1:1 Active-Active configuration.	53
Figure 3-5: Node status update for failover recovery coordination.	56
Figure 3-6: Interaction between the network entities during failover.....	58
Figure 3-7: Bearer-recreation and related message flows.....	59
Figure 4-1: Sample network configuration.	63
Figure 4-2: Message flows for UE originated sessions [53].	69
Figure 4-3: Uplink throughput for Case# 1.....	71
Figure 4-4: GTP traffic switchover from the active EPC to the backup EPC (Case# 1).	71

Figure 4-5: Uplink delay for Case# 1.....	72
Figure 4-6: EPS bearer delay of eNodeB_10 for Case# 1.	73
Figure 4-7: Uplink throughput for Case# 2.....	74
Figure 4-8: Uplink delay for Case #2.....	75
Figure 4-9: Link utilization of active (EPC_0) and backup (EPC_1) backhaul connections.....	75
Figure 4-10: EPS bearer traffic received by the eNodeB 10 for Case# 2.	76
Figure 4-11: Uplink delay for Case# 3.....	77
Figure 4-12: PDSCH channel utilization for Case# 3.....	77
Figure 4-13: Uplink delay for Case# 4.....	78
Figure 4-14: EPS bearer delay for Case# 4.....	78
Figure 4-15: Restoration delay comparisons for 1:1 Active-Backup configuration for different numbers of users.	81
Figure 4-16: Signaling message overhead due to session restoration.....	82
Figure 4-17: Signaling message overhead due to new UE arrivals during restoration time.	82
Figure 5-1: Connection links between the network entities for the self-healing system.	84
Figure 5-2: Message flows for UE context replication.	85
Figure 5-3: Failure notification message flows for the eNodeB.....	85
Figure 5-4: Failure notification message flows for the P-GW.....	85
Figure 5-5: Bandwidth usage due to UE context replication.	93

List of Tables

Table 1-1: LTE requirements [15]	10
Table 2-1: LTE standardized QCI [55].....	37
Table 2-2: NGMN SON use case definitions [48].....	41
Table 4-1: Network configurations in the simulation	62
Table 4-2: Simulation parameters and settings.....	64
Table 4-3: Measurement of restoration delay parameters.....	67
Table 4-4: Simulation results for average values of T_{rt} , T_{me} and T_{br}	79
Table 4-5: Simulation results for scenario 1 and 2.....	80
Table 5-1: Context replication message description.....	87
Table 5-2: Subscribers profile information.....	91
Table 5-3: Different user profile in UE arrival rate.....	92

List of Abbreviations

3GPP	Third Generation Partnership Project
AAA	Authentication, Authorization and Accounting
AuC	Authentication Centre
BS	Base Station
BSC	Base Station Controller
BTS	Base Station Transceiver
CAPEX	Capital Expenditure
CDMA	Code Division Multiple Access
CoMP	Coordinated Multiple Point
DL-SCH	Downlink Shared Channel
E-DCH	Enhanced DCH
EDGE	Enhanced Data Rates for GSM Evolution
EIR	Equipment Identity Register
eNodeB	Evolved NodeB
EPC	Evolved Packet Core
EPS	Evolved Packet System
E-UTRAN	Evolved Universal Terrestrial Radio Access
FDD	Frequency Division Duplex
FDMA	Frequency Division Multiple Access
GBR	Guaranteed Bit Rate
GERAN	GSM/EDGE Radio Access Network
GGSN	Gateway GPRS Support Node
GMSK	Gaussian Minimum Shift Keying
GPRS	General packet radio service
GSA	Global mobile Suppliers Association
GSM	Global System for Mobile Communications
GTP	GPRS Tunneling Protocol
GTP-C	GPRS Tunneling Protocol, Control Plane
GTP-U	GPRS Tunneling Protocol, User Plane
GUTI	Globally Unique Temporary Identity
GUMMEI	Globally Unique MME ID
HLR	Home Location Register

HPLMN	Home PLMN
HSDPA	High Speed Downlink Packet Access
HSPA	High Speed Packet Access
HSS	Home Subscriber Server
HSUPA	High Speed Uplink Packet Access
IMT-A	International Mobile Telecommunications – Advanced
ISDN	Integrated Services Digital Network
ITU	International Telecommunication Union
LTE	Long Term Evolution
LTE-A	LTE-Advanced
MAC	Medium Access Control
MBR	Maximum Bit Rate
MIMO	Multiple Input Multiple Output
MME	Mobility Management Entity
MS	Mobile Station
MSC	Mobile Switching Center
NAS	Non-access Stratum
OAM	Operation, Administration and Maintenance
OFDMA	Orthogonal Frequency Division Multiple Access
OPEX	Operational Expenditure
OSS	Operations Support System
PCRF	Policy and Charging Resource Function
PDCP	Packet Data Convergence Protocol
PDSCH	Physical Downlink Shared Channel
P-GW	Packet Data Network Gateway
PLMN	Public Land Mobile Network
PMIP	Proxy Mobile IP
QCI	QoS Class Identifier
QoS	Quality of Service
QPSK	Quadrature Phase Shift Keying
RLC	Radio Link Control
RNC	Radio Network Controller
RRC	Radio Resource Control
RRM	Radio Resource Management

RSRP	Reference Symbol Received Power
S1AP	S1 Application Protocol
SAE	System Architecture Evolution
SC-FDMA	Single Carrier Frequency Division Multiple Access
SCTP	Stream Control Transmission Protocol
SDU	Service Data Unit
SGSN	Serving GPRS Support Node
S-GW	Serving Gateway
SIM	Subscriber Identity Module
SON	Self Organizing Networks
S-TMSI	S-Temporary Mobile Subscriber Identity
TA	Tracking Area
TCP	Transmission Control Protocol
TDD	Time Division Duplex
TDMA	Time Division Multiple Access
UDP	Unit Data Protocol
UE	User Element
UMTS	Universal Mobile Telecommunications System
UTRA	Universal Terrestrial Radio Access
UTRAN	Universal Terrestrial Radio Access Network
VLR	Visitor Location Register
VOIP	Voice Over IP
WCDMA	Wideband Code Division Multiple Access
WLAN	Wireless Local Area Network
X1AP	X1 Application Protocol

Chapter 1

Introduction

1.1 Overview

In order to meet the tremendous growth of high speed mobile data traffic and quality of service requirements, new cellular wireless technology standards based on the third (3G) and fourth generation (4G) have been developed. The 3GPP (Third Generation Partnership Project) Long Term Evolution (LTE) is the latest standard with its origin in the GSM/GPRS legacy networks family which is currently being deployed worldwide. According to a report of the Global Mobile Suppliers Association (GSA), LTE is the fastest growing mobile technology ever [1]. As of August 31, 2011, 26 LTE networks launched in 18 countries and it is expected that by the end of 2012, 93 LTE networks will enter in commercial service [1]. LTE is an important evolution of the 3G UMTS/HSPA network and a significant step towards 4G technologies. It offers improved spectrum efficiency, higher data rates, low latency and increased system capacity. LTE provides downlink data rates up to 100 Mbps and uplink data rates up to 50 Mbps [2]. Further enhancement of LTE, known as LTE-Advanced (LTE-A) promises downlink data rates up to 1 Gbps and uplink data rates up to 500 Mbps [3].

LTE is the next step in the user experience, capable of supporting demands of bandwidth-intensive mobile broadband services and applications such as interactive TV, video streaming, online gaming etc. The first version of LTE was documented in 3GPP release 8 specifications [2]. The 3GPP LTE project was formerly known as evolved UMTS terrestrial radio access (E-UTRA) and evolved UMTS terrestrial radio access network (E-UTRAN) which defines a brand new physical layer radio access technology. The release 8 also specifies a new all-IP, packet based

core network called the Evolved Packet Core (EPC). The EPC is designed to be backward compatible with earlier 3GPP releases and capable of supporting non-3GPP radio access technologies [4]. 3GPP release 9 improves some specific features of LTE including emergency services, support of circuit switch calls over LTE etc. [5]. Furthermore, 3GPP release 10 specifies the evolution of LTE, referred to as LTE-Advanced (LTE-A). The LTE-A is a candidate technology of International Mobile Telecommunications – Advanced (IMT-A), which is a set of requirements defined by the International Telecommunication Union (ITU) for the 4G cellular wireless networks [3].

In 2010, an ITU press release stated that the number of mobile subscriptions was expected to surpass one billion by the first quarter of 2011 [6]. With this rapid advancement of wireless technology and continued subscriber growth, societal and business dependencies on the mobile services have increased in numerous forms. As a result, the impact of network infrastructure failure or service disruption has become critical for the subscribers as well as carriers in terms of emergency, revenue loss, market perception and economic impact. Examples can be drawn from; T-Mobile, Germany's nationwide network outage incident, occurred on April 23, 2009, that affected almost 40 million subscribers [7]. Estimated loss was approximately 100 million dollars. Furthermore, catastrophic events such as earthquake, fire, floods, explosions etc., can cause severe service disruption and data loss. The subscribers now expect improved QoS and reliable services comparable to the wired telecommunication networks. It is probably not an overstatement to say that user demands and expectations on the mobile networks have far more increased with the advent of 4G technologies like LTE. Therefore, ensuring high degree of reliability similar to the wire-line network infrastructure has become a necessity. In order to meet

this expectation, it is essential to design robust wireless mobile network so that it can deliver necessary performances for the communication services in the event of unexpected failures.

This thesis specifically focuses on the issues related to the LTE core network elements (NEs) failures and hence, proposes and implements self-healing solutions for the LTE EPC with a view to maintaining service continuity. It should be mentioned here that the research presented in this thesis is based on the 3GPP release 8 specifications.

1.2 Background

An overview of the evolution of cellular wireless networks is presented in the following section.

1.2.1 Evolution of Cellular Wireless Networks

The first Generation (1G) cellular wireless technology, developed by the Bell Labs, went into commercial service in the 1980s. A common 1G wireless standard was called Advanced Mobile Phone System (AMPS). The AMPS was created to produce a quick and efficient phone system. This effectively started a revolution of wireless technologies. It was one of the first wireless phone system running analog transmission.

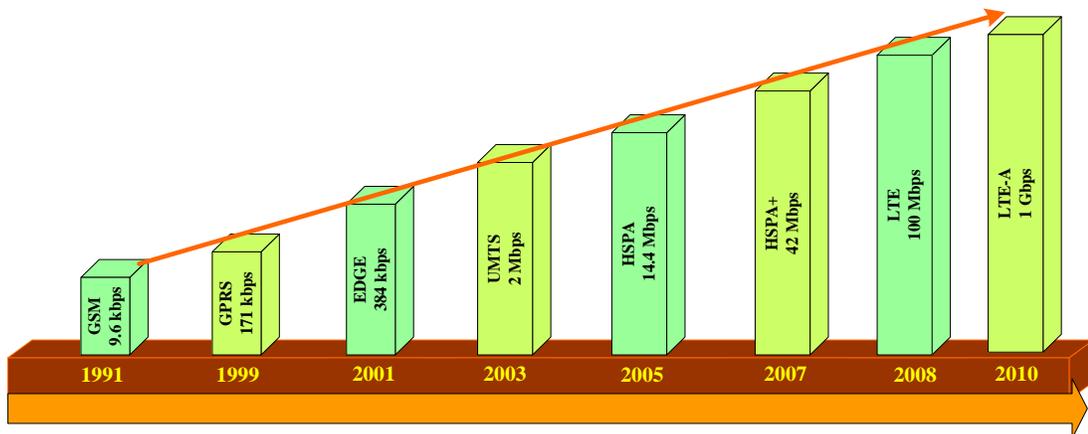


Figure 1-1: Evolution of GSM based cellular wireless network.

In the subsequent sections, an overview is presented on the 2G, 3G and 4G technologies of the GSM family and their evolution. Figure 1-1 illustrates the evolution path of the GSM. Each of the generation has entirely new technology developed or enhancements of previous technologies are adopted.

1.2.2 Global System for Mobile Communications (GSM)

In the 1990s the second generation (2G) of mobile system was introduced. One of the main characteristics of the 2G mobile system is that the phone conversation become digitized and compressed which allows more subscribers to be accommodated in the air spectrum. The Global System for Mobile Communications (GSM) was one of the dominant 2G standards. It was standardized in Europe by the European Telecommunications Standard Institute (ETSI). Figure 1-2 depicts the GSM network architecture. Generally, the GSM operates at 900 MHz and 1800 MHz frequency band. As access technology, it uses combination of Time Division Multiple Access (TDMA) and Frequency Division Multiple Access (FDMA), and Gaussian Minimum Shift Keying (GMSK) for modulation [8]. It provides basic and extended phone services such as international roaming, interoperability with Integrated Services Digital Network (ISDN), authentication and security.

In GSM networks, the user device is called a mobile station or MS. The MS includes a card that allows the user to be uniquely identified (Subscriber Identity Module, SIM), and a device (phone). The SIM card allows each user to be identified regardless of the terminal used while communicating with the base transceiver station (BTS). When the user engages the mobile device, the communication path occurs through a radio link between the mobile station (users' phone) and the base station (network transmitter). Each BTS is assigned a set of channels, and neighboring base stations are assigned different sets to avoid interference. A group of BTS in a

cellular network are connected to a base station controller (BSC) that determines resource allocations.

The BSC and the BTS are together referred to as the Base Station Subsystem (BSS). Each BSC is connected to a Mobile Switching Center (MSC), which belongs to a Network Station Subsystem (NSS) responsible for call control and establishing communication with other users.

The NSS also contains the network elements HLR, VLR, AuC and EIR. The HLR or Home Location Register retains permanent information of users including subscriber identity number and subscribed services, location information. The Visitor Location Register (VLR) maintains temporary information of users, for example, current location of a subscriber's service area. The AuC or Authentication Centre authenticates the registered users and generates encryption key for security purpose. The Equipment Identity Register (EIR) database stores the identity of the MS.

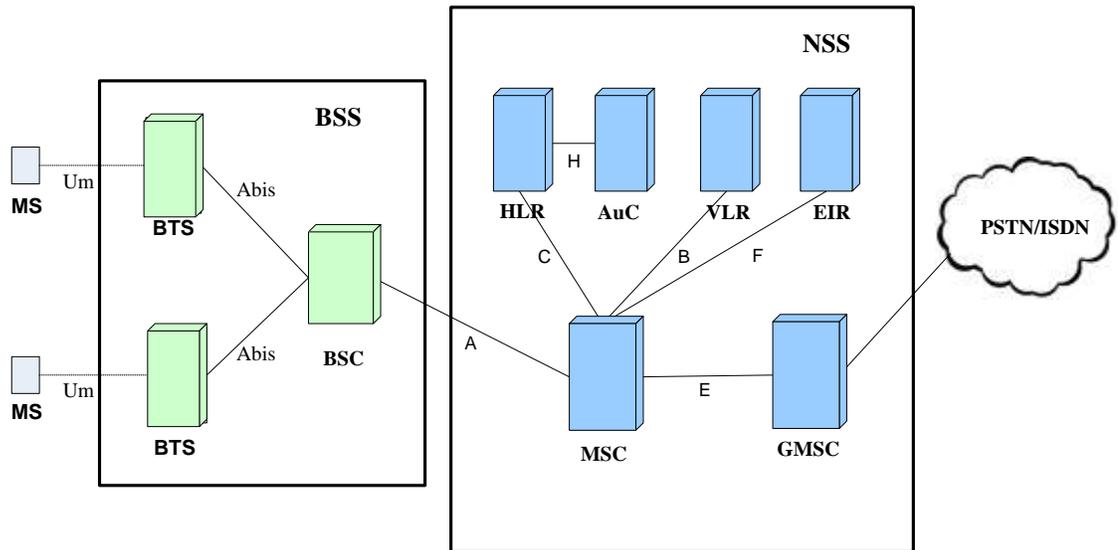


Figure 1-2: GSM network architecture.

As users move around, geographically, the mobiles leave the transmission range of one BTS to enter the range of another. The user's mobile periodically check the signal levels of surrounding BTSs and based on the signal strength chooses a new BTS to obtain service. This process is

called handover. GSM also supports roaming; the movement of users from one operator network to another.

1.2.3 General Packet Radio Service (GPRS)

Due to the tremendous growth of Internet in the late 1990s, subscribers demand for high speed wireless data services were increasing in the cellular network. But the GSM data service capability was limited to 9.6 kbps [9]. Furthermore, resource utilization in the circuit switching was inefficient for the bursty data traffic. To address these issues, a packet oriented data service called GPRS was introduced in the GSM network by the ETSI. The GPRS is sometimes referred to as 2.5G technology. In the access network part, GPRS uses the same air interface or radio access technology as in GSM.

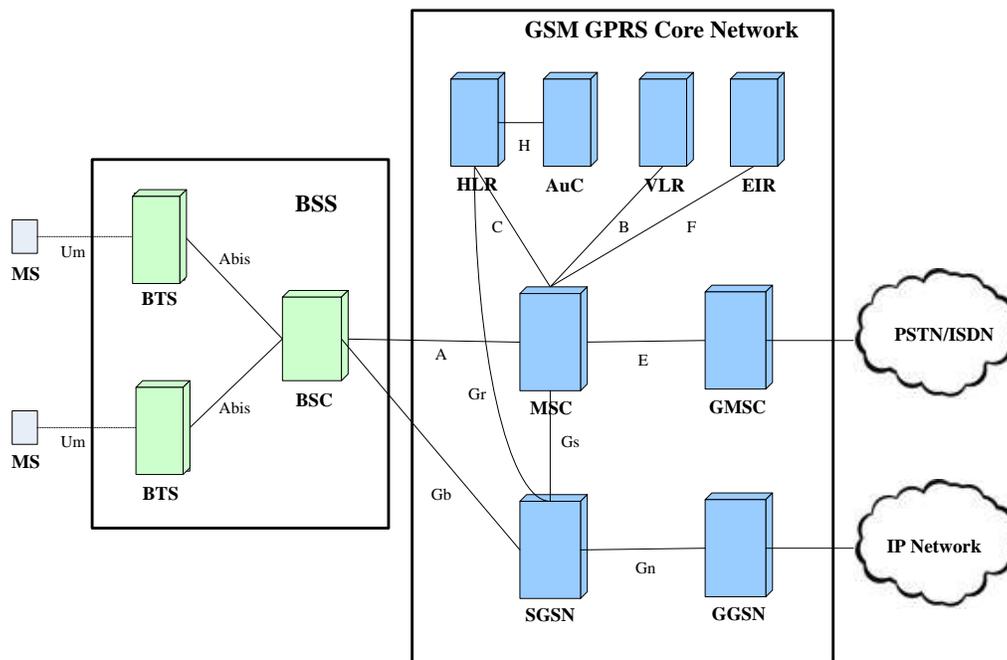


Figure 1-3: GSM-GPRS network architecture.

Two new functional elements - the Serving GPRS Support Node (SGSN) and Gateway GPRS Support Node (GGSN) are added to the GSM core network architecture to support the packet

switching. GPRS also facilitates efficient tariff system for various services using the volume based charging as opposed to the time duration charging used in GSM. Figure 1-3 shows the GPRS system architecture.

The SGSN is responsible for establishing and controlling communication between the MS and GPRS network. It handles ciphering, session management, mobility management functions such as paging, roaming. It routes data to the relevant GGSN when the external network connection is required. The GGSN is the external gateway of the GPRS network and routes data to and from the external network and SGSN. The GPRS offers up to 171 kbps downlink speed [10]. A further enhancement of GPRS, called the EDGE improved the data rates up to 384 kbps [5].

1.2.4 Universal Mobile Telecommunications System (UMTS) and High Speed Packet Access (HSPA)

UMTS is the 3G evolution of the GSM/GPRS networks and specified by 3GPP in Release 99 [5]. Driving forces behind the 3G evolution included increased demands for faster data rates, higher capacity, improved quality-of-services (QoS), and mobile roaming between different cellular technologies. UMTS is capable to provide data rates up to 2 Mbps in picocellular environments and 384 kbps in microcellular environments [11]. It introduced a new access network, called UMTS Terrestrial Radio Access Network (UTRAN) using Wideband CDMA (WCDMA) radio technology. Two new functional elements are introduced in the UTRAN - the Node B and Radio Network Controller (RNC) which replace the functionalities of the BTS and BSC of earlier GSM network. The functionalities of existing core network elements, such as, MSC, SGSN, HLR are extended to adopt the UMTS standards. The GSM/UMTS dual-mode user element (UE) is used to connect to the GSM networks via the GSM air interface and UMTS network via the UMTS

radio interface. UMTS requires an entirely new spectrum range and supports both Frequency-Division and Time-Division Duplex. Figure 1-4 illustrates the UMTS architecture.

The network access technologies are further improved by the introduction of High Speed Packet Access (HSPA). HSPA includes two evolutions of W-CDMA packet data access technology - High Speed Downlink Packet Access (HSDPA) and High Speed Uplink Packet Access (HSUPA). HSDPA, standardized in 3GPP release 5, is the first evolution of W-CDMA. It increases downlink data rates up to 14 Mbps which is a significant performance improvement of 384 kbps data transfer rate of WCDMA [12-13]. HSDPA uses a new transport channel, known as High Speed Downlink Shared Channel (HS-DSCH).

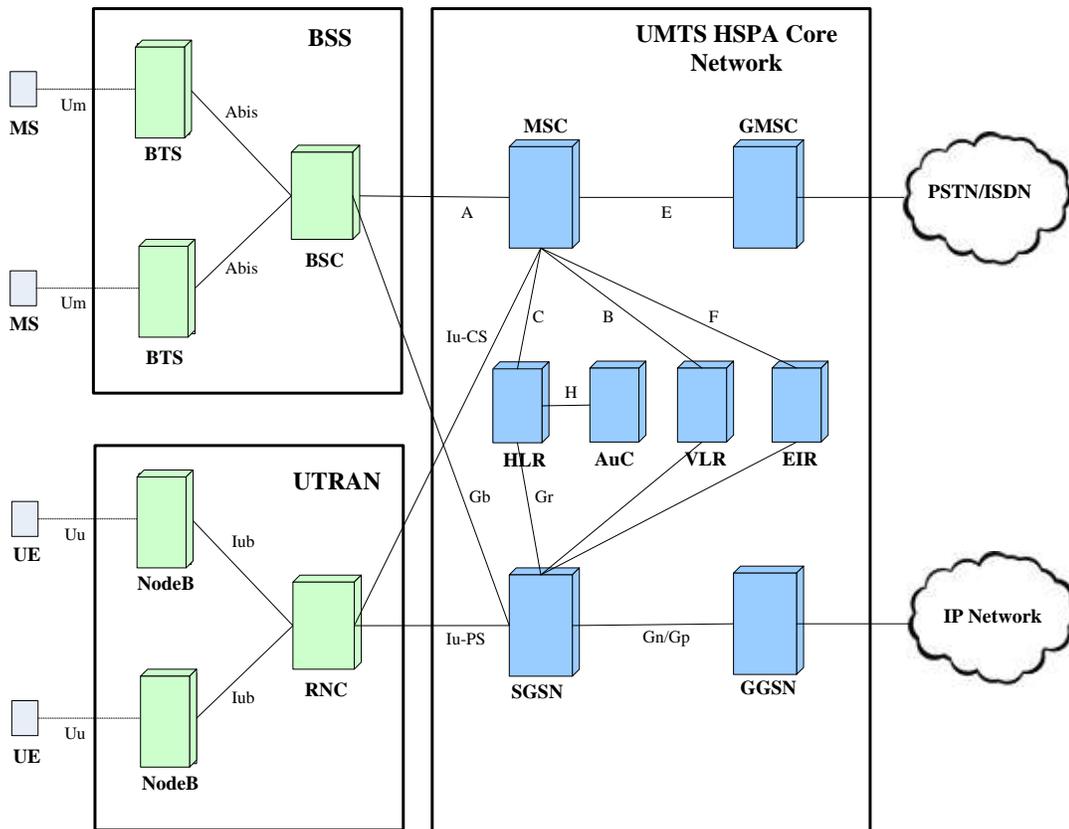


Figure 1-4: UMTS/HSPA network architecture.

HSUPA was specified in 3GPP release 6 specifications. It increases the uplink data rates up to 5.8 Mbps and improves network coverage. Like HSDPA, HSUPA includes a new transport channel to WCDMA, called Enhanced Dedicated Channel (E-DCH) in order to improve the data rates at the uplink side. The HSPA+ is a term used to represent the evolution of the HSPA technologies. It promises higher data rates up to 42 Mbps for downlink and 28 Mbps for uplink, and further reduces cost of the voice and data services [14]. The HSPA+ was defined by 3GPP in release 7.

1.2.5 Long Term Evolution (LTE)

The LTE is the evolution of the UMTS/HSPA network. Motivation for LTE is to further improvement of data rates and coverage provided by the HSDPA/HSUPA technology. As mentioned in the beginning, LTE aimed to achieve peak data rates up to 100 Mbps for downlink and 50 Mbps for uplink, improved spectrum efficiency, less latency time and reduced cost for the user as well as for the operators [2]. Table 1-1 shows the major requirements of LTE.

It introduced a new access network architecture known as Evolved Universal Terrestrial Radio Access Network (E-UTRAN). In LTE, Orthogonal Frequency Division Multiple Access (OFDMA) is adopted as the access technology for the downlink while Single Carrier Frequency Division Multiple Access (SC-FDMA) is used for the uplink. Unlike other 3GPP technologies, the LTE E-UTRAN is a simple and flat architecture that consists of only one node, eNodeB which performs functionalities of the UMTS/HSPA Node B and Radio Network Controller (RNC). Figure 1-5 shows the LTE network architecture. The main benefit of this flat architecture is that it reduces the number of levels and nodes in the network, which in turn reduces call processing time, latency, and cost.

The OFDMA uses a large number of narrow-band and parallel sub-carriers for transmitting data. Each of these sub-carriers transmits data using 64-QAM, 16-QAM or QPSK modulation scheme considering radio link condition [13, 16]. OFDMA significantly reduces multipath interference

Table 1-1: LTE requirements [15]

Bandwidth (MHz)	Scalable bandwidths of 1.25, 2.5, 5, 10, 15, 20 MHz
Data Rates	Downlink peak data rates of 100 Mbps (5 bps/Hz) and uplink data rates of 50 Mbps (2.5 bps/Hz) within 20 MHz uplink and downlink spectrum.
Latency (ms)	Control plane latency less than 100 ms. User plane latency less than 5 ms.
Capacity (control plane)	At least 200 users for spectrum allocation of 5 MHz.
Mobility	Up to 500 km/h depending on frequency band and optimized support for 0 to 15 km/h.
Spectrum efficiency	3 to 4 times higher spectrum efficiency in downlink and 2 to 3 times higher spectrum efficiency in uplink compared to Rel. 6 HSDPA

effects and highly suitable for advanced antenna techniques. Flexible use of the spectrum is one of the important features of LTE radio access technology. It supports carrier bandwidth from 1.4 MHz to 20 MHz. Cellular communication network's frequency bands may vary depending on the areas or countries. LTE is flexible enough to operate in different frequency bands. Multiple antenna schemes plays a vital role to increase the coverage, capacity and data rates in uplink and downlink which cannot be gained in the single antenna system. Various multiple antenna schemes

are supported in LTE such as, multiple-input-multiple-output (MIMO), beam forming etc. As already mentioned in the beginning, 3GPP also defined a flat core network architecture, called the Evolved Packet Core or EPC along with LTE. A detail description of the EPC architecture and protocols is given in chapter 2.

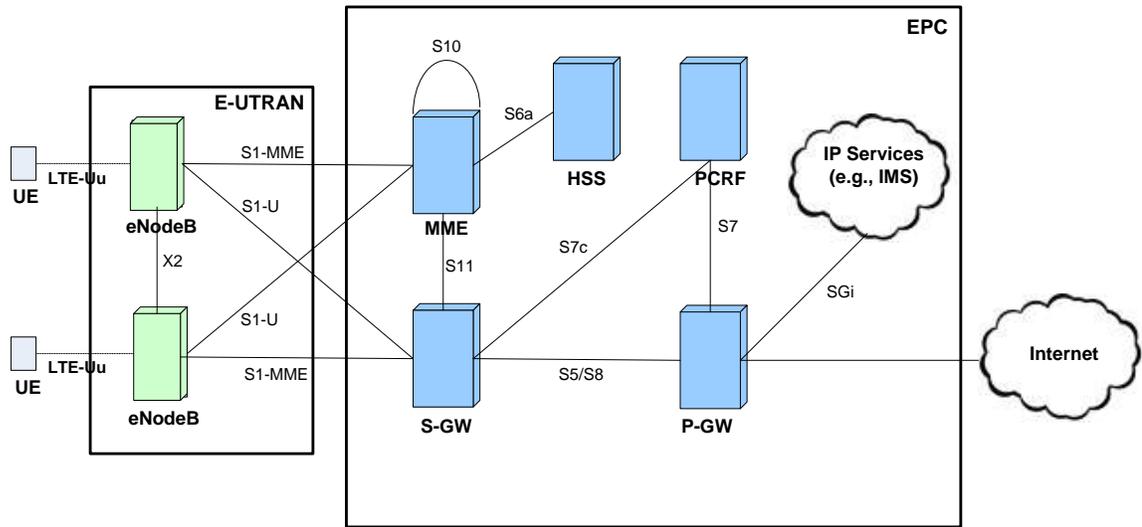


Figure 1-5: LTE Network Architecture.

1.2.6 Long Term Evolution – Advanced (LTE-A)

The LTE Advanced (LTE-A), standardized by 3GPP Release 10, is the evolution of LTE first release (3GPP release 8) for 4G cellular networks. The ultimate goal of LTE-A is to provide significant improvements in terms of capacity and coverage, lower latency and deployment cost in comparison to other 3GPP releases. The LTE-A is backward compatible with the LTE E-UTRAN. New features in LTE-A include [3, 17]:

- **Multi-carrier Aggregation:** The transmission bandwidth expansion from 20 MHz to 100 MHz has been considered in LTE-A as an effective way to ensure high data rates and capacity. The carrier aggregation technology facilitates this extension by combining multiple available carriers.

- **Multi-Antenna Technology:** Improved antenna technologies continues to be a dominant part in LTE-A. The LTE-A includes the multi-user MIMO (MU-MIMO) along with the single-user MIMO (SU-MIMO) technique. Using the MU-MIMO technology, it is possible to transmit data to more than one user through the same frequency band at a time.
- **Coordinated multipoint transmission or reception (CoMP):** This concept has been introduced in LTE-A where the cell edge user's data transmission and reception are coordinated by multiple eNodeBs. This technique increases performance and coverage.

LTE-A network coverage and capacity can be further improved by deploying low power nodes scattering around the network. The Femtocell, picocell and relay stations can be useful in this purpose.

1.3 Motivation

The Self-Organizing Network (SON) concepts are introduced in 3GPP release 8 along with the LTE standards. It is considered as one of the important areas of research to improve the efficiency of the network Operation, Administration and Maintenance (OAM). In fact, SON is envisioned to be the new model for the next generation OSS (Operations Support System) [18]. Subsequent 3GPP releases have expanded the scope of SON. Main functionalities of SON are: i) self-configuration, ii) self-optimization and iii) self-healing. This research focuses on the self-healing feature of SON. In general, the purpose of the self-healing functions is to detect failure and applies healing mechanisms in order to recover services from failures. Some self-healing approaches for LTE are available in 3GPP specifications, technical draft report and literatures [19-20] for access network but to the best of our knowledge, core network self-healing issues have not been widely studied from the perspective of service continuity and disaster recovery. It

should be mentioned here that the LTE access network consists of eNodeBs which cover a small geographic area. As a result, failures in the access networks affect small numbers of subscribers. On the other hand, infrastructure failures in the core network affect large number of subscribers in comparison to the access network and the impact of network disruption is significant. By minimizing the impact of failure in the core network, the overall performance of the network can be greatly enhanced. In fact, the vision of this thesis is to increase the degree of network availability by adopting the failure restoration process discussed here as a function of self-healing aspect of SON.

1.4 Literature Survey and Related Works

The term self-healing refers, in general, to the ability of a system to cope up with unexpected situations. Ghosh et al. [21] provides a more specific definition of self-healing systems which describes that self-healing characteristics enables a system to recognize the faults and make necessary adjustments in order to automatically recovering itself to the normal operation from malfunction. Ganek and Corbi [22] include self-healing as one of the main four fundamental features for future autonomic computing systems. They state that the purpose of self-healing is to minimize the effect of outages in order to maximize reliability and availability of the system. Kant [23] defines self-healing from networks points of view which describes self-healing as the capability of the network to provide continuous service to the end nodes by restoring affected services during network links and nodes failure incidents. There are other terms in the literature that are closely related to the self-healing concept – most importantly, fault-tolerance and survivability. A fault-tolerant system is designed to provide service seamlessly or with no loss of service in the presence of component failure [21]. Most fault-tolerant systems rely on replicating all functions so that another system may takeover in the case of failure. Survivability has been

defined as the ability of a system to maintain a certain level of service while the system experience failure [24]. There is a debate among the researchers whether self-healing should be considered as an independent study area or should be considered as a part of the traditional fault-tolerance system concepts. Gosh et al. [21] argues that though the fundamental concepts of these terms are highly similar but the self-healing system elaborates the recovery process to a greater degree which neither of these two systems does.

As the above discussion indicates, though each of these fields - fault-tolerance, survivability, reliability and self-healing are different and has their own domain and applicability in networking but their underlying purpose is almost similar – to increase the availability of the network services. In the literature review, we attempt to explore the knowledge, methods and mechanisms used in other fields along with the self-healing in order to understand the failure recovery process of the mobile systems.

Earliest efforts that address fault-tolerance of wireless networks can be found in Tipper et al. [25]. Their study proposed a multi-layer framework based on the second generation (2G) network architecture and carried out simulation for different kinds of failure scenarios. The concept of framework is an important factor in order to analyze the failure issues and design fault-tolerant network in a systematic way. Authors divide the whole network into three layers: access, transport and intelligent layer. Each layer has different functionalities and consists of specific nodes, links and transmission medium. The access layer comprises of nodes, links and functionalities that are responsible for the physical air interface. Important components of this level include, BS, BS-BSC links. Due to the tree-leaf structure of the 2G mobile wireless networks, failures of the BS or BS-BSC link affect a certain amount of subscribers resides in the cells. Additionally, one BS failure may impact adjacent BSs in the form of increased traffic due to

the user mobility. Moreover, congestion can be increased as failed or disconnected users may try several times to connect. The transport layer consists of network elements that are responsible for the call management functionalities, for instance, connection setup, call re-routing, and mobility management. The network elements in this layer include MSC, MSC-BSC links. Failure of the MSC affects subscriber of a wide area. The BSC-MSC link failure impacts entire subscribers under a particular BSC. The elements of the intelligent layers are HLR/VLR, service data management nodes (SCP, SDP). Loss of the VLR node impacts roaming functionalities of the users. Along with identifying the cause and effects of failures, authors also provided various metrics such as call blocking probability, call termination probability, call setup delay etc., to quantify the survivability of the network. Their simulation results suggest that the mobility plays a vital role in the network performance degradation during failure.

Snow et al. [24] proposed architectural changes to improve the reliability and survivability of wireless networks. Fault-tolerant SONET ring architecture can be very useful to connect the Mobile Switching Center (MSC) and Base Stations (BS). The SONET ring provides robust failure restoration mechanism and minimizes the impact of the single fiber cut or circuit system failure. Besides, the multimode or multifunction devices capable of operating in different technologies can be a useful way to improve network survivability. Based on the subscriber location information and availability of different access networks, devices can adapt with a particular network technologies to transmit data. In addition, universal access point based overlay networks can be a very useful way to ensure connectivity of the network. In this method, whenever a subscriber connects to the network, they first initiate connection with the universal access point. Depending on the network availability and other issues, the access point selects appropriate network connection for the subscribers.

In [23], a self-healing mechanism was presented and evaluated based on the cost-performance complexity, specifically for links failure restoration of different components connectivity of a GSM/GPRS network. The metrics that are used to assess the performance of the network in a post failure scenario includes blocking probability of existing users, blocking probability of new requests and restoration delay. Blocking probability of the existing users is defined by the number of affected requests that are not successfully recovered. Similarly, blocking probability of the new request is the fraction of new service requests that cannot be served by the network after failure. The performance data of the self-healing methods were collected through simulation.

In [26], the author proposes solutions for improving the scalability of the GGSN in a GPRS/UMTS network. The mechanism allows the GGSN to re-direct existing MSs to another GGSN in the event of excessive load. A new device is introduced, called the GGSN controller which task is to monitor the load of a GGSN and triggers re-direction of the MSs by first transferring the PDP contexts to a new GGSN if necessary. In this way, the GGSN controller dynamically adjusts the load of the GGSN. New protocol level messages were introduced to coordinate the load balancing and PDP contexts transferring process. The author also discusses various scenarios that may occur during the MS re-direction process. Unfortunately, this study didn't provide any performance evaluation data. Recently, Kustos et al. [27] conducted the performance evaluation of the above mentioned approach on a 3G/UMTS testbed, simplifying the architectural requirements and protocol messages. Their results show that the proposed scheme improves throughput and scalability of the packet switch domain as the average packet latency is decreased due to dynamically balancing the load.

3GPP specification [20] pointed out self-healing features and requirements of the OAM to facilitate the SON functionalities. This document mostly contains general information of the self-

healing concept, functions and logical architectures. Various strategies that can be adopted in recovery process of software and hardware faults are outlined. Some use cases related to the access network are mentioned and stages of the healing process are identified. The SOCRATES project objectives include developing methods for LTE SON [28]. However, use cases and scenarios for the self-healing that are discussed still limited to the LTE access network, for instance, coverage restoration in the case of cell site failure.

3GPP standards [29-30] discuss procedures to limit the effects of the UMTS and EPC core network elements failures and possible restoration mechanisms. In the case of EPC, standard [29] basically defines the functionalities of the peer nodes of the failed network elements in order to restore subscriber status and data to a consistent state. On the other hand, the technical draft [30] contains study on various EPC network elements failure scenarios. It analyzes the consequences of the failures and evaluates possible recovery procedures. So far the proposed solutions are based on the assumption that the affected node (MME, S-GW and P-GW) will be restarted after failure and recovery mechanism will take place when the failed node returns into the service. This study does not answer yet what would happen if the failure node does not restart or remain out of the service for a long period of time. In this context, proposed solutions can be described as reactive by nature.

In a recent paper [31] by Taleb and Samdanis, provides a service restoration for the MME failures. Restoration procedures are different for the active and idle mode UEs. This paper discusses an enhancement of the paging process where connected eNodeBs start sending paging message to the IDLE mode UEs when the MME fails. In order to minimize the effects of large numbers of UEs re-attachment, authors introduced bulk signaling method where eNodeBs aggregate signaling message of the UEs and send them to the MME using one single message to

avoid congestion in the MME side. It is assumed that parsing of bulk message will be faster than individual message which in turns increases efficiency. For the re-attachment, suitable MME is selected from the MME pool based on the load balancing feature of the eNodeB. For ACTIVE mode UE service restoration, contextual information availability of the failed UEs in the new MME is essential. It should be noted that this study does not include the S-GW node failure. As the S-GW already contains UE context information, the new MME can recover UE context from the S-GW. A message sequence diagram is presented which describes the UE re-attachment procedures. The bulk signaling method can also be applied in case of ACTIVE mode UE re-attachment. The performances of the proposed approach were evaluated through a set of metrics including signaling overhead against inter-arrival rate of UEs, TAU (Tracking Area Update) drop rate against inter-arrival rate of incoming UEs and the MME failure duration against the rate of incoming UEs.

1.5 Research Objectives

As stated earlier, this thesis proposed and analyzed self-healing approaches for LTE core network with a focus on the MME and S-GW failures. Along with ensuring automatic recovery from failures, one of the main design objectives of the proposed solution is to maintain seamless service continuity. We are also interested in failure scenarios where a major failure (e.g., a natural disaster) disables a site and any on-site redundancies. The self-healing proposal is comprised of independent/dependent functional entities such as fault-tolerant network architecture, fault detection and isolation, information dissemination, availability of crucial data, failover and recovery procedures. The proposed solutions address all the issues, particularly our main contribution is in developing methods and procedures with regard to service restoration

mechanism for LTE core network. The bandwidth requirements for the signaling traffic incurs by the design and its ramification to the LTE core network are also analyzed.

The specific goals of this research include:

- To determine core network elements failure effects and propose solution in order to recover services.
- To identify the efficient self-healing architectural configurations in order to achieve a certain degree of reliability based on the operator implementation policy.
- To implement associate failure recovery mechanisms which conform to existing LTE architecture and protocols.
- To identify possible complexities and available technologies to overcome the hurdle of real life implementation.

1.6 Approach

In the literature survey a detail study of the LTE network protocols and architecture was carried out for relevant information gathering. Existing self-healing, fault-tolerance and service survivability mechanisms of similar cellular wireless and wire-line technologies were investigated. It should be noted that inherently, the cellular wireless network has some unique characteristics and influenced by the factors like air interface, interference, limited frequency spectrum, population density, random user mobility etc. [25]. As a result, well-investigated self-healing mechanisms of the wire-line technologies may not be directly applicable here or the outcome of those mechanisms may significantly differ in wireless networks.

For the deployment scenarios of the EPC network in this thesis, it is assumed that the MME and S-GW functions are implemented in the same platform. This type of combined solution strategy is common among the vendors. In the failure recovery model, it is assumed that both MME/S-

GW fail ungracefully at the same time without informing other connected NEs and remain out of the service. In order to tackle the large scale failure scenarios, the proposed self-healing solution utilized remote backup or active MME/S-GW infrastructure in a centralized or distributed manner. The recovery mechanisms relied on the replication of critical context information of the subscribers in order to provide service resilience by minimizing the disruption time. Due to the limitation of the simulation model, the proposed approach addresses only the ACTIVE mode terminals or UEs which have radio control connection to the network.

A study was conducted to investigate the state of the art disaster recovery solutions and the off-site data replication methods as remote data protection is an important requirement of the proposed solution. For the performance evaluation, the proposed approach was modeled in OPNET simulation environment [32]. Simulation experiments were focused on the centralized approach. MATLAB [33] was used to measure the bandwidth requirements and signaling message overhead induced by the self-healing design.

1.7 Thesis Contribution

- The limitation of current EPC network element restoration methods is that it does not ensure service continuity. In fact, described solutions work based on the assumption that the MME or the S-GW restarts after failure. The service disruptions during the restart time or the service outage for prolonged period of time due to unavailability of the affected MME/S-GW have not been addressed yet. The self-healing solution proposed here maintains service continuity in the case of the MME/S-GW failures by proactively restoring the affected services.

- Various aspects of the self-healing mechanisms are discussed in details and solutions are designed to address these issues from LTE network and protocol perspective. The effectiveness of the proposed approach is demonstrated by simulation.

1.8 Thesis Outline

The structure of this thesis is organized as below:

The first chapter gives a brief introduction and background to the problem, motivation for this research, literature survey and related works, research objectives, approaches and contribution.

The second chapter provides an overview of the LTE core network architecture and protocols, and some important concepts that are used throughout the thesis.

The third chapter describes various aspects of the self-healing proposals for the LTE core network.

The fourth chapter discusses the simulation model and provides performance evaluation of different configurations and scenarios of the self-healing system.

The fifth chapter analyzes the bandwidth requirements induced by the proposed design.

Finally, the sixth chapter concludes the thesis with some directions for the future research.

Chapter 2

LTE EPC Architecture and Protocols

The objective of this chapter is to present an overview of the EPS architecture, protocols and some related concepts. A detail description of the network elements and interfaces, especially the core network or EPC, are described along with the protocol stacks. We have outlined some other important key concepts, with particular focuses on the end-to-end EPS bearers and QoS, user context information, Self-organizing Network (SON) that are relevant to this thesis.

2.1 The EPS Architecture

In 2005, 3GPP started investigation of next generation network under two study programs [34]. The LTE program focused on the radio network and air interface evolution known as Evolved UMTS Terrestrial Radio Access Network (E-UTRAN). The Service Architecture Evolution (SAE) program started investigation on the evolution of packet core network, commonly known as the Evolved Packet Core (EPC). Though LTE is the dominant term in the literature but in 3GPP specifications, the term Evolved Packet System (EPS) is used instead of LTE to refer to the system which consists of E-UTRAN, EPC and the terminals [35]. The terms LTE and EPS have been used interchangeably throughout this thesis.

The higher level objectives of the EPC are [36-37]:

- The EPC should provide a converged packet core that supports various access technologies and ensures seamless mobility based on operator implementation strategy, subscriber's choice etc. For example, both 3GPP access (LTE-EUTRAN, UMTS-UTRAN, GPRS-GERAN) and non-3GPP access technologies (WiMAX, CDMA2000,

WLAN) are supported by the EPC. This feature is not supported by other 3GPP technologies.

- Simplified flat architecture with improved performance in terms of higher data rates, lower latency and reduced costs for the operators.
- End-to-end IP based connectivity as well as maintaining negotiated QoS across the whole system.
- The EPC should ensure efficient support for different kinds of services in the PS (packet switched) network.

Figure 2-1 illustrates overall architecture of the EPS along with the E-UTRAN, EPC, interfaces and interfaces to other access networks (e.g., UTRAN, GERAN). A detailed description of the network elements and the interfaces is given in subsequent sections.

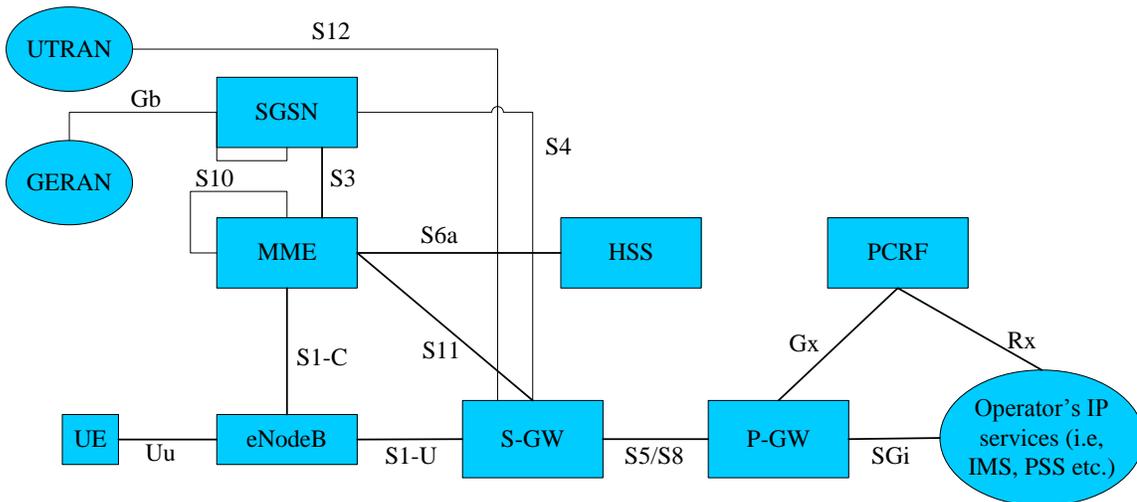


Figure 2-1: Evolved Packet System (EPS) architecture.

2.1.1 The Evolved UMTS Terrestrial Radio Access Network (E-UTRAN)

The LTE radio access network consists of eNodeBs. An eNodeB is a radio base station that is in control of all radio related functions. It inherits functionalities of the 3G NodeB. In addition, most

of the functionalities or protocols implemented in the 3G Radio Network Controller (RNC) are transferred to the eNodeB. Benefits of the RNC and Node-B merger include reduced latency with fewer hops in control and data plane. The eNodeB is also responsible for header compression, ciphering and reliable delivery of packets. On the control plane, functions such as admission control and radio resource management (RRM) are also incorporated into the eNodeB. A new interface, called X2 has been defined to interconnect eNodeBs. The main purpose of this interface is to facilitate seamless user mobility by reducing data loss during handover. Figure 2-2 shows E-UTRAN architecture.

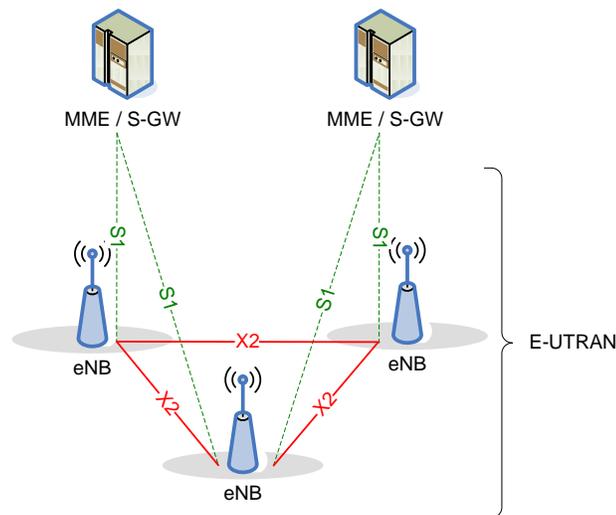


Figure 2-2: LTE E-UTRAN architecture [2].

2.1.2 The Evolved Packet Core (EPC) Network

One of the most important features of the EPC is that the circuit-switched based voice service and packet-switched based data service of earlier 2G and 3G networks are unified under one single IP mobile domain [36] which is an evolution of the packet-switched domain of earlier GPRS/UMTS network. Furthermore, the control plane and data plane have separate interfaces and network

entities in the EPC. This feature provides flexibility to the carriers or network operators to optimize their signaling and data traffic capacity independently.

The LTE-EPC consists of the following primary functional elements [36, 38-39]:

- Mobility Management Entity (MME)
- Serving Gateway (S-GW)
- Packet Data Network (PDN) Gateway

The EPC also includes other nodes such as the Home Subscriber Server (HSS) and Policy Control and Charging Rules Functions (PCRF). The HSS contains user subscription information and acts as authorization, authentication and accounting (AAA) server and executes QoS and charging rules via the PCRF.

Mobility Management Entity (MME): In the control plane of the LTE EPC, main network element is the MME. It handles functions related to signaling, resource assignment, user mobility and session management. There is a direct control plane logical connection between the MME and UE. The MME sends paging messages to the UE when the UE is in ECM (EPS Connection Management) idle state. It also manages tracking area (TA) list and provides handover support. In order to ensure authentication of the subscriber, the MME communicates with the HSS and initiate security and ciphering/integrity protection procedures. It maintains the privacy by assigning a temporary identity called GUTI (Globally Unique Temporary ID) to the UE instead of using permanent UE identity or IMSI [40]. It is the responsibility of the MME to select appropriate S-GW and P-GW for a UE in initial network attachment. The MME performs bearer establishment and release, and coordinates issues related to the bearer activation/deactivation, negotiation of bearer QoS etc. It provides control functionalities required for the mobility between the LTE and other 2G/3G networks (e.g., GSM, UMTS).

Serving-Gateway (S-GW): S-GW is a data plane network element of LTE EPC. The main function of the S-GW is to relay user data to and from the eNodeB and P-GW. It allocates resources as per the requirements of the MME and P-GW. It works as mobility anchor for the user plane during inter-eNodeB handover as well as in the case of mobility between the LTE and other 3GPP technologies. The S-GW buffers subscriber data when the UE is in ECM-IDLE state. It transfers buffered data to the UE when it connects to the network or in ECM connected state after successful paging.

Packet Data Network (PDN) Gateway (P-GW): P-GW is the exit and entry point between the EPC and the external packet data network. Typically, it is responsible for allocating IP address to the UE. The P-GW establishes and maintains GPRS Tunneling Protocol (GTP) tunnels to the S-GW. It can also establish and delete tunnels to the SGSN for inter-RAT mobility scenario. The P-GW can enforce policy for resource allocation, usage and packet filtering etc. In order to do policy control it can be connected to a Policy and Charging Rule (PCRF) via the Gx interface. P-GW also acts as a mobility anchor between the 3GPP and other non-3GPPP networks such as WiMAX, CDMA etc.

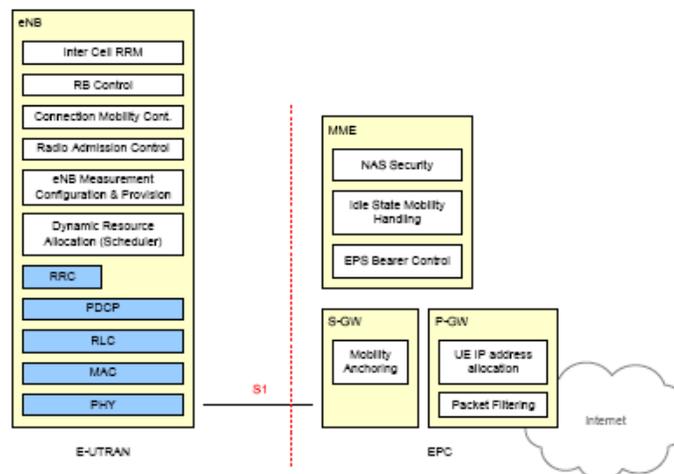


Figure 2-3: Functional split between the E-UTRAN and the EPC [2].

2.1.3 LTE Interfaces

LTE interfaces are often referred to as the reference points. LTE introduces some new interfaces and some interfaces are incorporated from existing UMTS/HSDPA network. Figure 2-1 shows interfaces and corresponding network elements. A description of important interfaces of the access and core network relevant to this thesis is given below [4, 38]:

S1: S1 interface is used to connect the eNodeB to the EPC network elements. An IP based transport stack is implemented on the S1. This interface is divided into two groups - the S1-U (user or data plane) interface and S1-C (control plane) or S1-MME interface. The purpose of the S1-U interface is to transport user data between the eNodeB and S-GW. This interface uses the GTP over UDP/IP. On the contrary, S1 control interface (S1-C) is a signalling interface which is used to transport control signals between the eNodeB and MME.

X2: The X2 interface connects an eNodeB to other eNodeBs. Like the S1 interface, it is divided into X2-U (X2 user plane interface) and X2-C (X2 control plane interface). The X2-U interface is used to transport user data between eNodeBs during handover. It uses the same GTP tunneling protocol as used on the S1-U interface. The X2-C is a signaling interface which supports a set of functions and procedures between the eNodeBs.

S10: The MME connects to other MMEs over S10 interface. It is used for MME relocation scenarios. The S10 is considered as an extension of S1 interface.

S11: The MME and S-GW communicates over the S11 interface. GPRS Tunneling Protocol for the control plane (GTP-C) is used at this reference point.

S5: It is the reference point between the S-GW and P-GW for transferring both user plane and control plane data. GTP-C and GTP-U are used as control plane and user plane protocol, respectively.

S8: It is also used as reference point between the S-GW and the P-GW. The difference between the S5 and S8 interface is that the S8 interface is used for only roaming subscribers. It is used for transferring data between the S-GW and P-GW where S-GW belongs to a Visited PLMN (VPLMN) and P-GW located in a Home PLMN (HPLMN). The GTP-U and GTP-C protocols run over this interface.

S6a: The MME communicates with the HSS over the S6a interface in order to transfer subscriber's authentication and authorization data. The DIAMETER protocol is used in this reference point.

Gx: The QoS policy and charging related rules are transferred from the PCRF to the P-GW over this interface. At the time of bearer establishment, the PCRF provides the P-GW with required information for charging service data.

2.2 LTE Network Protocol Architecture

Logically, LTE network protocols can be divided into control plane and user plane [38]. The main purpose of the control plane is to set up the user plane. It establishes, modifies, and releases the user plane. The user plane transfers user data between users or application servers. Figure 2-4 and Figure 2-5 show the end-to-end control plane and user plane protocol stacks in the EPS.

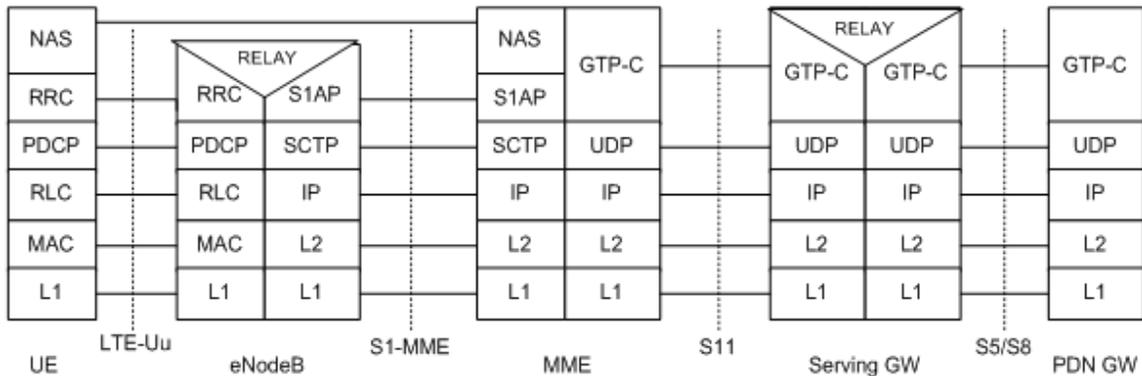


Figure 2-4: Control plane protocols stack.

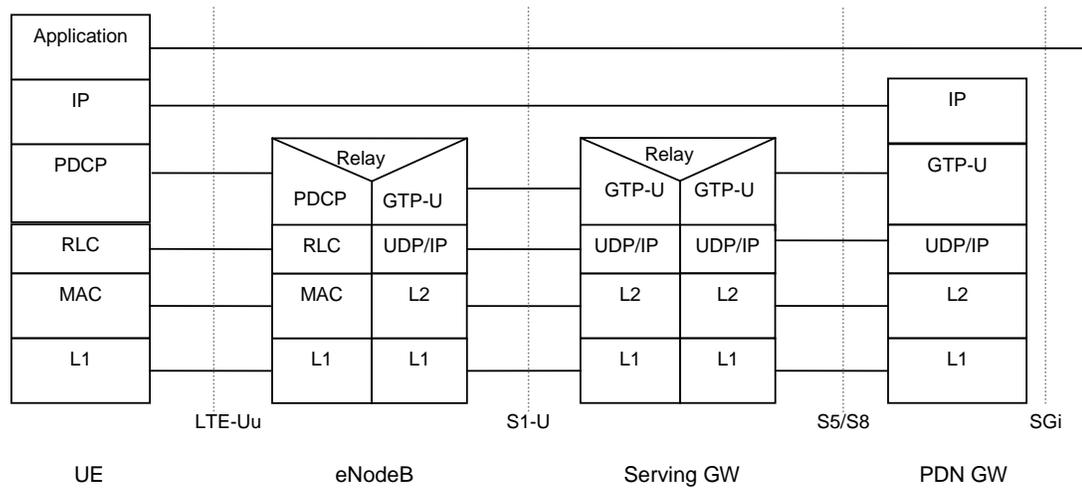


Figure 2-5: User plane protocols stack [38].

2.2.1 Description of Protocols Used in Control Plane and User Plane

Though focus of this thesis is on the protocols which are related to the core network or the EPC, a brief description of the radio protocol architecture is included in the discussion.

Non-access Stratum (NAS): Non-access Stratum or NAS is the top most control plane layer which spans over the access network and the core network. Signaling messages are directly transported from the UE to the MME via NAS. NAS layer communications between the UE and the MME are transparent to the eNodeB, it just forwards the NAS messages to the MME. The NAS layer protocols are divided into two parts [41]:

- EPS Mobility Management (EMM):** The EMM protocols manage issues regarding UE mobility within the E-UTRAN. It is responsible for the functions related to the connection management such as UE attachment or detachment to the network. Other functions of EMM layer include handling the UE initiated service request or network initiated paging procedures; authentication and the UE privacy, ensuring NAS layer security.

- **EPS Session Management (ESM)** The ESM protocol is responsible for handling bearer context between the UE and MME. It works along with the E-UTRAN bearer management functions and supports UE initiated bearer procedures such as creation, modification and release of EPS bearers with specific QoS.

On the radio interface, the control plane uses the PDCP, RLC, MAC and PHY stack to transport both RRC (Radio Resource Control) and core network NAS signaling. The PDCP, RLC, MAC and PHY layers support the same functions for both the user and control planes.

Radio Resource Control (RRC): The RRC layer provides the E-UTRAN Radio signaling connections to the upper layers in the protocol stack, for example, Radio Resource Management is used to facilitate the transferring of the upper layer's message flow [42]. This signaling connection is used between the UE and core network. The ACTIVE mode UE maintains the RRC connection with the eNodeB but IDLE modes UEs do not have a RRC connection.

Packet Data Convergence Protocol (PDCP): The PDCP provides data transfer, header compression as well as ciphering services to both control plane (RRC) and user-plane (application) entities.

Radio Link Control (RLC): The RLC layer is responsible for segmentation and reassembly as well as the error correction procedures.

Medium Access Control (MAC): In the eNodeB, The MAC layer performs scheduling which distributes the available bandwidth to a number of active UEs. It also implements HARQ operation for the retransmission of data.

The S1 and X2 interfaces are used to connect the E-UTRAN to the EPC and E-UTRAN to another E-UTRAN, respectively. Following protocols are involved:

S1-AP: The S1-AP (S1 Application Protocol) protocol is used on the control plane to transfer signaling message between the E-UTRAN and EPC. Figure 2-6 illustrates protocol stack over the S1-C interface including the S1-AP. Following are the typical functionalities of the S1-AP [43]:

- NAS signalling message transferring between the MME and the UE.
- EPS bearer setup, modify and release.
- S1 based handover management functions such as handover preparation, resource allocation and notification, path switch request.
- Paging and location reporting activities related to the UE.
- Network management functionalities, for example, S1 setup, eNodeB and MME configuration update, overload etc.

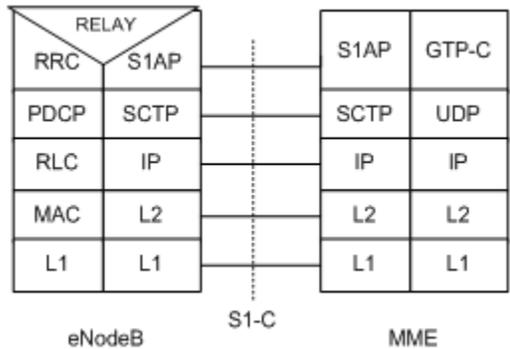


Figure 2-6: Protocol architecture between the eNodeB and the MME.

X2AP: The X2AP is primarily responsible for the X2-based handover. Figure 2-7 shows protocol stack over the X2 interface including X2AP. Following are the main functionalities of this protocol [44]:

- X2AP facilitates UE mobility between E-UTRANs.

- Using the management functionalities, adjacent eNodeBs exchange resource status which provides information regarding network configuration, traffic load etc. In this way, an overloaded eNodeB can move UEs to a lightly loaded eNodeBs.
- X2 Setup procedures allow exchanging necessary data for the eNodeB to setup the X2 interface.

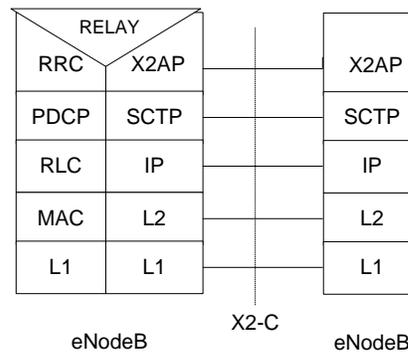


Figure 2-7: Inter-connected eNodeB protocol architecture.

Stream Control Transmission Protocol (SCTP): In order to ensure high reliability in terms of message re-transmission and delay, LTE S1-C uses SCTP as a transport protocol between the E-UTRAN and EPC. SCTP is a connection-oriented protocol similar to the TCP. It provides flow and congestion, re-transmission mechanism, detection and data corruption etc., as TCP. In addition, it provides two important features, multi-homing and multi-streaming which TCP does not support [45].

GPRS Tunneling Protocol-Control Plane (GTP-C): This protocol is used on the core network part of the network. It is an IP-based protocol, which carries traffic in GPRS and UMTS networks. The GTP is also used in LTE. In telecommunication, tunneling represents a two-way point-to-point communication path established between two entities [38]. GTP-C implements necessary procedures to create, maintain and delete tunnels. Another important task of the GTP-C

is to forward relocation information of the moving users. Figure 2-8 shows the GTP-C in the protocol stack over S11 and S5/S8 interfaces.

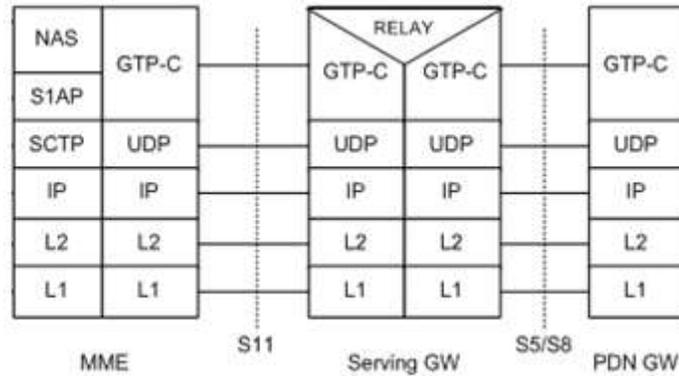


Figure 2-8: GTP-C over S11 and S5/S8 interfaces.

GPRS Tunneling Protocol-User Plane (GTP-U): The GTP-U is used to carry user data between the E-UTRAN and EPC as well as within the E-UTRAN and core network. The protocol runs over S1-U, S10, X2-U, S5 and S8 interfaces. Figure 2-9 shows GTP-U protocol stack over S1-U and S5/S8 interfaces.

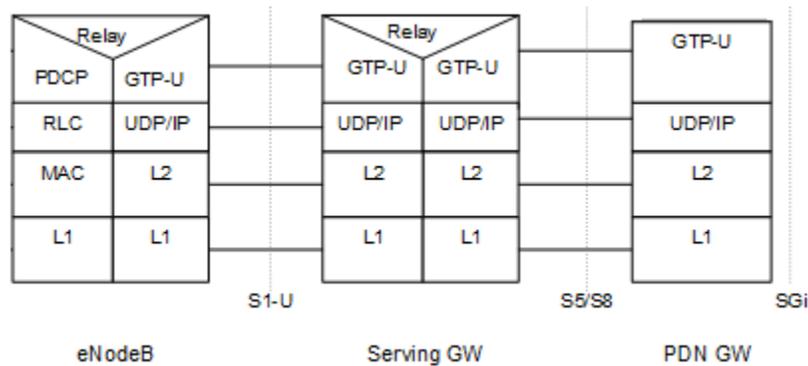


Figure 2-9: GTP-U over S1-U and S5/S8 interfaces.

2.3 LTE EPS Bearer and QoS

In LTE, an EPS bearer is established through the EPS network (core network and E-UTRAN access network) to transport user data to and from the UE and PDN. It has a defined set of data transmission characteristics such as, specific QoS and flow control etc. A default EPS bearer with default QoS is established when a UE connects to a PDN. The UE maintains this connection throughout the lifetime of the PDN connection which provides the user with always-on IP connectivity to that PDN. Furthermore, additional EPS bearers can be established in the course of time in response to specific service requests which are called dedicated bearers. In fact, one default and several dedicated bearers can be established within one connection at the same time.

An EPS bearer between a UE and PDN, has three segments [2]:

- Radio bearer between the UE and eNodeB.
- Data bearer between the eNodeB and S-GW (S1 bearer).
- Data bearer between the S-GW and P-GW (S5/S8 bearer).

Figure 2-10 shows parts of the EPS bearer and corresponding interfaces as well as network elements.

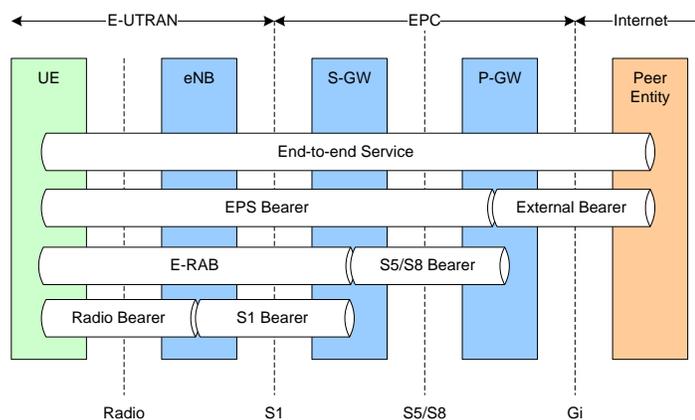


Figure 2-10: LTE EPS bearer architecture [2].

The end-to-end QoS management is very important in LTE in order to provide bandwidth intensive multi-media rich services. A UE can run different applications at the same time that may need specific QoS requirements based on the nature of service. For instance, multi-media streaming session requires more stringent QoS requirements in terms of delay than an http session. In order to support such scenarios, multiple bearers can be set up assigning different QoS. Furthermore, data plane traffics are transferred over a bunch of virtual connections within these bearers which are called SDFs or Service Data Flows. As illustrated in Figure 2-11, one bearer is comprised of such several SDFs in order to transport different data that requires same QoS. For example, a VOIP session initiated by a UE can be mapped to a SDF of a particular bearer. Similarly, a streaming session from a server can be mapped to another SDF of the same bearer. The IP packets directed to a particular bearer receive same QoS treatment in the EUTRAN and EPC in terms of scheduling, queue management etc. In order to map the IP packets to a particular bearer a filtering mechanism should be implemented in both sides (uplink and downlink). Each of the SDFs is defined by 5-tuple (source IP address, destination IP address, source port, destination port, protocol used above IP) [46]. The appropriate IP packets of corresponding SDFs are mapped and delivered based on this definition.

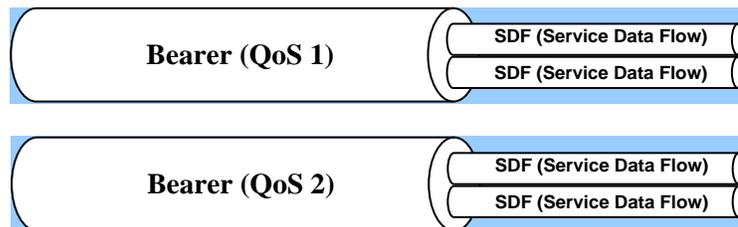


Figure 2-11: LTE EPS bearers and corresponding service data flows (SDFs).

Again, each bearer is associated with the **TFT** (Traffic Flow Template) which defines a set of packet filters. The UE and the PDN filters the packet into different SDFs in a bearer based on the

TFTs. The IP header information such as source and destination IP addresses and TCP ports are used by the TFTs in order to filter packets to the bearers with appropriate QoS. Two kinds of TFTs are used – UL TFT for uplink direction and DL TFT for downlink direction. Both of these TFTs associated with bearers in the UE and P-GW send packets to the respective directions. Figure 2-12 illustrates the relationship between a TFT and a bearer. The Tunnel End Point Identifier (TEID) is used to identify a S1 and S5 bearer in the EPC part which corresponds with a Radio Bearer ID (RBID) in the access network.

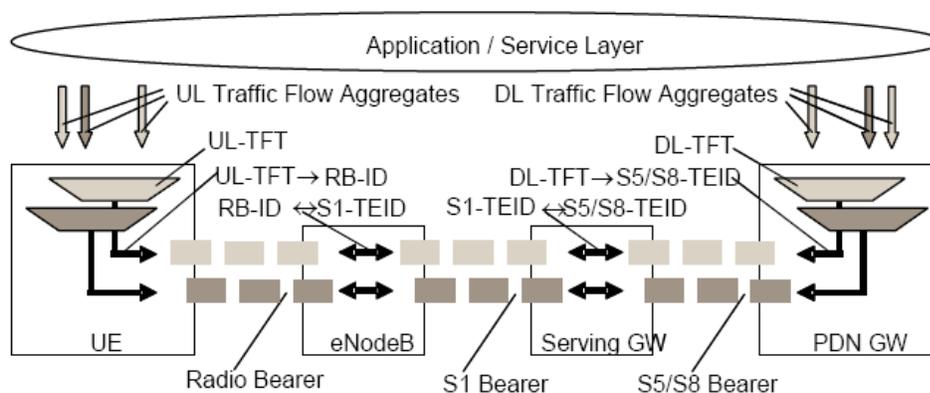


Figure 2-12: Two unicast EPS bearers (GTP based S5/S8) [38].

Based on the QoS, bearers can be divided in two categories [36]:

- Minimum guaranteed bit rate (GBR) bearers that can be used for specific QoS oriented applications. The network can provide higher bit rate for a bearer in condition of resource availability. In that case, a maximum bit rate or MBR can be assigned for a GBR bearer. User can get higher bit rate up to the limit of MBR from that particular GBR bearer.
- Non-guaranteed bit rate (GBR) bearer does not provide any specific bit rate. Less QoS oriented applications such as, http session or file transfer can be used in these bearers. Generally, dedicated bearer is Non-GBR.

The eNodeB provides required QoS for a bearer in the radio interface. Each bearer has two parameters, QCI (QoS Class Identifier) and ARP (Allocation and Retention Priority). As indicated in Table 2-1, QCI has specific characteristics in terms of priority, packet error loss rate, delay budget. The eNodeB treats a bearer based on the QCI label associated with it. So far a dozen QCI have been standardized which facilitates a common understanding among the vendors about the requirements of the services.

Table 2-1: LTE standardized QCI [55].

QCI	Resource Type	Priority	Packet Delay Budget	Packet Error Loss Rate	Example Services
1		2	100 ms	10^{-2}	Conversational voice
2		4	150 ms	10^{-3}	Conversational video (live streaming)
3	Guaranteed Bit Rate (GBR)	3	50 ms	10^{-3}	Real-time gaming
4		5	300 ms	10^{-6}	Non-conversational video (buffered streaming)
5		1	100 ms	10^{-6}	IMS signalling
6		6	300 ms	10^{-6}	Video (buffered streaming) TCP-based (e.g., www, e-mail, chat, ftp, p2p file sharing, progressive video, etc.)
7	Non-GBR	7	100 ms	10^{-3}	Voice, video (live streaming), interactive gaming
8		8	300 ms	10^{-6}	“Premium bearer” for video (buffered streaming), TCP-based (e.g., www, e-mail, chat, ftp, p2p file sharing, progressive video, etc) for premium subscribers
9		9	300 ms	10^{-6}	“Default bearer” for video, TCP-based services, etc. for non-privileged subscribers

2.4 UE Context Information

When a UE connects to the network, the MME creates a UE context. UE context contains subscription information of user that is retrieved from the Home Subscriber Server (HSS). This process ensures that MME does not need to communicate with the HSS every time when there is a request for new bearer establishment or other similar queries. This facilitates faster execution of

procedures. The MME assigns a short temporary identity which is called SAE Temporary Mobile Subscriber Identity (S-TMSI) to the UE in order to identify the context. In addition, UE context contains information elements (IEs) such as International Mobile Subscriber Identity (IMSI), UE security contexts, UE network capability, selected core network operator id, EPS bearer context(s) which in turn includes S-GW and P-GW addresses and TEIDs of the EPS bearers, access point name (APN) etc. [38]. The S-GW, eNodeB and P-GW hold UE context data and EPS bearer context table. Detail description of each of the information elements (IEs) are provided in chapter 5. UE contexts are important piece of information in proposed self-healing mechanism.

2.5 S1-Flexibility (S1-Flex)

In traditional 2G and 3G networks, core network and access network connections have one-to-multi hierarchical relationship. In this concept, core network (CN) elements like MSC or SGSN handles a group of radio controllers (BSC or RNC) which are responsible for a set of base stations (BTS or NodeB) [45]. In these hierarchical architectures, access network elements can be connected to only one core network element which means a core network elements, for example, MSC has its own set of BSC and BTS. Thus, these access nodes, BTS or BSC, cannot be connected to another MSC at the same time.

3GPP release 5, first introduced the concept where multiple core network nodes can share one common geographical areas or access network. LTE standards, adopted this feature where the S1 interface is used to connect one eNodeB to multiple MME/S-GWs through a mesh network. This is known as S1-flexibility or S1-flex. Figure 2-13 illustrates a typical LTE S1-flex configuration. The MME/S-GWs that serves a common geographical area are known as MME/S-GW pool. There are three important advantages of this feature. First, it provides network redundancy.

Second, the core network elements may execute load sharing of traffic in time of overload situation. Third, multiple operators can share same radio access network.

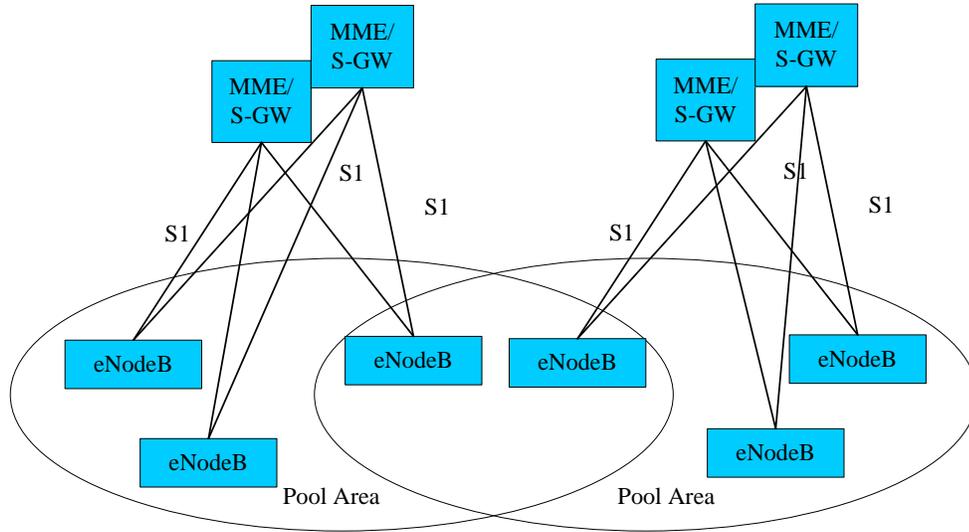


Figure 2-13: LTE S1-flex configuration.

2.6 LTE Self-Organizing Network (SON)

In current GSM/UMTS wireless network, network elements and most of the corresponding parameters are manually configured. In LTE, the number of access network NEs (eNodeB and Home eNodeB known as Femtocell) and tuneable network parameters become large and complex due to its flatten architecture. Furthermore, growth of diversified services with stringent QoS requirements and parallel support of multiple radio access networks have increased the complexity in terms of resource utilization and performances. As a result, maintenance and deployments of such network become labor-intensive and expensive for the network operators. In the above context, operator's forum and international standardization bodies have paid particular attention to the research of self-organizing network (SON) in order to automate the planning, configuration and management of the network. In fact, SON is the key driver for improving the

LTE OAM [47]. Generally, two important benefits can be achieved by reducing the manual intervention from the network management tasks:

- **Reduction of CAPEX and OPEX:** Study [47] shows that 17% of operators CAPEX goes to the engineering and installation service purpose. In addition, 24% of the revenue is spent for the OPEX or network operation and maintenance purpose. Automation can significantly reduce operator's expenditure and enhance resource utilization.
- **Fast response to network change:** Manual intervention is inefficient and results in long delays to adjust in fast-changing network conditions. Thus, automation can significantly improve end users experience, network quality and availability.

LTE SON development was initially promoted by the industry forum NGMN (Next Generation Mobile Networks) [48]. In 2006, NGMN established a set of preliminary requirements for SON and later defined some use cases covering various aspects of the network including planning, optimization and maintenance. 3GPP adopted these concepts in its LTE SON standardization process, which started in release 8. Table 2-2 shows some NGMN use cases.

2.6.1 Main functionalities of SON in LTE

The main aspects of the LTE SON are: Self-Configuration, Self-Optimization and Self-Healing. A short description of each of this feature is given below.

Self-configuration

The self-configuration mechanism enables a network element to configure itself to perform necessary task of the network. From LTE network perspective, self-configuration feature is desirable during deployment and subsequent expansion phase where NEs like eNodeBs are configured by automatic installation procedures [49]. The self-configuration aims to reduce the CAPEX expenditure. Though pre-planned configuration is allowed but it is desirable that eNodeB

should have the capability to automatically discover and configure neighbour lists, physical cell id and other RF parameters as mentioned in the use cases described in [19].

Table 2-2: NGMN SON use case definitions [48]

Planning	Optimization
Planning of eNodeB	Support of centralized optimization entity
Planning of eNodeB Radio parameters	Neighbor list optimization
Planning of eNodeB Transport parameters	Interference control
Planning of eNodeB data alignment	Handover parameter optimization
	QoS parameter optimization
	Load Balancing
	Home eNodeB optimization
	RACH load optimization
Deployment	Maintenance
Hardware installation	Hardware/capacity extension
eNodeB/network authentication	Automated NEM upgrade
O&M Secure tunnel setup	Cell/Service outage detection and compensation
Automatic inventory	Real-Time Performance management
Automatic Software download to eNB	Information correlation for fault management
Transmission setup	Subscriber and equipment trace
Radio parameter setup	Outage compensation for higher level network elements
Self Test	Fast recovery of unstable NEM system
	Mitigation of outage of units

Self-Optimization

The objective of self-optimization is to auto-tune network parameters dynamically in response to rapid network conditions and traffic change. It takes account of live measurement data, and based on the analysis continuously selects and adjusts various parameters of the network in order to achieve optimal system performances. The self-optimization is important for operational state of the network as currently lots of human efforts are required for optimization activities. The main optimization use cases that are discussed in [50] are mainly related to the coverage and capacity, interference, load balancing, energy savings, mobility robustness.

Self-Healing

The purpose of self-healing is to automatically detect and localize failures, and applies appropriate recovery mechanisms. The self-healing use cases discussed in [20] are mainly related to cell outage detection and compensation, recovery from software and hardware failures of the eNodeB.

2.6.2 LTE SON Architecture

SON architecture can be divided into three classes based on the location of SON procedures: *i)* centralized, *ii)* distributed and *iii)* hybrid SON [51]. A brief description of each of this architecture is given below:

Centralized SON

In centralized SON architecture, SON algorithms are executed at the network management level or in the OAM. In this case, small number of higher level network elements contains SON functionalities. The advantage of this approach is that it involves less complexity in deployment and management tasks. But it has several drawbacks. First, in centralized SON architecture, SON algorithms and procedures reside at the network management level. So, SON decisions are made at smaller numbers of higher level network elements in the OAM. Failures of central network elements can be critical on overall performances of the SON. Second, different vendors have their own OAM system. And co-ordination and optimization supports among the vendors are low. There is an interface called Itf-N to facilitate multi-vendor management which should be standardized and extended to increase the support between different vendors. Finally, huge volume of data processing and computational complexity can be a bottleneck issue. Figure 2-14 shows centralized SON architecture.

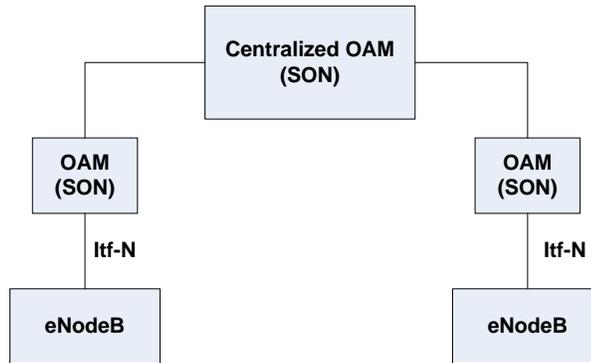


Figure 2-14: LTE centralized SON architecture [51].

Distributed SON

In distributed SON architecture, SON decision making takes place at low levels of the network, for example, access network elements or eNodeBs. Different eNodeBs can share information over the standardized X2 interface. This approach is fast and flexible in comparison to the centralized SON and reduces the complexity of multi-vendor SON support. But deployment cost is high as many network elements are involved in this process. Though optimization task that requires involvement of smaller numbers of eNodeB yields better result in this case, complex optimization schemes that require co-ordination among many eNodeBs are difficult in distributed SON. Figure 2-15 shows distributed SON architecture.

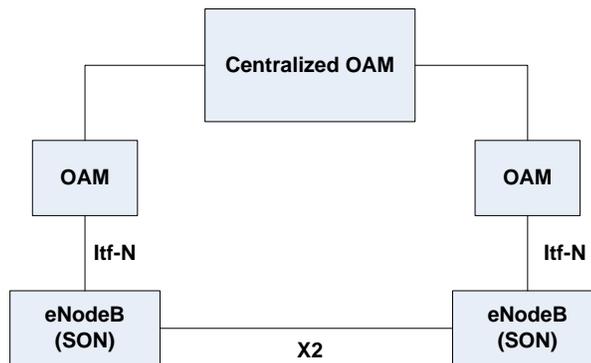


Figure 2-15: LTE distributed SON architecture [51].

Hybrid SON

In hybrid SON architecture, some of the SON algorithms are executed in the access networks or eNodeBs and other functionalities are executed in the management level or the OAM. It is flexible and has the advantages of both centralized and distributed approaches where simple and quick optimization tasks take place in the eNodeBs and complex optimization tasks are implemented in the OAM. The X2 interface can be used to support optimization across different vendors. But the disadvantage is that the deployment and implementation costs, as well as interface extension efforts are high. Figure 2-16 illustrates the hybrid SON architecture.

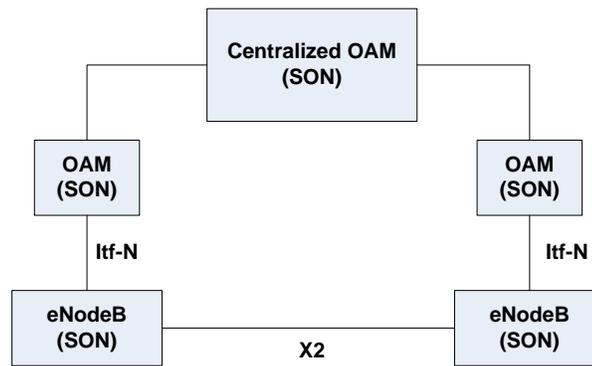


Figure 2-16: LTE hybrid SON architecture [51].

Chapter 3

LTE Self-healing Core Proposal

This chapter provides a detailed description of various proposals for implementing a self-healing system in the LTE core. Service restoration mechanisms are presented in the case of a combined MME and S-GW failure which are implemented in the same platform. The MME/S-GW failure affects both the control and data plane of the core network and may impact large number of UEs. The main concept of the proposed approach is to ensure service continuity by triggering the MME/S-GW relocation and recovering the failed states of the affected UEs by re-creating corresponding bearers. Details of message flows between different network elements are presented to implement these procedures. Other aspects of a typical self-healing system are utilized and explained from LTE network architecture and protocols perspective.

3.1 Self-healing System Overview

In order to design the self-healing system, the centralized and distributed SON architecture concepts described in chapter 2 have been adopted in this thesis. As mentioned, the SON procedures are executed in the management level in the centralized architecture. But in the distributed architecture, SON algorithms are executed in the lower level nodes. In the centralized approach, modules of the self-healing system discussed in this chapter are implemented in the OAM. On the other hand, distributed self-healing system functionalities are implemented in the MME/S-GW. The pros and cons of each of the architecture are presented in later sections.

Figure 3-1 and 3-2 show critical components or modules of the self-healing system in a centralized and distributed approach. It includes:

- Fault-tolerant system architecture.

- Configuration and status update.
- Failure detection, notification and fault isolation.
- Failover coordination and decision making process.
- Service recovery mechanism.

Crucial data replication and maintaining component redundancy are some of the mechanisms that are commonly used in the design of fault-tolerant systems [21]. In such designs when one component fails or malfunctions, it can be replaced by another component with equivalent functionalities. But complete redundancy of all equipments for tackling large scale failures requires expensive implementation and resources. Alternate cost-effective approaches can be implemented using the centralized or distributed redundancy strategy; for example, where one redundant backup system supports multiple active systems or a group of active systems provide redundant support to each other in time of need. Information redundancy or availability of critical data is essential to ensure service continuity. These principles have been adopted in the proposals. The backup system should be aware of configuration changes that occur in the active nodes in order to ensure proper and efficient recovery. The states of the active nodes must be monitored periodically for timely detection of failure and faulty device should be isolated to minimize the effects of faults. The purpose of the failure coordination and decision making process is to start appropriate recovery process based on the overall system status and operators chosen policy. Finally, the recovery mechanism triggers healing techniques. It is expected that the restoration process should be quick enough so that continuous availability of the service can be maintained. Collaboration of above mentioned tasks is necessary for the proper functioning of self-healing schemes. In order to facilitate co-ordination among the different modules, a recovery manager

module is introduced in the OAM, active and backup MME/S-GWs. Main self-healing modules are discussed in details throughout this chapter.

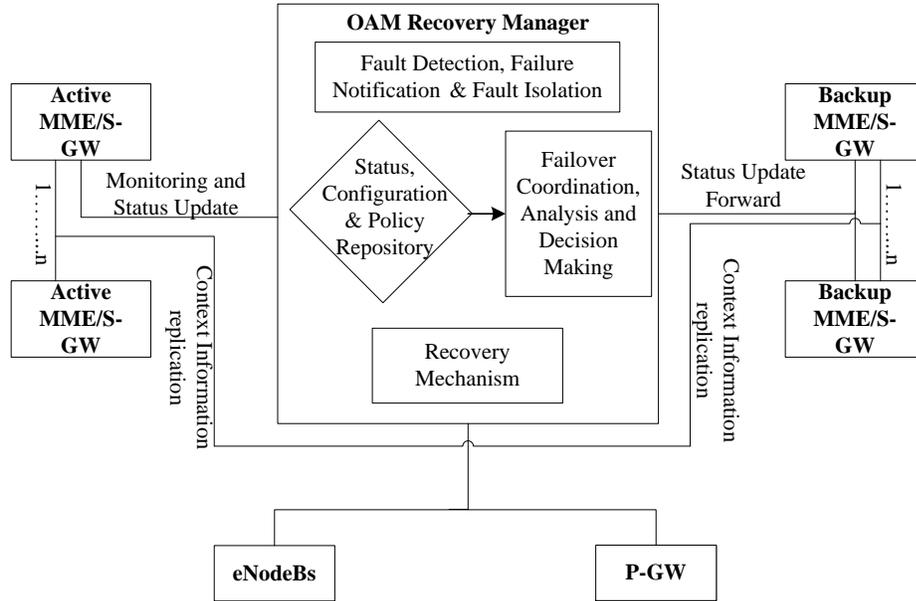


Figure 3-1: Logical architecture of the Self-healing system in the centralized approach.

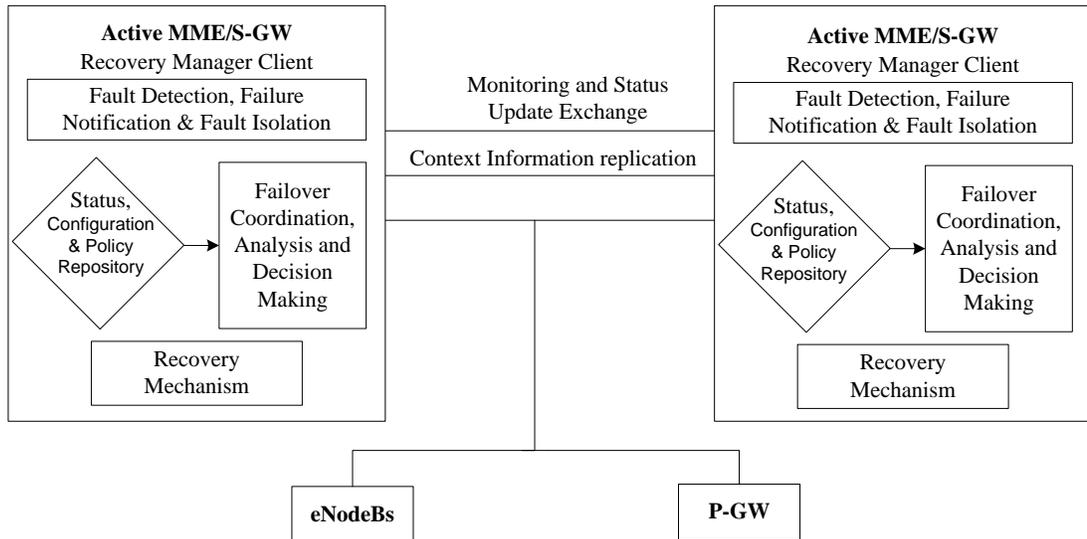


Figure 3-2: Logical architecture of the Self-healing system in the distributed approach.

3.1.1 Fault-tolerant System Architecture

In the self-healing solution, a combination of site and node redundancies must be employed to maintain service in the case of core component failures. Based on the discussion in previous section, the following self-healing designs are considered for achieving the required redundancy:

- N:M Active-Backup Configuration: M backup MME/S-GW nodes for N active MME/S-GW nodes. This configuration follows the centralized approach.
- 1:1 Active-Active Configuration: Two active MME/S-GW nodes acting as backup for each other. This configuration is an example of the distributed approach.

Note that the backup and active nodes in the above architectures are assumed to be located in different geographical areas so that if a major failure or natural disaster occurs in primary sites, backup sites remain unaffected and would be able to take over the whole or part of the functionalities of the primary site to provide the required service continuity. Some design requirements for implementing this strategy include high speed link connections for data transfer, protocol-level support of off-site data protection technology and failover and recovery mechanism coordination. Well-known high speed transport technologies such as optical fiber or Metro/Carrier Ethernet are capable to provide required bandwidth and latency over connected links. Furthermore, synchronous or asynchronous remote mirroring and shared storage network technology, for instance, Storage Area Network (SAN) extension solutions such as FCIP, iSCSI are now available to facilitate real time data replication over a long distance link in a WAN [52]. It should be noted that the backup MME/S-GW must contain active MME/S-GW's latest configuration information, for example, supported eNodeBs, Tracking Area (TA) List to ensure effective recovery of the services.

3.1.1.1 N:M Active-Backup Configuration

In the N:M active-backup configuration, a pool of inter-connected remote backup MME/S-GWs are available to support multiple active MME/S-GWs in the event of nodes failure. In general, one backup MME/S-GW supports some selected active MME/S-GWs. But in multiple MME/S-GWs failure scenarios, other backup MME/S-GWs can also participate in recovery process considering resource constraints and performance issues. Proper recovery co-ordination measures should be implemented in the OAM side to facilitate this support. This strategy improves the availability and scalability of the backup system. As the centralized SON architecture is used in this case, the OAM recovery manager is responsible for failure detection, monitoring and failover recovery tasks as depicted in Figure 3-3.

As the dedicated backup node is involved, this approach adds more costs in the self-healing design. However, it provides efficient support in terms of capacity and resource in the case of one or more MME/S-GW failures. In this configuration, when a UE is connected to the active MME/S-GWs, the UE contexts information is replicated from active nodes to the backup node. The backup node uses these UE contexts to re-create the UE bearers if corresponding active node or nodes fail. Other configuration and dynamic data of active nodes can be replicated as well to the remote site for the backup purpose to tackle disaster scenarios. Dedicated or cost effective shared WAN links can be used to connect the primary sites with the backup site. However, appropriate bandwidth provisioning must be done to prevent performance degradation. It should be noted that the backup node is an independent node and has its own configuration.

The number of MME/S-GWs that can be supported during a failure depends on the capacities and resources of the backup and active MME/S-GWs and corresponding backhaul links. Node's capacity and resource can be defined in terms of number of concurrent active bearers and the

CPU processing power, respectively. If capacities and available resources of the backup nodes (e.g., N:1 Active-Backup configuration) or backhaul are not greater than the combined capacities and resources of the active nodes, it can only support a fraction of the combined capacities when multiple nodes become unavailable. Typically, in this approach, the backup node is dedicated to support only the high priority traffic on one or more nodes. As the MME/S-GW nodes are generally setup in geographically separated sites, simultaneous multiple node failures would have to be a result of a rare event such as a massive natural disaster.

Two issues should be resolved in multiple failure scenarios where the backup node resources are limited (i.e., N:1 Active-Backup configuration). First, a load balancing or load re-distribution strategy must be devised for the scenario when more than one node fails. One approach is to have an OAM node to keep track of the load status of each of the active nodes periodically. When a failure occurs, the OAM instructs the backup MME/S-GW to take over the subscribers of a particular node or nodes if capacity requirements are satisfied. Second, there can be a priority policy of supporting one or more preferred MME/S-GW nodes over others if the capacity of the backup node is limited.

Involvement of other backup MME/S-GWs in multiple active nodes failures may ensure effective supports. But it may increase restoration time as the OAM needs to select additional backup MME/S-GWs and transfer UE context information to those nodes. Besides, due to the centralized nature failover procedures require more message exchange between the active and backup nodes, and the OAM which in turns increases the service restoration time. Moreover, this approach suffers from single point of failure considering the OAM failures. But it should be noted that the reliability and availability of the OAM is much higher than other ordinary nodes.

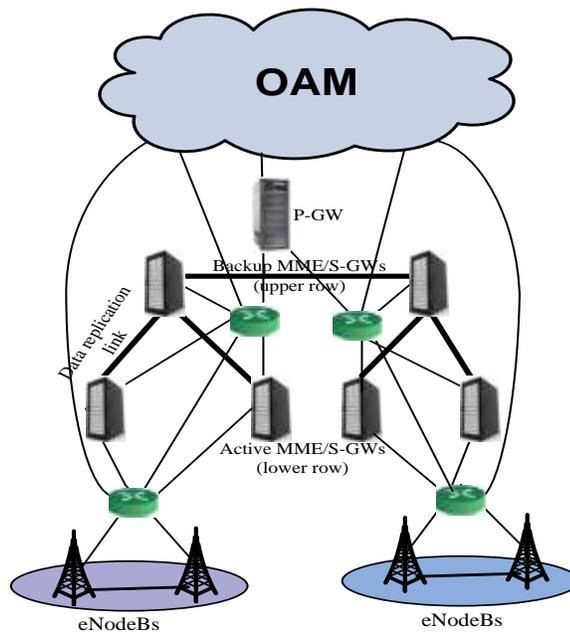


Figure 3-3: N:M Active-Backup configuration.

3.1.1.2 1:1 Active-Active Configuration

In this approach two active MME/S-GWs work independently but at the same time can act as each other backup in time of failure. The distributed SON architecture is adopted in this case where the active MME/S-GWs has recovery manager client which gathers configuration and status information of the peer nodes. Each MME/S-GW is responsible for monitoring and failure detection of other MME/S-GW, and dissemination of failure notifications to the access and core network elements. Figure 3-4 shows 1:1 Active-Active configuration.

As both active and backup nodes are in operation, each node must set aside some spare resources to handle the traffic of its counterpart in the case of failure. Three strategies can be adopted to resolve the problem. First, subscriber provisioning can be done in both nodes such a way that each node reserves sufficient resources for the other. The obvious disadvantage of this strategy is that the spare resources of a node will be idle until a failure occurs. Second, the recovery manager

client can keep track of the load status of the backhaul links and MME/S-GWs. At the time of failure, once it determines that the backup node has sufficient resources to support the failed node, only then it initiates the failover procedures. Third, based on the replicated UE contexts of other active node and its own load status, the backup node can forecast how many of these UEs it can support without affecting its normal operation. A simplified implementation method of this strategy is that when the active node fails, the backup MME/S-GW will start allocating resources based on the stored UE contexts and continuously keeps track of used and available resources. It will stop allocating resources when the utilization reaches a threshold limit. In this way, the backup MME/S-GW can support as many UEs as possible avoiding any risk. To accommodate more UEs, the backup node may create default and guaranteed bit rate (GBR) bearers for the failed node's UEs discarding the Non-GBR bearers. Moving GBR or Non-GBR bearers traffics to the default bearer can be another option. But this strategy will degrade the performance of the services on those bearers that requires specific QoS. The traffic restoration could also be restricted to high priority traffic or premium customers, if spare capacity is limited.

The advantage of this approach is that it is cost-effective. Furthermore, adoption of distributed architecture may decrease restoration time as it does not use the OAM to perform failover recovery tasks which require more message exchange and processing time. However, as mentioned above, adequate resources must be allocated to avoid performance degradation. Some of the major disadvantages of this approach are lack of co-ordination and implementation complexities considering large numbers of network elements involvement in failure recovery process. However, for the sake of proper and efficient failover strategy implementation, the centralized approach can be adopted in this configuration.

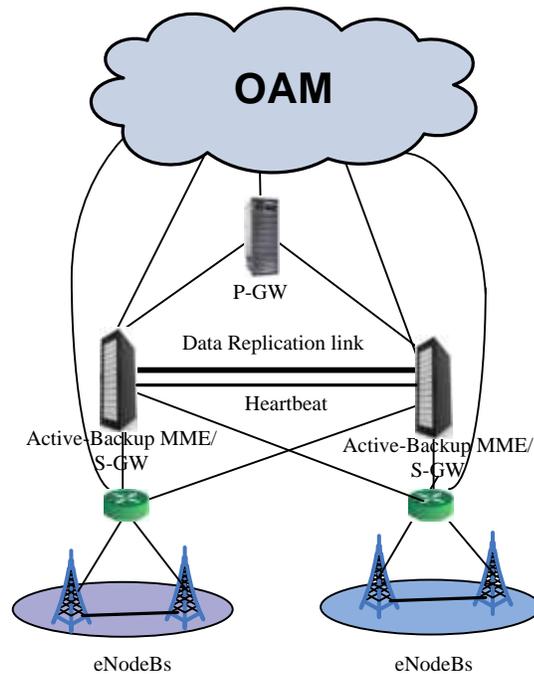


Figure 3-4: 1:1 Active-Active configuration.

3.1.2 Failure Detection, Notification and Fault Isolation

Failure detection is the first step in a self-healing solution. The active nodes may encounter different kinds of faults - hardware or software, partial or complete outage etc. For the proposed solutions, our assumption is that the active nodes will completely shutdown, i.e., there will be a complete outage. Failure detection has significant impact on overall performances of the healing mechanisms in terms of data loss, restoration time and user perceived quality of experience (QoE). In the N:M active-backup and 1-1 active-active configuration, the network OAM and the peer active node perform the failure detection of operational MME/S-GWs, respectively. Periodic heartbeat message in the form of "hello" or "I am alive" can be exchanged between the OAM and target active nodes or between the peer nodes. If the monitoring node does not receive response within a predefined threshold time, it considers the node as failed or down. It should be noted that several indicators such as SCTP keep-alive message which are exchanged over the S1 interface

between the eNodeB and MME can be employed to detect MME/S-GW failure from the eNodeB side [31]. By analyzing the heartbeat information, the OAM or the peer node detects the failure and notifies the affected eNodeBs and P-GW about this failure incident and order the backup MME/S-GW to take over the functionalities of the failed node. Usually, in a large network one MME or S-GW supports hundreds of eNodeBs. In this context, a centralized approach of sending failure notification to large number of network elements may cause link delay and congestion based on the network configuration and, result in performance degradation. To improve the scalability and reliability, a hybrid approach can be adopted where eNodeBs perform the task of disseminating or forwarding the notification message to neighbor eNodeBs over the X2 interface. It is the responsibility of the eNodeBs to inform corresponding active UEs about failure incidents. Upon reception of the notification message, the UE initiates the recovery procedures.

Faulty device isolation is another important step of the fault-tolerant and self-healing system. The purpose of fault isolation is to minimize the consequences of the problem and thereby, limiting the impact of the failure. Though failure recovery process has significant impact on the performance of the network but fault in the MME/S-GW does not propagate to other active MME/S-GWs or cause deterioration of services provided by those MME/S-GWs. Newly arrived UEs that are registered to the failed MME/S-GW before the failure incident and are not aware of the current status of the MME/S-GW must be prevented from establishing connection to it. The S1-flex functions of the eNodeB play a vital role in this regard. By utilizing this feature the eNodeB can forward the request of these UEs to the backup MME/S-GW.

3.1.3 Failover Recovery Coordination

Failover recovery coordination process is described here from the centralized architecture's points of view. When a new MME/S-GW is added to the network, the OAM assigns a dedicated

backup node for it from the backup MME/S-GW pool. The active node may start replicating UE contexts data to the backup node immediately. Requirements for selecting a backup node by the OAM may vary based on the type of deployment scenario and network status. For example, one of the main requirements is that the selected backup node should have sufficient resources to support this new node. However, considering real life deployment scenario, the backup node may not have enough capacity or resources to support more than one node in multiple failure scenarios. Other backup nodes from the backup MME/S-GW pool which are dedicated to some other groups of active MME/S-GWs may take part in providing service. However, it will take more time in restoration process in this case since UE contexts information needs to be copied to that backup node as well. To facilitate effective failover strategy, the OAM should keep track of all the changes that take place in the network and updates related nodes accordingly. Various changes in terms of load, configuration and policy may occur in the network. The capacity of any active MME/S-GW may increase in the course of time due to the eNodeB addition or the subscriber provisioning policy. The OAM may have priority policy for supporting specific MME/S-GWs in simultaneous large scale failures. Even within the eNodeBs, priority can be given to particular eNodeBs during recovery process considering subscriber based or revenue generation statistics. All of this dynamic information is needed to be exchanged between the active and backup nodes, and the OAM for efficient network operation. To facilitate the coordination and decision making process, a recovery manager module is introduced in the active and backup nodes, and the OAM which operates in client-server manner. A sample message sequence diagram is presented in Figure 3-5 that can be used in this purpose. Structures and contents of the messages may vary based on the operator's implementation policy and needs.

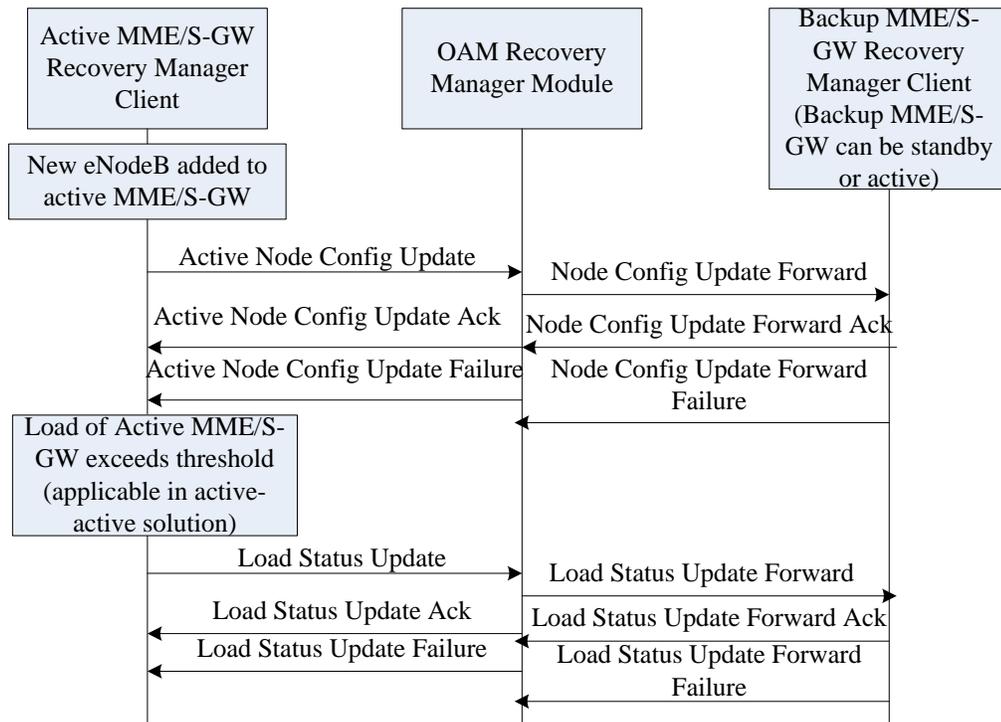


Figure 3-5: Node status update for failover recovery coordination.

3.1.4 Recovery Mechanism and Service Continuity Procedures

This section focuses on the failover and service continuity procedures for the centralized architecture where the OAM is responsible for overall self-healing tasks. When an active MME/S-GW fails, the OAM recovery manager module detects the failure from the disappearance of the heartbeat message. When the failure is confirmed, it quickly adopts the failover procedures based on the latest configuration and status update information and the recovery policy. It then sends notification messages to the backup MME/S-GW, all eNodeBs, P-GWs connected to the failed node about this incident. Upon receipt of the message, the eNodeB releases its S1 connection to the failed MME/S-GW and initiates standard S1AP connection setup [43] procedure with the backup MME/S-GW. It provides information of the MME/S-GW failure to the UEs by releasing the RRC connection with the release cause *EPC re-attachment required*. In the

meantime, the eNodeBs and P-GW start to buffer active UEs uplink and downlink data, respectively. However, the nodes may drop packets if buffer size reaches a certain threshold. When a UE receives the failure notification message from the eNodeB, it immediately stops sending application data and suspends the bearer activity. It keeps the bearers in the suspended state until further notification arrives from the eNodeB. The UE then initiates re-attachment procedures with the backup MME by sending the NAS attach request [41] message. It should be mentioned here that all the terminals or UEs should be aware of the identity of the backup MME/S-GW. This can be achieved in three ways. First, the identity of the backup MME/S-GW can be assigned to the UEs during initial network registration process along with the active MME/S-GW. The UE stores this information and uses the identity of the backup MME/S-GW to attach to the network during failover procedure. The disadvantage is that if the backup MME/S-GW changes afterwards by the OAM, the network needs to send additional messages to the UE to reconfigure the backup node information. Second, the eNodeB can perform the task of redirecting the attach request of the UEs to the backup MME based on the GUTI information. The GUTI is given to the UE by the active MME when it first registers to the network. The GUTI contains the Globally Unique MME ID (GUMMEI) which the UE always uses whenever it connects to the network. Thus, using the MME id in the GUTI information, the eNodeB can redirect the request to the appropriate backup MME/S-GW. In that case, the OAM should provide the eNodeB with the backup MME/S-GW information before failure incident. Third, the OAM can send backup MME/S-GW information to the eNodeBs along with the failure notification message. The eNodeB, in turn, includes this information with the RRC connection release message which it sends to the UE after failure. Figure 3-6 shows the failover procedures and the communications between the network elements.

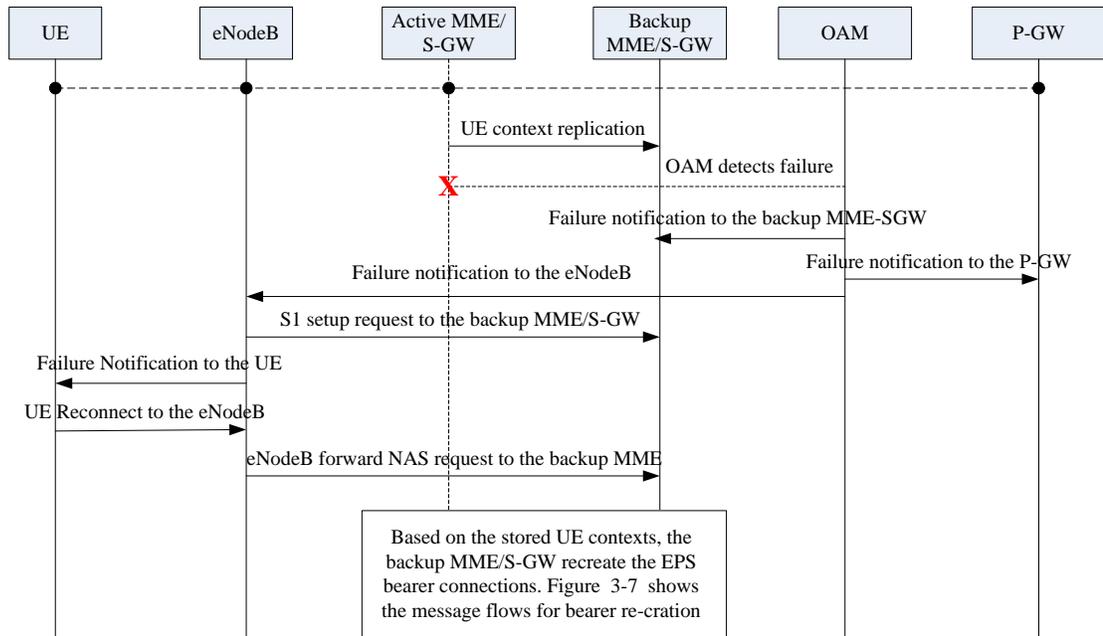


Figure 3-6: Interaction between the network entities during failover.

The Backup MME receives connection request message from the UE and verifies its context information based on the IMSI value. Then it starts the bearer re-creation procedures as depicted in Figure 3-7. As the EPS bearer information is already available in the backup MME for the UE, the MME sends UE's bearers list along with the IMSI, the eNodeB address and eNodeB TEIDs (each bearer is identified by a TEID) to the S-GW by sending the *Create Bearer Request* message. The S-GW updates corresponding EPS bearer table entry for this UE, and forwards this message to the associate P-GW including the S-GW addresses and S-GW TEIDs for the data and control plane. The S-GW may select only the default and GBR bearer to re-create if there is any resource shortage issue. A situation like this may arise in multiple MME/S-GW failure situations. After updating the EPS bearer context table for this UE, the P-GW sends the *Create Bearer Response* message to the S-GW along with the data plane and control plane P-GW addresses and P-GW TEIDs. It is assumed that the P-GW maintains same IP address for the UE during the

failover and bearer re-creation process as allocated before the failure. The S-GW forwards the S-GW and P-GW addresses, the S-GW and P-GW TEIDs for which bearer re-creation is successful to the MME via the *Create Bearer Response* message. In the meantime, P-GW sends downlink data for this UE to the S-GW. The MME forwards the bearer re-creation information to the eNodeB by sending the *Bearer Setup Complete* message. It may send the *Bearer Setup Failure* message if the bearer re-creation is unsuccessful.

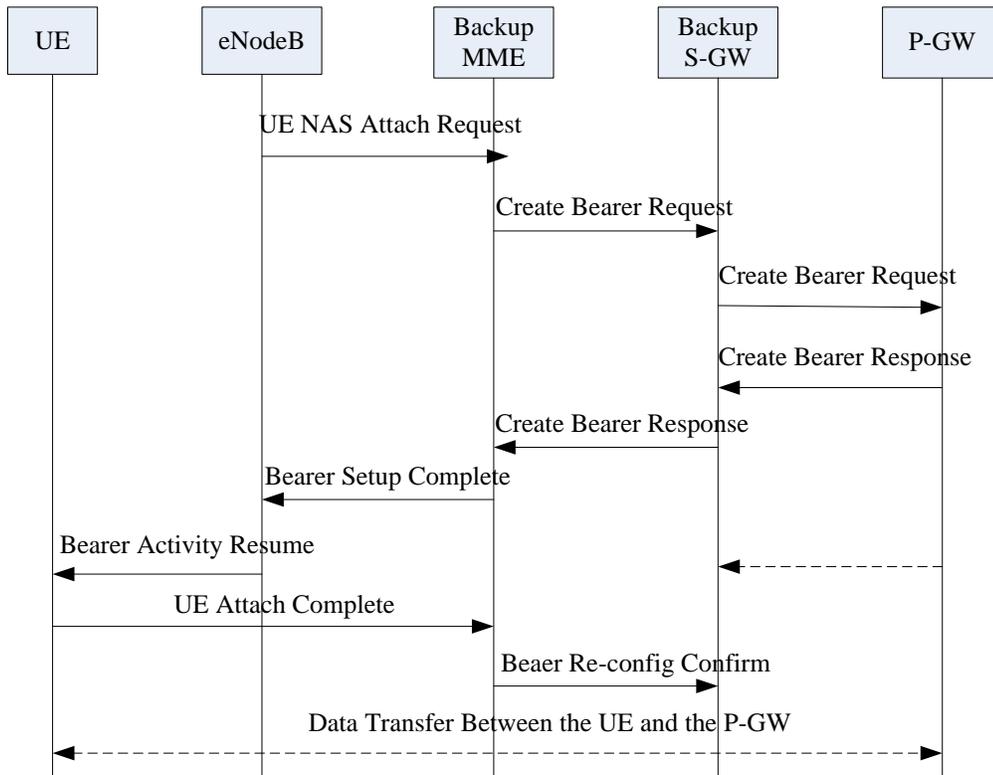


Figure 3-7: Bearer-recreation and related message flows.

Subsequently, the eNodeB informs the UE to start transfer application data by sending the *Bearer Activity Resume* message where it provides the list of bearers that are accepted by the backup MME/S-GW. In response, the UE sends the *UE Attach Complete* message to the backup MME. Upon receipt the confirmation message from the UE; the MME sends the *Bearer Re-config*

Confirm message to the S-GW. In the meantime, the eNodeB transmits buffered uplink data to the S-GW, and the S-GW sends downlink data to the eNodeB.

This procedure can be easily extended for the distributed architecture. In either the centralized or the distributed architecture, main aspects of the self-healing schemes have been clearly defined in this section.

Chapter 4

Performance Evaluation of the Self-healing Scheme

In this chapter the performances of the proposed self-healing solutions are evaluated and analyzed under various failure scenarios. The system was modeled in OPNET simulation environment. In addition, MATLAB was used to measure signaling message overhead. A series of simulation experiments were performed in OPNET for different network configurations and scenarios. Some important quantitative results that are collected from the simulation are also discussed.

4.1 Simulation Setup

The network considered here is composed of UEs, eNodeBs, EPCs, monitoring nodes and the server. The centralized approach where the OAM is responsible for failure recovery process execution is adopted in the simulation environment. An eNodeB manages a hexagonal cell where the UEs are randomly dispersed. In our model, a single EPC node represents a combined network of the MME, S-GW and P-GW, which is often a realistic scenario in LTE deployments. The external server represents a Packet Data Network (PDN). The monitoring nodes are connected to the EPCs and eNodeBs. The failure detection and notification modules of the self-healing system are implemented in the monitoring nodes. The UEs select a cell or eNodeB based on the 'First Suitable eNodeB' policy which dictates that the UE chooses the first eNodeB that matches cell selection criteria. The selection criteria for a cell are evaluated based on the comparison between the signal strength parameters of the cell, RSRP (Reference Signal Received Power) and the minimum required strength of a UE. A UE is initially configured to be connected to one active EPC. An eNodeB may communicate with one or more EPCs.

To simplify the simulation model, the EPS bearer configurations of the UEs were made globally accessible so that the UE context replications were not required in current simulation setup. A failure scheduling module was also configured to model the active nodes failure. In this model an active EPC fails at certain time of simulation and remains failed until the end of the simulation.

To simulate the N:1 Active-Backup and 1:1 Active-Active configurations, four networks setup were considered in terms of different number of UEs, eNodeBs and EPCs. Table 4-1 illustrates the simulation assumption for network elements and terminals for these four network configurations. Simulation times for each of these network configurations are also mentioned.

Table 4-1: Network configurations in the simulation

N:1 Active-backup Configuration					
	Case No.	Active node UEs	Number of eNodeBs	UEs per eNodeB	Simulation Time
1:1 Active-Backup	#1	500	10	50	175 seconds
1:1 Active-Backup (with load in backhaul)	#2	500	10	50	200 seconds
2:1 Active-Backup	#3	300	10	60	175 seconds
1:1 Active-Active Configuration					
	#4	200	10	40	175 seconds

In the network, each UE selected an OPNET defined bronze type EPS bearer traffic model to send an IP Unicast traffic flow of 1,000,000 bits/second to the server along with the default bearer. All eNodeBs, EPCs, monitoring nodes and external servers were connected to the

backbone network with 1000Base-X Gigabit Ethernet links. Monitoring nodes send heartbeat messages to active EPCs every 2 seconds. This heartbeat message contains simply a “Hello” string. In return the active EPC sends “Hello OK” message. The monitoring node starts a timer when it sends a heartbeat message. This setting determines how long the monitoring node waits for the “Hello OK” message after sending the “Hello” message. In the simulation, the expiration of the timer was set to 1 second, therefore, if the monitoring nodes does not receive response message from an active EPC within one second after sending a heartbeat message, it considers that active EPC as failed. For all scenarios, the failure scheduling module configures single or multiple EPC failure events at t=147 seconds. It should be mentioned here that the simulation time includes the UE power up time and provides adequate time for IP routing protocols to converge.

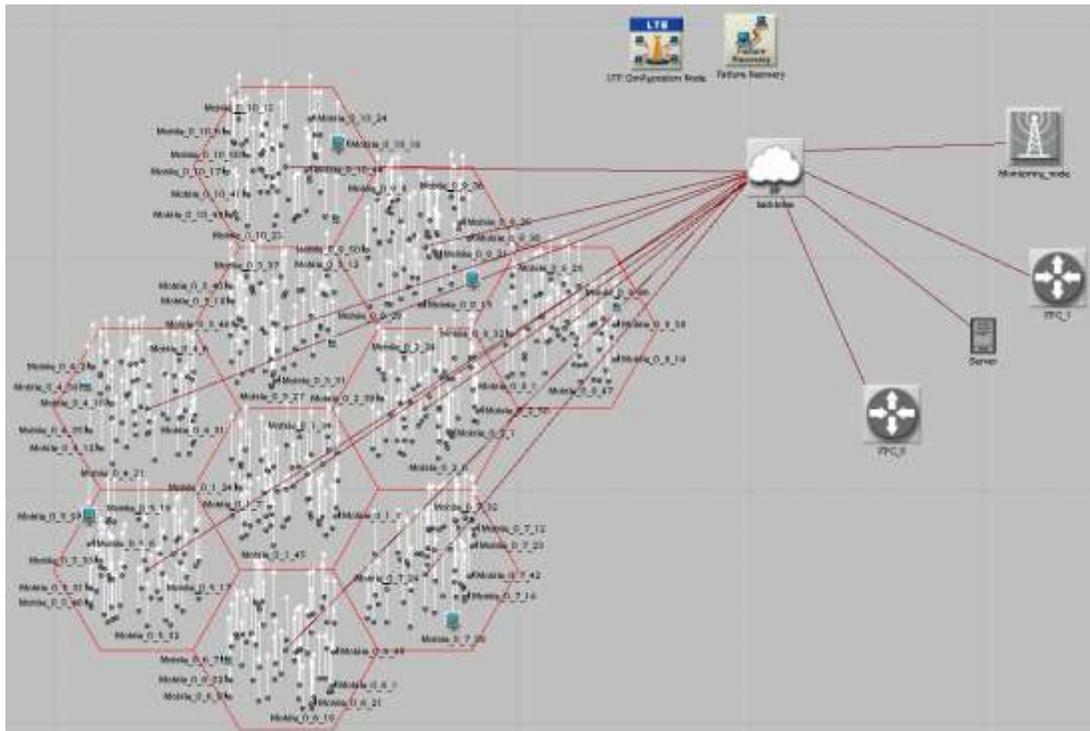


Figure 4-1: Sample network configuration.

Figure 4-1 shows a sample network setup and corresponding network elements in the simulation. The main simulation assumptions are given in Table 4-2. These settings are used to obtain results presented in the remaining sections of this chapter. Some of the limitations of the OPNET simulation model are as follows [32]:

- The sizes of RRC messages are approximated.
- The sizes of S1 messages are approximated.
- Stream Control Transmission Protocol (SCTP) over S1-C interface is not supported.

Table 4-2: Simulation parameters and settings.

Entity	Parameters	Value	Unit
UE	Antenna Gain	-1	dBi
	Battery Capacity	5.0	Watt-Hour
	Maximum Transmission Power	0.0101	Watt
	Modulation and Coding Scheme	9 (specified in 3GPP TS 36.213)	
	Multipath Channel Model	ITU Pedestrian B	
	Number of Receive Antennas	1	
	Number of Transmit Antennas	1	
	Pathloss Parameters	Free Space	

	Mobility	Random Waypoint	
	Speed	5	meter/second
eNodeB	Physical Profile	Frequency Division Duplex (FDD)	
	Channel Bandwidth	LTE 20 MHz FDD	MHz
	Cell Radius	2	Km
	Antenna Gain	15	dBi
	Battery Capacity	Unlimited	
	Maximum Transmission Power	0.0222	Watt
	Number of Receive Antennas	2	
	Number of Transmit Antennas	2	
	Path loss Model	Free Space	
	Terrain Type (Suburban Fixed)	Terrain Type A	
	eNodeB Selection Threshold	- 110	dBm
	Traffic Model	Traffic Flow	IP Unicast
Traffic Intensity		1,000,000	bits /second
Traffic Intensity		100	packets /second

	Type of Service	Best Effort	
	Traffic Duration	End of Simulation	Second

4.2 Performance Metrics

The main performance metric used to study the performance of the EPC network failure recovery was the service restoration time. Another important metric was the signaling traffic overhead resulting from the proposed failover recovery mechanism and related message flows. Various other performance factors such as throughput, uplink or downlink delay, EPS bearer delay that are associated with the recovery process were also examined. In the following sections, an analysis of these parameters for the LTE self healing solution is discussed and then simulation results are presented.

4.2.1 Service Restoration Time

Based on the discussion in chapter 3, service restoration time for a particular UE is defined as the time interval between the moment when the active EPC goes down and the moment when the UE sets its suspended bearers status to active again. The following notations are used in our analysis for various delay components:

- T_{dm} : Failure detection delay by the monitoring node.
- T_{ne} : Notification delay between the monitoring node and eNodeB.
- T_{ur} : Time required for sending the connection release message to the UE with appropriate release cause.
- T_{me} : Time required for the connection re-attachment towards the backup EPC.
- T_{br} : Delay for the bearer re-creation procedure.

T_{rt} which denotes the average restoration time is derived below:

$$T_{rt} = T_{dm} + T_{ne} + T_{ur} + T_{me} + T_{br} \quad (1)$$

Some of the events occurrence times (Δt) pertaining to the measurement of the above mentioned delay components have been defined as follows:

$t1$: Δt EPC failure time specified by the failure schedule module.

$t2$: Δt Failure detection by the monitoring nodes.

$t3$: Δt Failure notification message received by an eNodeB.

$t4$: Δt Failure notification message received by the UE with the re-attachment instruction.

$t5$: Δt Network attachment accept message received by the UE.

$t6$: Δt Bearer activation complete message received by the UE.

Table 4-3 shows the measurements of delay components.

Table 4-3: Measurement of restoration delay parameters

Parameters	Delay Measurement
T_{dm}	$t2 - t1$
T_{ne}	$t3 - t2$
T_{ur}	$t4 - t3$
T_{me}	$t5 - t4$
T_{br}	$t6 - t5$

It should be noted that the failure detection and the failure notification procedures are implemented in the OAM side. The performance of these two parameters does not depend on the LTE network. Furthermore, the eNodeBs adopt random notification delay (T_{ur}) while sending failure notification message to the UEs in order to prevent concurrent network re-attachment

process. Random uniform distribution with the mean of 10ms is used in this purpose. As a result, main parameters in these simulation experiments are T_{me} and T_{br} delays as these parameters are directly related to the LTE access and core network architectures and protocols.

4.2.2 Signaling Message Overhead

In this section, an analysis is presented to measure the signaling load at the EPC induced by the proposed approach during the recovery time. These signaling messages are divided into two categories: *i*) messages that are required to recover the failed sessions, and *ii*) messages that are generated due to the new UE arrivals during the recovery process. Therefore, signaling procedures between the E-UTRAN and EPC network, more specifically, across the S1-MME/S1-C, S11 and S5/S8 interfaces are the focus of this analysis. It should be mentioned here that other control messages related to the failure recovery management such as UE context replication, MME/S-GW monitoring, failure notification to the eNodeB and P-GW, which involve the OAM, E-UTRAN and EPC are described in chapter 5.

In [53] a method to quantify the signaling load at the MME for the UE originated session is discussed. The UE-originated session establishes an EPS bearer to the P-GW, via the S-GW and eNodeB. It should be noted that in the proposed approach, the UE initiate network re-attachment procedure after failure, which in turns triggers the bearer re-creation procedure at the EPC.

Let λ_A be the average arrival rate of an application session originated at a UE. The application can be voice call, http request etc. In Figure 3-7 in chapter 3, it is shown that the backup MME/S-GW requires a total of 8 messages (incoming and outgoing) to reconfigure the failed bearers or sessions originated from active UEs. For simplicity, it is assumed that all the messages have the same size. Let T_d denote the average duration of the sessions. Then the average number of active sessions, N_S at the time of failure is:

$$N_S = \lambda_A * T_d \quad (2)$$

Thus, if the total number of active subscribers at the time of failure is C then the total number of messages generated at the MME in order to recover failed session,

$$N_T = 8 * N_S * C \quad (3)$$

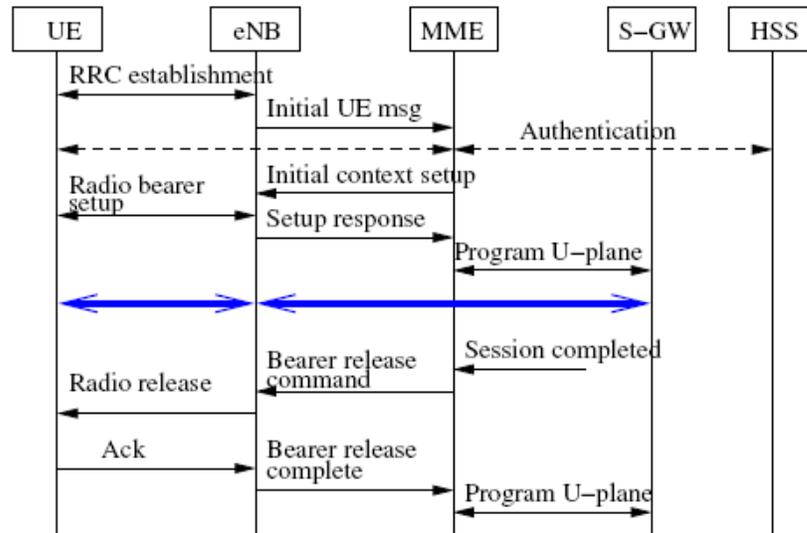


Figure 4-2: Message flows for UE originated sessions [53].

For each new UE originated sessions during recovery time, MME/S-GW processes 10 messages [53] as depicted in Figure 4-2. In this procedure, the IDLE mode UE establishes RRC connection to the eNodeB. Then the UE creates a service request. The eNodeB forwards the request to a MME through initial UE message which contains the identity of the UE and the NAS packet data unit (PDU). Upon reception of the session, MME creates an identity for the UE which is used within the MME and provides bearer level information such as, bearer identity and corresponding QoS, IP address, GTP tunnel ID between the S-GW and the eNodeB through the initial context setup request. If this process is successful, the MME initiates bearer setup procedures between the eNodeB and the S-GW. Then, data transfer takes place between the UE and S-GW. The MME also can release a bearer by sending the bearer release command message to the eNodeB.

When new UEs arrive at the eNodeB during the restoration time, it is the responsibility of the eNodeB to re-direct the request to the backup MME/S-GW as these UEs are not aware of the active MME/S-GW failure incident. As mentioned in chapter 3, using the GUTI information of the UE, an eNodeB can forward attach request to the backup MME/S-GW.

Lets T_R denotes recovery time in seconds and λ_{Ak} be the number of session originated by new arrival of a UE. The arrival of new UE per second is denoted by λ_{UE} . Thus, the total number of messages generated at MME/S-GW for new sessions,

$$N_T = 10 * T_R * \lambda_{Ak} * \lambda_{UE} \quad (4)$$

4.3 Results and Analysis

At first, some performance data related to the cases mentioned in Table 4-1 are presented in this section. Then restoration time and signaling message overheads described in previous section are evaluated.

4.3.1 N:1 Active-Backup Configuration

In this section, simulations results for 1:1 and 2:1 Active-Backup configurations for Case# 1 to Case# 3 mentioned in Table 4-1 are discussed. Figure 4-3 presents the uplink throughput of 1:1 active-backup configuration (Case#1) before and after failure. In this configuration, the active EPC handles 500 users. The backup EPC must re-create the bearers of these users after failure. A decline in the traffic indicates the failure of the core network which is followed by the traffic rise around at $t = 149$ seconds. In this case, the backup EPC and the backhaul link have sufficient resources and the session recovery rate is 100%. The active EPC fails at $t = 147$ s and it takes 2s for the OAM to detect the failure due to the adopted detection policy discussed in section 3.1.2 of chapter 3. Then it triggers recovery procedures by sending notification message to the eNodeBs

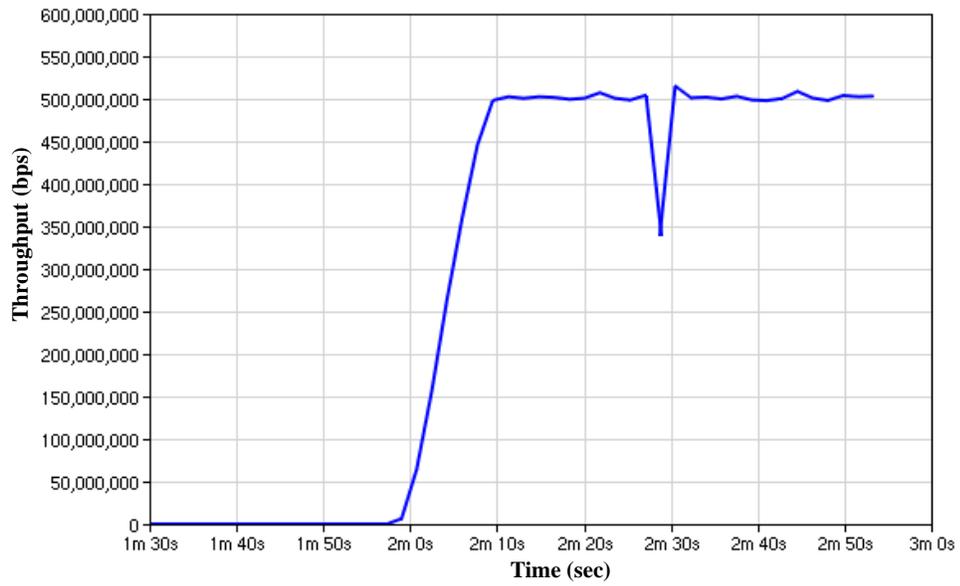


Figure 4-3: Uplink throughput for Case# 1.

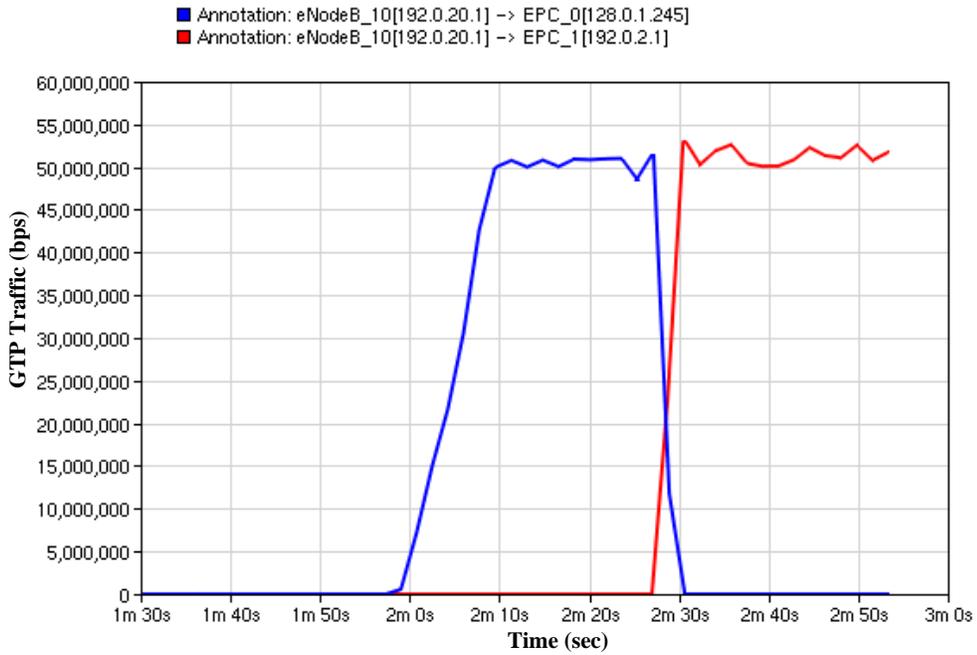


Figure 4-4: GTP traffic switchover from the active EPC to the backup EPC (Case# 1).

around at $t = 149s$. Event details and measurements of delays related to the restoration procedures are discussed later in this thesis.

Figure 4-4 illustrates the switchover of the GTP (GPRS Tunneling Protocol) layer data traffic over the S1-U interface from the active EPC (EPC_0) to the backup EPC (EPC_1) network for an eNodeB. There is a GTP tunnel for each EPS bearer and the corresponding UEs can be identified with each GTP tunnel. The traffic presented in this figure represents the aggregate GTP traffic of all the UEs of the eNodeB 10.

Figure 4-5 depicts the uplink delay of the traffics in Case# 1. The uplink delay is defined in seconds and measured from the time the traffic arrives to the LTE layer of the UEs until it is delivered to the higher layer of the corresponding eNodeBs. As depicted in the graph uplink delay increases during recovery process which occurs around at $t = 2m\ 30s$ in the graph.

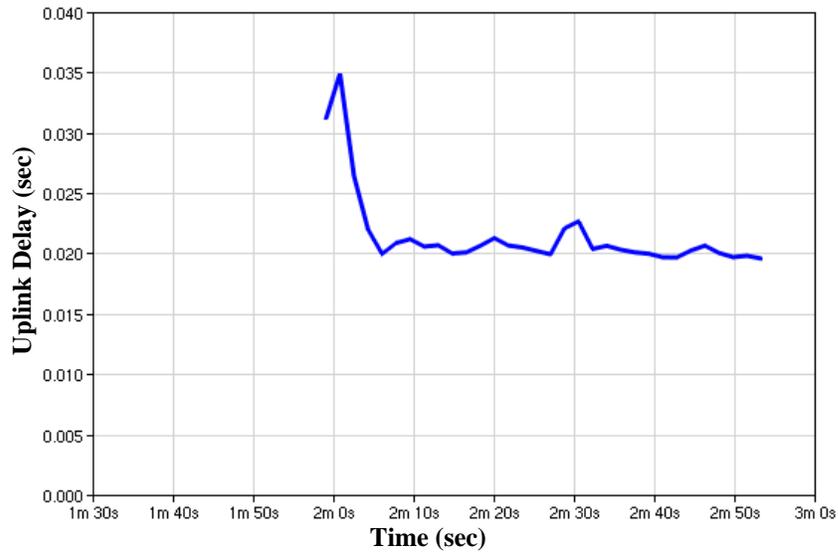


Figure 4-5: Uplink delay for Case# 1.

The EPS bearer delay for all the traffic received at eNodeB 10 for the UEs is plotted in Figure 4-6. It can be seen from the graph that the bearer delay at the time of recovery process increases.

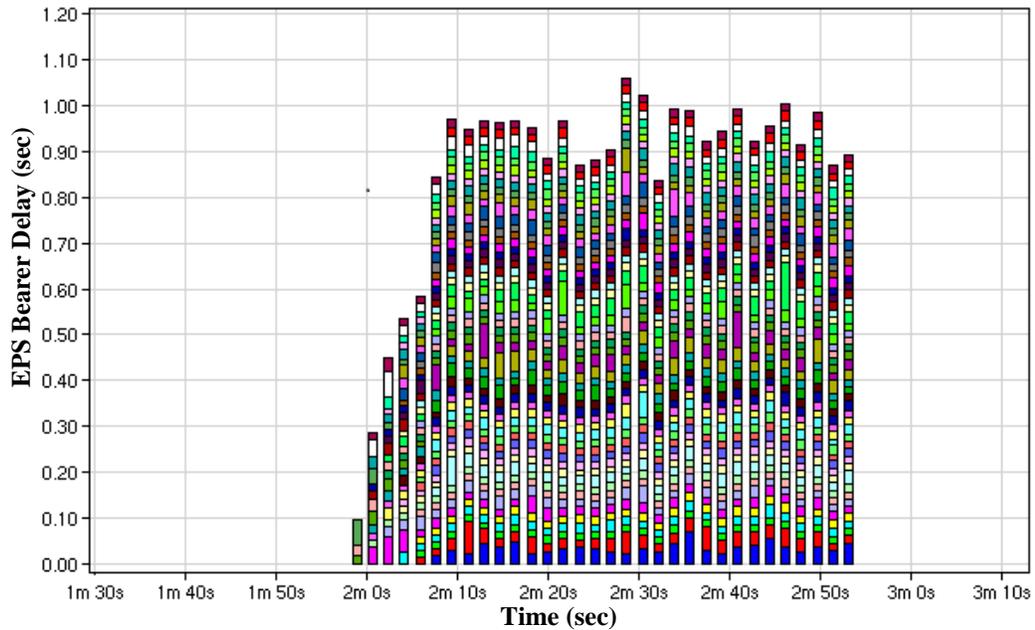


Figure 4-6: EPS bearer delay of eNodeB_10 for Case# 1.

The Case# 2 network configuration is similar to the Case# 1 except that 700 Mbps background load was configured on 1 Gigabit backup backhaul link. The background load results in congestion during failover and causes significant delay in restoration time. A quantitative analysis of this delay is presented later in this chapter. In this case, user sessions recovery rate is 98.8%. All the users successfully connected to the network after failure but corresponding bearer re-creation procedures were not completed for all the UEs within the simulation time due to the congestion and increased packet waiting time in the backhaul network. It should be noted that data packets are always routed into the corresponding EPS bearer. So, unless failed bearer is activated UEs cannot transmit data through the network. Furthermore, as illustrated in Figure 4-7, user sessions were restored at different period of time. Around 5.8% user's bearer re-activation process took unusually long time. More specifically, 93% user's bearers were restored around at t

= 150s of the simulation time. Approximately 2.4% user's bearers restored around at t = 166s and 3.4% user's bearers restored around at t = 180s.

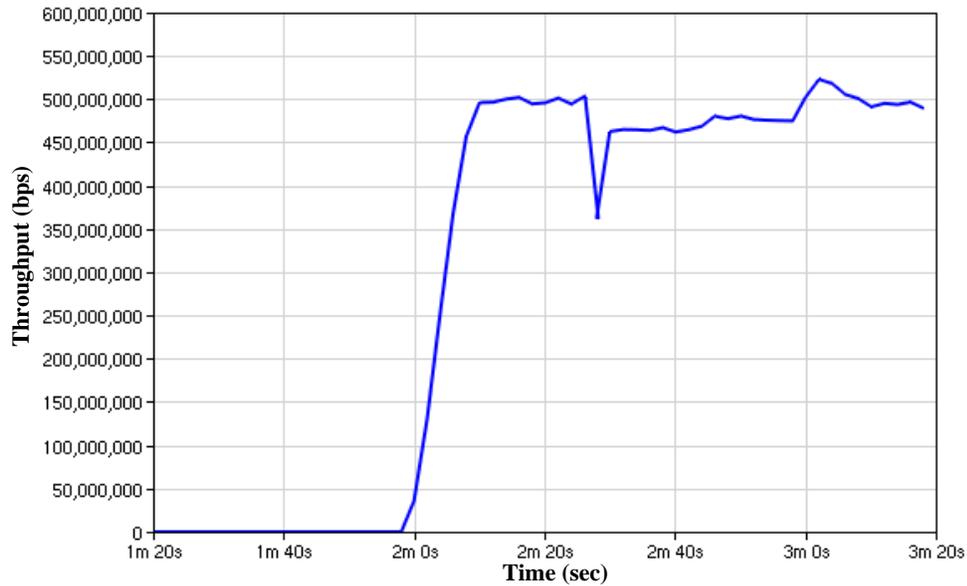


Figure 4-7: Uplink throughput for Case# 2.

It can be seen from Figure 4-8, the uplink delay for Case# 2 significantly increases around at t = 180 seconds and this delay reaches 0.27 seconds which is nine times higher than average uplink delay for the timely restored user sessions. Figure 4-9 shows the backhaul link utilization of the active and backup EPCs. The link utilization of the backup backhaul connection reached 100% after failure due to the presence of configured background load (700 Mbps) and re-directed UE throughput which in turns increases packet waiting time. The EPS bearer traffics received by the eNodeB 10 for some of the UEs are shown in Figure 4-10. As it can be seen from this graph, the bearer re-creation process was not completed for all the UEs at the same time due to the backhaul network congestion. For instance, bearers for the UE with IMSI 470 were not re-created during the simulation time. The bearer restoration process for the UE with IMSI 485 and 486 started

around at $t = 180\text{s}$ (3m). However, bearer restoration was completed for the UE with IMSI 464 just after failure.

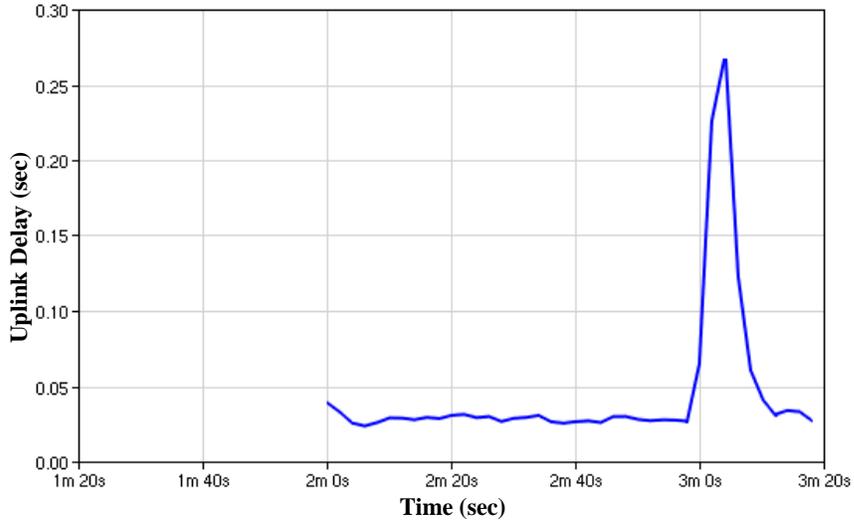


Figure 4-8: Uplink delay for Case #2.

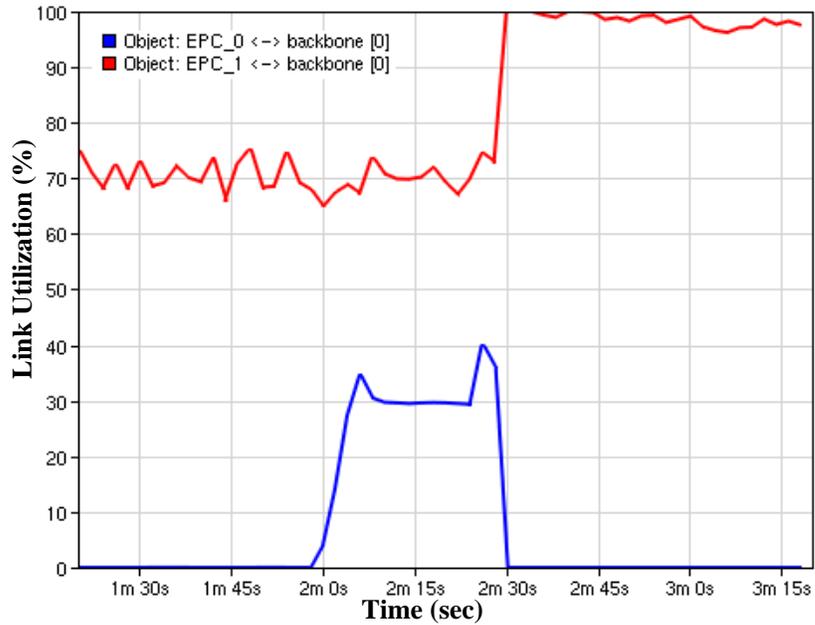


Figure 4-9: Link utilization of active (EPC_0) and backup (EPC_1) backhaul connections.

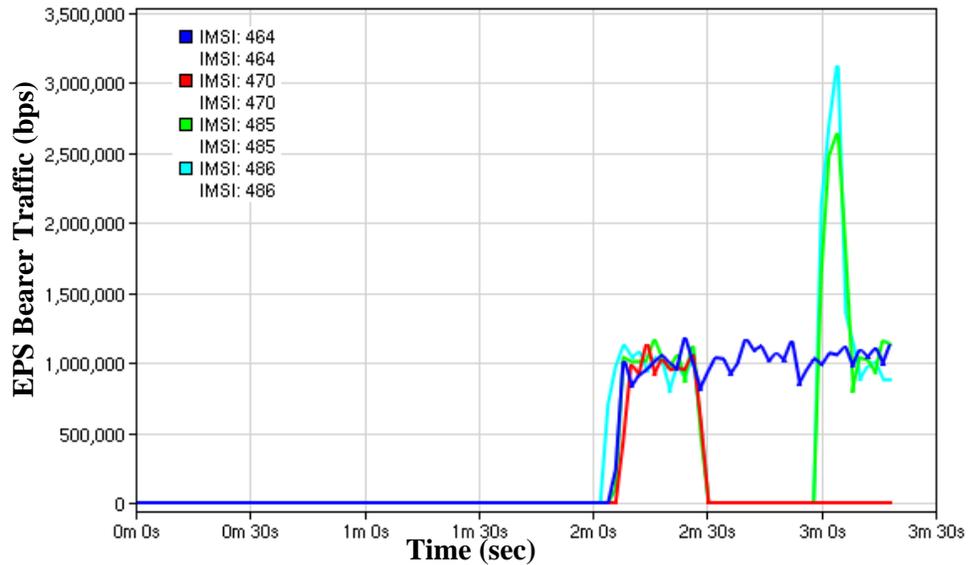


Figure 4-10: EPS bearer traffic received by the eNodeB 10 for Case# 2.

Case# 3 represents a 2:1 Active-Backup configuration where each EPC handled 300 users. The backup EPC was configured to handle 600 users. However, the simulation output shows that the success rate of session recovery was 99.5% in this case. Though the network attachment and default bearer re-creation were successful for all the users in this case but the bronze bearers re-creation were not successful for 3 users. The reason of the failure is that the MME failed to activate bronze bearers for the UE with IMSI 306, 401, 494 during simulation time. Figure 4-11 shows the uplink delay for this configuration. Figure 4-12 presents PDSCH (Physical Downlink Shared Channel) channel utilization. The PDSCH channel uses DL-SCH or Downlink Shared Channel that carries downlink data and signaling messages from the eNodeB to the UE.

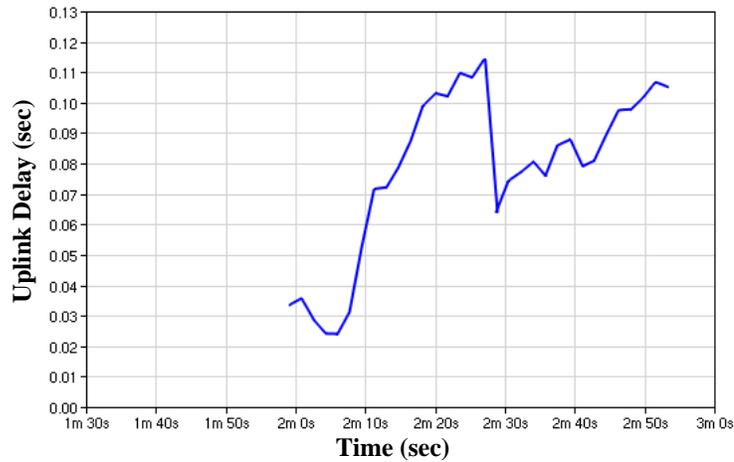


Figure 4-11: Uplink delay for Case# 3.

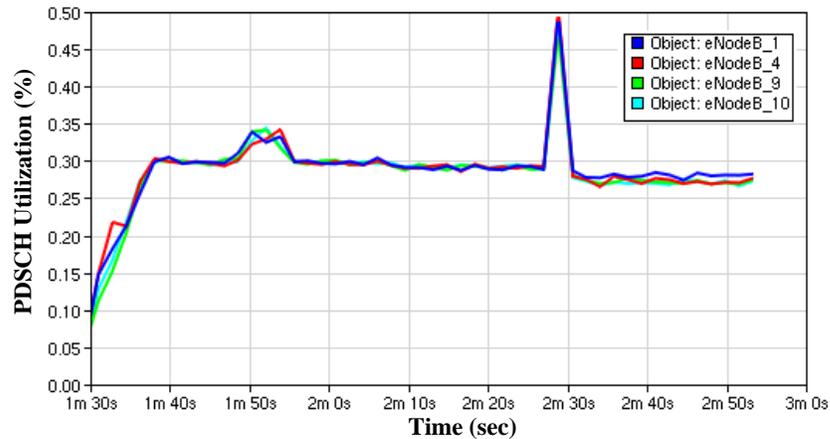


Figure 4-12: PDSCH channel utilization for Case# 3.

4.3.2 1:1 Active-Active Configuration

As discussed in chapter 3, both EPCs in the active-active configurations work as backup node for each other. Case# 4 in Table 4-1 represents an active-active configuration where each EPC has 200 UEs. When one active EPC fails, the peer active EPC takes over the failed UE's sessions. A detail comparison between the 1:1 Active-Active configuration and N:1 Active-Backup configuration is presented in the next section in terms of delays involved in restoration

procedures. Figure 4-13 and 4-14 shows uplink delay and EPS bearer delay of the eNodeB 10, respectively for this configuration. These graphs show that both the uplink delays and bearer recreation delays for the Case# 3 are relatively higher than the other cases.

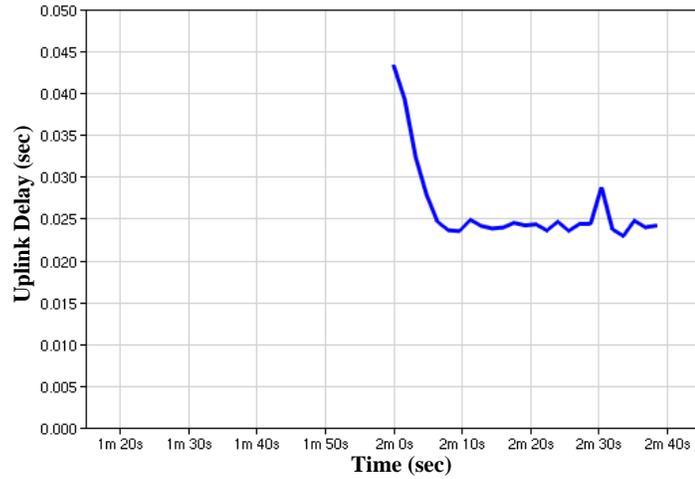


Figure 4-13: Uplink delay for Case# 4.

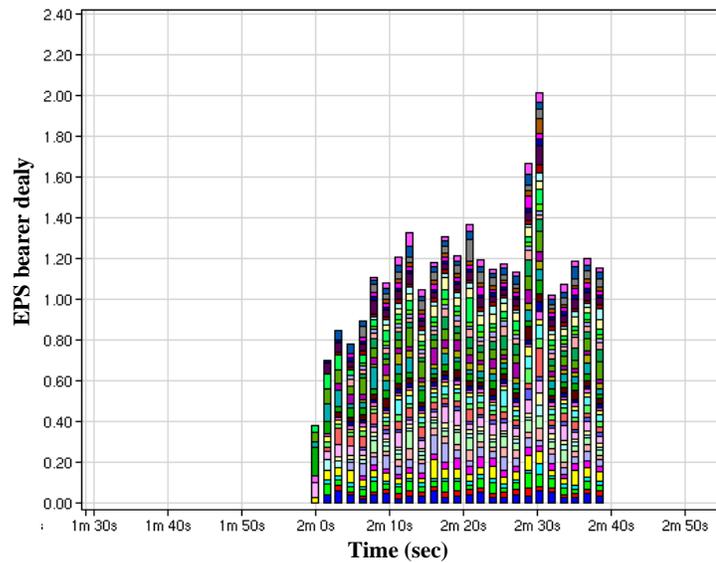


Figure 4-14: EPS bearer delay for Case# 4.

4.3.3 Service Restoration Time Calculation

Assuming that the processing delay is negligible, our analysis discussed in section 4.2.1 shows that the time elapsed in failure detection (T_{dm}) and notification (T_{ne}) which is a little more than 2s in current simulation configuration is a significant portion of the total restoration time. It means that the restoration time can be significantly improved if the failure detection time or *hello* message interval time is reduced in current simulation. It should be mentioned here that though the notification delay is typically small, it may increase due to the traffic congestion, particularly in a network with a large number of eNodeBs. Furthermore, if failure strategy revision and decision making scenarios are invoked as described in chapter 3, it will add more delays in the failure notification and faulty device isolation procedures. However, as mentioned in performance metric section, T_{me} and T_{br} delays are the focus of this analysis. Excluding T_{dm} and T_{ne} delays, equation 1 can be redefined as follows:

$$T_{rt} = T_{ur} + T_{me} + T_{br} \quad (5)$$

Table 4-4 shows average values of T_{rt} , T_{me} and T_{br} for the cases mentioned in Table 4-1.

Table 4-4: Simulation results for average values of T_{rt} , T_{me} and T_{br} .

Case	T_{rt} (ms)	T_{me} (ms)	T_{br} (ms)	Restoration Success Rate
#1	582	511	62	100%
#2	2015	899	1106	98.5%
#3	675	559	107	99.5%
#4	853	676	168	100%

It can be seen from the results of Table 4-4 that the restoration time, connection re-attachment and bearer re-creation delays in Case# 4 take longer time than the corresponding delays in other scenarios where more users are involved in recovery process and sufficient resources are available (Case# 1 and Case# 3). In this case, 200 users of the failed EPC are restored through the backup EPC. It should be noted that 500 user's default and bronze bearers were re-created in the backup EPC in Case# 1 while the numbers of re-created bronze bearers for the users in the backup EPC were 597 in Case# 3. For better comparison, we further collected data from two scenarios: 1) 1:1 Active-Backup configuration with 200 UEs and 2) 2:1 Active-Backup configuration with 300 UEs where each active node handles 150 UEs. Table 4-5 provides performance data. Simulation assumptions described at the beginning of this chapter for four cases were also applied in these scenarios. Comparing the values in Table 4-4 and Table 4-5 for Case# 4 scenario, it can be observed that the restoration delays are higher in 1:1 Active-Active configuration than other N:1 Active-Backup configurations. As the re-attachment delay (T_{me}) may vary based on the radio network condition, we can look into the bearer re-creation delays where it took 168 ms for 200 users (Case# 4). As the backup node is active in this configuration and engaged in handling the bearers for 200 users of its own, re-creation of additional bearers for 200 users of the failed node causes some delays.

Table 4-5: Simulation results for scenario 1 and 2.

Scenarios	T_{rt} (ms)	T_{me} (ms)	T_{br} (ms)	Restoration Success Rate
Scenario# 1	562	530	22	100%
Scenario# 2	535	506	22	100%

Further simulations were conducted to check how the restoration delay components change in the 1:1 Active-Backup configuration (e.g., Case#1) when number of users varies. Figure 4-15 shows T_{rt} , T_{me} and T_{br} delays for 300, 400 and 500 users, respectively. As the graph shows, overall restoration time increases with the number of users.

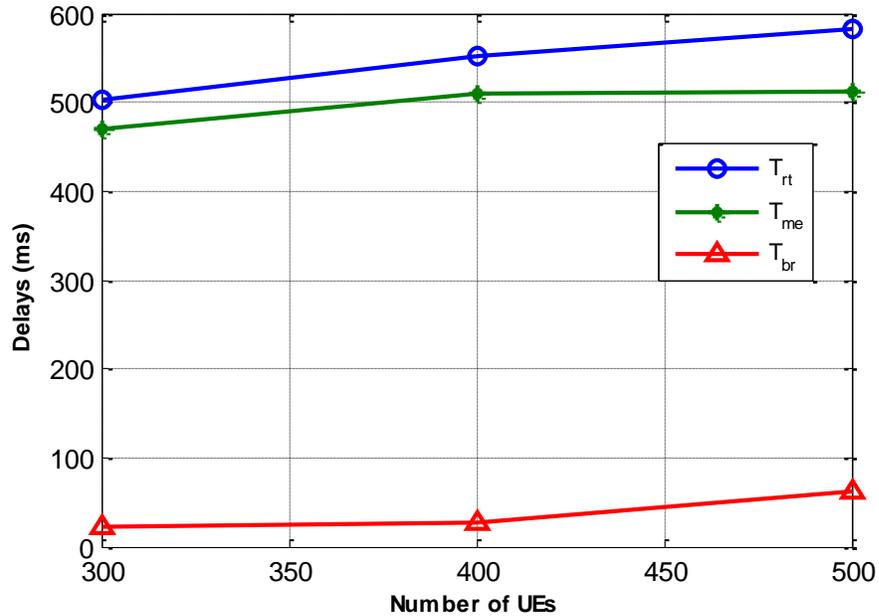


Figure 4-15: Restoration delay comparisons for 1:1 Active-Backup configuration for different numbers of users.

4.3.4 Signaling Message Overhead Estimation

For signaling message overhead calculation in the core network, a network is assumed where an active EPC can handle 100000, 200000 and 300000 active subscribers. Based on the analysis in section 4.2.2, the average signaling overhead associated with the failover procedure against different session arrival rate is shown in Figure 4-16. Average sessions duration is assumed to be 15 minutes and the session arrival rate follows the Poisson distribution. Figure 4-17 shows signaling load for new arrival of UEs during recovery time which is calculated according to the

equation 4. In this case, the average restoration time is assumed to be 500 ms and the session arrival rate follows the Poisson distribution.

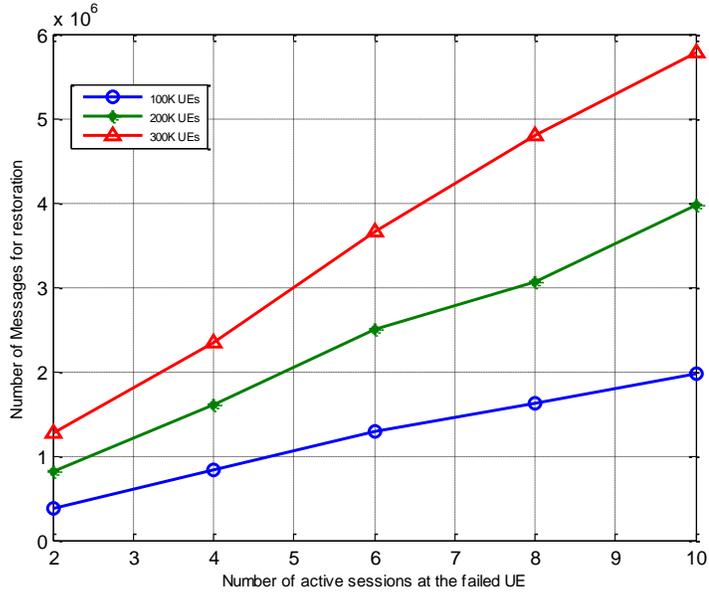


Figure 4-16: Signaling message overhead due to session restoration.

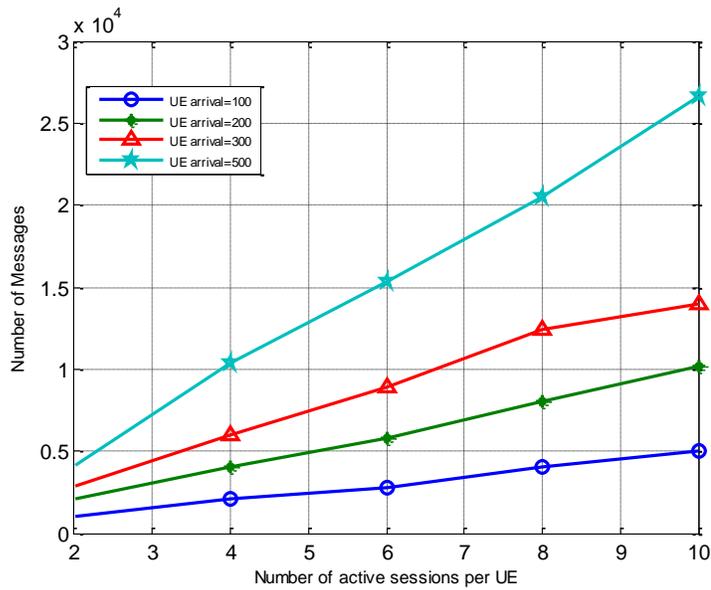


Figure 4-17: Signaling message overhead due to new UE arrivals during restoration time.

Chapter 5

Bandwidth Analysis

Signaling interaction is necessary in order to replicate critical UE context information, failure notification procedures and exchanging status or configuration update of the active and backup MME/S-GWs. In this chapter several logical links are defined which are needed for the aforementioned tasks. The bandwidth requirements for the signaling messages over these links are measured from analytical and numerical perspectives. It should be noted that the proper bandwidth dimensioning will help in determining the backhaul network design for the proposed self-healing system as direct or dedicated links between the network entities may not be available. In fact, links that are referred to in the self-healing architecture are available in the OAM as it is involved in monitoring and administrative tasks of all the network entities. The objective of this chapter is to present the bandwidth analysis incurs by the self-healing system design so that these customized links discussed here can be provisioned over the existing OAM links or LTE standard interfaces like S1 or X2.

5.1 Descriptions of the Logical Connections

In the proposed self-healing system, the following types of logical connections are used for transferring signaling messages and data.

- Replication link: The link between the active MME/S-GW and backup MME/S-GW to replicate the UE context information.
- Configuration update link: The link between the active MME/S-GW and the OAM, as well as the backup MME/S-GW and the OAM to exchange configuration and load status update.

- ENodeB notification link: The links between the OAM and eNodeBs to send the failure notification message. Notification messages can be forwarded over the X2 interface by the informed eNodeBs along with the OAM which improves reliability.
- Core network notification link: The links between the OAM and P-GW to send failure notification messages.

The block diagram in Figure 5-1 shows the links along with the corresponding network entities.

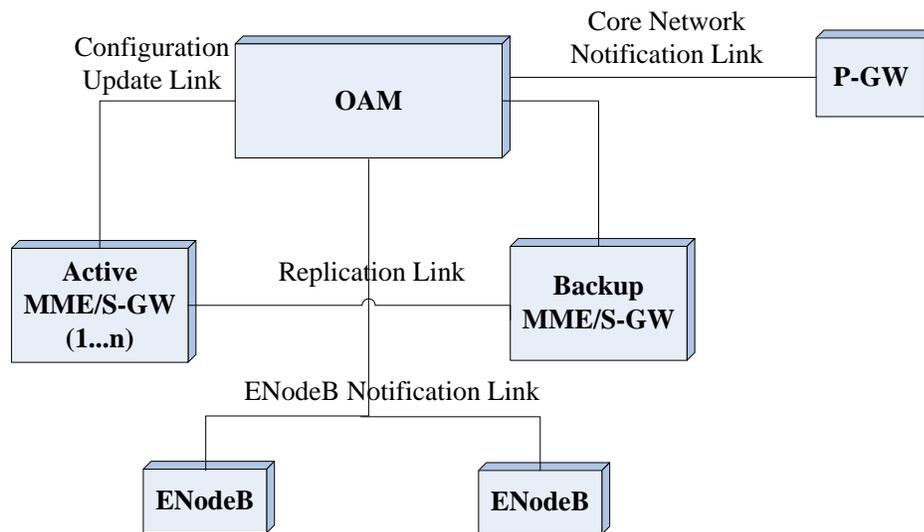


Figure 5-1: Connection links between the network entities for the self-healing system.

5.2 Description of the Signaling Messages

Figure 5-2 shows the signaling message flows over the replication link in order to transfer UE context from the active MME/S-GW to the backup MME/S-GW.

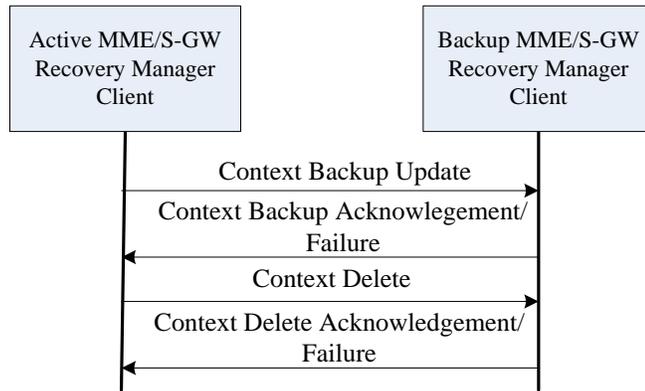


Figure 5-2: Message flows for UE context replication.

An example of signaling message flows over the configuration update link are presented in chapter 3 Figure 3-5.

Signaling messages over the eNodeB notification link are presented in Figure 5-3.

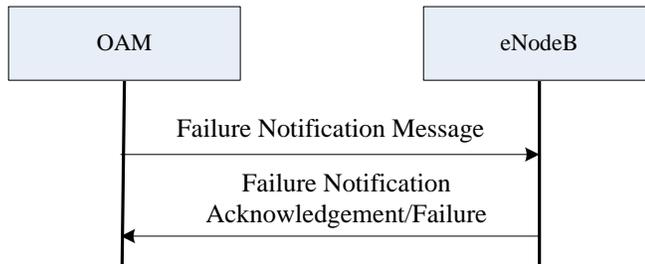


Figure 5-3: Failure notification message flows for the eNodeB.

Figure 5-4 depicts message flows over the core network notification links.

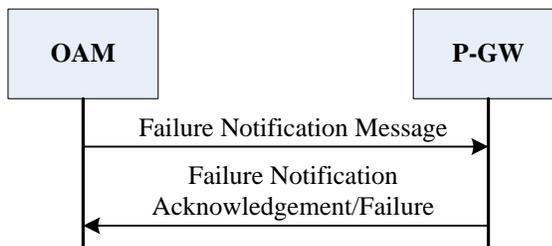


Figure 5-4: Failure notification message flows for the P-GW.

It should be mentioned here that the failure notification messages are invoked when a failure incident occurs in the network. In the event of a failure, number of these messages generated by the OAM depends on the number of network entities supported by a particular MME/S-GW. Generally, a MME/S-GW can supports hundreds of eNodeBs and these nodes are geographically dispersed. So, it can be assumed that the OAM and eNodeBs inter-connected links are not shared. Of course, a group of eNodeBs based on the same location or adjacent areas may be connected through shared link to the OAM. But all the eNodeBs will not be connected through a shared link. The total traffic requirements over the inter-connected links between the OAM and eNodeBs are negligible. Same conclusion can be drawn for the P-GW. As a result, the bandwidth analysis for the eNodeB notification link and the core network notification link are not performed.

In the thesis message flows regarding the configuration update link are implemented partially. For example, periodic heartbeat messages or the configuration or status update messages that the OAM sends to the backup MME/S-GW. It should be clarified at this point that other messages, shown in Figure 3-5 in chapter 3, depend on network configurations, operator's implementation policy. Some example scenarios in this regard have already been discussed in chapter 3. Except the periodic heartbeat message, these messages are event based and infrequently invoked. For example, the OAM exchanges load status information when load of the system exceeds certain predefined threshold in the 1:1 Active-Active configuration. Considering the smaller size of the heartbeat message, the overall bandwidth requirement of this link is also negligible.

However, the messages over the replication link are very frequently exchanged. It should be noted that the context information will be copied from the active MME/S-GW to the backup MME/S-GW in two cases. First, whenever a UE attaches to the network and the default bearer is created and second, whenever the UE creates an additional bearer. Similarly, when the UE

disconnects from the network or terminates a bearer, the active MME/S-GW sends context delete related messages to the backup MME/S-GW. In this way, it can be ensured that the backup MME/S-GW contains most recent context data of a subscriber which facilitate proper recovery of user sessions after the failure incident.

Table 5-1 shows description of the messages over the replication link. Information Elements (IEs) and average size of each of the messages are also mentioned.

Table 5-1: Context replication message description.

Message Type	Information Elements	Average size
Context Backup Update	<i>Message ID, IMSI, GUMMEI, UE security context, UE network capability, Selected CN operator ID, APN, APN restriction, APN-AMBR, Serving GW addresses, TEID for control signalling, TEIDs for uplink traffic, PDN GW addresses, GTP-based S5/S8 TEIDs, APN, Serving GW addresses, TI.</i>	200 Bytes [54]
Context Backup Acknowledgment	<i>Message ID, GUMMEI, ME Identity, IMSI</i>	10 Bytes
Context Backup Failure	<i>Message ID, GUMMEI, IMSI, Time to wait, Release Cause</i>	12 Bytes

Context Delete	<i>Message ID, GUMMEI, ME Identity, IMSI</i>	10 Bytes
Context Delete Acknowledgement	<i>Message ID, GUMMEI, ME Identity, IMSI</i>	10 Bytes
Context Delete Failure	<i>Message ID, GUMMEI, IMSI, Time to wait, Release Cause</i>	12 Bytes

5.2.1 Definitions of Information Elements (IEs)

Message ID: Message ID identifies a message.

IMSI: International Mobile Subscriber Identity (IMSI) is unique number that identifies a UE in GSM, UMTS or LTE network.

GUMMEI: Global unique MME identity (GUMMEI) uniquely identifies a MME in LTE network.

UE security context: UE security context contains temporary UE identity, ciphering and integrity key used for security and protection.

UE network capability: This parameter contains general UE characteristics. It provides network information regarding various aspects of UE related to EPS. Based on this information, network may handle UE operation in different ways.

APN: Access point name (APN) determines an IP packet data network (PDN) for user sessions. It also defines the service types that can be provided to the user in a session.

APN restriction: APN restriction parameter is used to verify whether a UE is allowed to establish bearers with other APN's in the network.

APN-AMBR: APN-Aggregate maximum bit rate (AMBR) is used to determine the aggregate bit rate that can be provided to the user over all non-GBR bearers established through a particular APN.

Selected CN operator ID: This parameter identifies core network operator.

Serving GW addresses: The S-GW data plane and control plane addresses.

PDN GW addresses: The P-GW data plane addresses.

TEID: An EPS bearer is identified on different interfaces using the GTP Tunneling End ID (TEID) and corresponding IP addresses.

Time to wait: This parameter defines minimum allowed waiting time.

Release Cause: Release cause is used to identify the unexpected events during signaling.

5.3 Analytical Model for Bandwidth Calculation

In this section, the bandwidth requirement for the replication link is estimated. There are four messages exchanged by the active and backup MME/S-GW in a successful operation: *Context Backup Update*, *Context Backup Acknowledgment*, *Context Delete* and *Context Delete Acknowledgement*.

It is assumed that the context replication operation is successfully completed and no failure arises during signalling exchange. As a result, *Context Backup Failure* and *Context Delete Failure* messages are not included in the calculation. As mentioned in Table 5-1, the payload of *Context Backup Update* which contains UE context data is estimated to be on the order of 200 bytes [54] and the payload of the other three messages are on average 10 bytes, excluding lower layers overheads. Therefore, average message size for the context backup is $[(Context\ Backup\ Update + Context\ Backup\ Acknowledgment)/2]$ or $[(200+10)/2]$ or 105 bytes and average message size for

the context delete is $[(Context\ Backup\ Delete + Context\ Delete\ Acknowledgement)/2] [(10+10)/2]$ or 10 bytes.

Following parameters are considered to formulate the equation for bandwidth calculation:

- BW_R : Average required bandwidth for UE context replication.
- λ_{AUE} : This parameter denotes the average UE arrival rate in the network. We can assume Poisson distribution to model it [31].
- λ_{TUE} : The average UE disconnect rate from the network. It can also be modeled with Poisson distribution.
- N_B : Average number of bearers created by a UE.
- L_{CB} : Average length of context backup related messages.
- L_{CD} : Average length of context delete related messages.
- N_{CB} : Number of messages involved in the backup of context information.
- N_{CD} : Number of messages involved in the deletion of context information.

Thus, average required bandwidth is calculated using following equation,

$$BW_R = N_{CB} * \lambda_{AUE} (1+N_B) * L_{CB} + N_{CD} * \lambda_{TUE} (1+N_B) * L_{CD} \quad (5)$$

As discussed in chapter 2, the default bearer is always established when a UE is connected to the network. The default bearer is non-GBR type. The UE can request for a dedicated bearer which can be a GBR or non-GBR depending on the requirements of the application and subscription status of the user. The number of bearers - default or dedicated that can be opened by a UE is limited. A UE can have a maximum of 11 EPS bearers [56]. In order to measure N_B or the average number of bearers created by a UE, we are depending on the user's application usage characteristics. For instance, subscribers who use basic services such as e-mail, text messaging or web browsing do not require specific QoS requirements or enhanced user experience. For this

kind of services, the GBR bearer is not needed. That means the default bearer alone is sufficient to provide these services. On the contrary, multimedia-rich services such as online games, video conferencing demand stringent QoS that may require the dedicated bearers. The time and frequency of usage is also important in measuring N_B . For example, concurrent QoS based session may require multiple dedicated bearers.

In order to decide the average application usage behavior, subscribers have been divided into three categories: *i)* light use subscribers, *ii)* medium use subscribers, and *iii)* heavy use subscribers. The light users infrequently use the basic services and these services are less QoS constrained which can be served by the default bearer. The medium users more frequently use the QoS based services along with the basic services. The heavy users frequently and concurrently use the QoS based services and the basic services. Table 5-2 shows services used by the subscribers of each profiles.

Table 5-2: Subscribers profile information.

Subscriber type	Services	Frequency of dedicated bearer usage	Average number of bearers
Light use	E-mail, FTP, HTTP, SMS	less frequent	1 (default)
Medium use	VOIP, VPN	more frequent	2 (default and dedicated)
Heavy use	online games, video streaming	very frequent	6 (default and dedicated)

Based on the assumption in Table 5-2, there can be subscribers of different profiles in the average UE arrival rate. That means the average number of bearers opened by those subscribers varies. Three different scenarios are presented in Table 5-3 for different mixture of user profiles.

Table 5-3: Different user profile in UE arrival rate.

Category	Light Use	Medium Use	Heavy Use
Light Mix	70%	20%	10%
Balanced Mix	50%	30%	20%
Heavy Mix	35%	35%	30%

5.4 Numerical Results for Bandwidth Calculation

This section presents some numerical results on the bandwidth analysis based on the proposed methods for various UE arrival rates and profiles. Figure 5-5 shows the bandwidth usage at the EPC against different UE arrival rate for different subscribers mix. It is noted that the presence of higher percentage of heavy users increases the bandwidth requirements significantly.

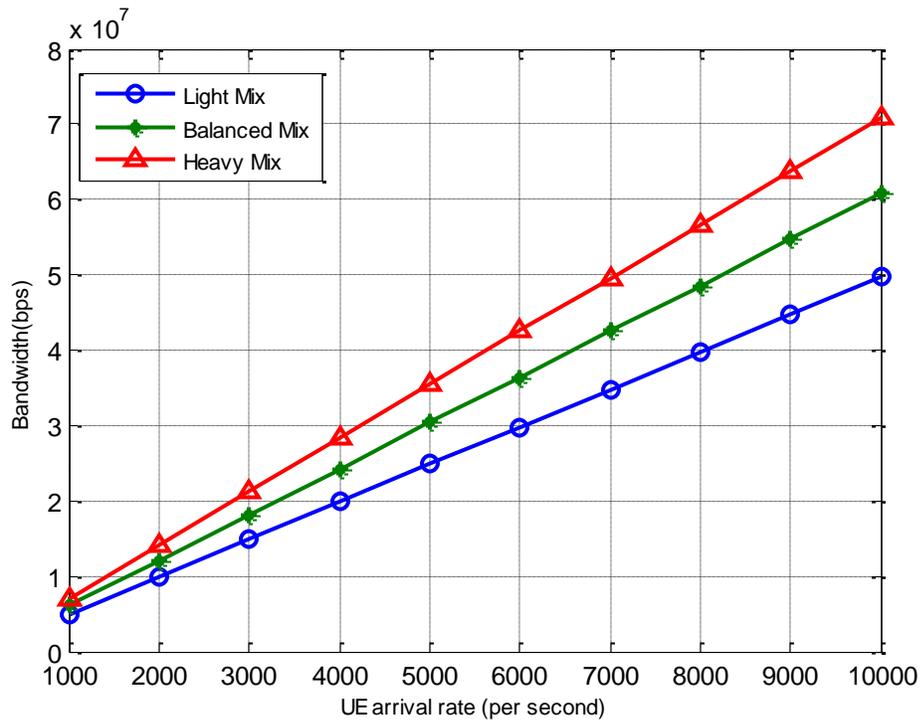


Figure 5-5: Bandwidth usage due to UE context replication.

Chapter 6

Conclusion and Future Works

With increasing reliance on mobile networks, service or infrastructure failures have severe economic impact on the operators and customers. As a result, more attention has been paid to the issues of service survivability and the failure recovery in the cellular networks. But failure recovery with the strategy of service continuity is complex in telecommunication world and demands stringent requirements and costly implementation plan. Motivated by the SON functionalities of LTE, this thesis presents a self-healing approach for the control and data plane network elements of LTE EPC – the MME and S-GW failures. The assumption about the nature of failure is that the MME/S-GW does not fail gracefully. In order to ensure seamless service continuity, the self-healing scheme discussed two fault-tolerant EPC architectures, the centralized active-backup configurations and distributed active-active configurations. Cost-performance trade-offs can be an important criteria to the operators for selecting an appropriate architecture.

In order to demonstrate the performance, the self-healing schemes were modeled for the active-backup and active-active scenarios using the centralized approach. Simulations were conducted to investigate various factors in terms of restoration time, message overhead, bandwidth usage and different kinds of delays (e.g., EPS bearer delay) involved in two architectural configurations with different failure scenarios and subscribers base. The simulation results showed that the proposed approach performs reasonably well in terms of recovering services if resources are properly provisioned. We configured congestion or overload in the backhaul network in the simulation which results in significant bearer re-creation delay and caused performance degradation. The obtained results from other scenarios indicate that the dedicated backup

configuration (N:1 active-backup configuration) performs better than the active-active configuration in terms of restoration delay.

An analysis was presented on the signaling load at the EPC due to simultaneous re-attachment of the failed UEs during the recovery process. Numerical results show that number of sessions or bearers in a UE and corresponding bearer re-creation process have significant impact in terms of signaling congestion.

Moreover, operational procedures of the proposed schemes incur signaling interactions between different network elements. The bandwidth requirements for the signaling are analyzed to understand their potential implication for the backhaul network design. The UE context replication process between the active and backup MME/S-GWs requires more bandwidth than other procedures. To estimate the bandwidth requirements in this regard, it is necessary to determine the frequency of the bearer usage by the subscribers. In order to facilitate that, subscriber's base was divided into several user profiles based on the frequency of usage. Simulation results show that the bandwidth requirements increase significantly with the presence of subscribers who frequently uses dedicated bearers.

As future plan, it is suggested to study the self-healing solutions in the LTE-A network where the EPC supports large number of relay nodes which are used to increase the capacity and coverage of the network. Another extension is to design efficient self-healing scheme for the heterogeneous access networks supported by the EPC, for example, WiMAX, WLAN etc.

References

- [1] http://www.gsacom.com/downloads/pdf/GSA_evolution_to_lte_report_310811.php4, retrieved Sept 1, 2011.
- [2] 3GPP TS 36.300 V8.12.0: 3rd Generation Partnership Project; Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2.
- [3] A. Ghosh, R. Ratasuk, B. Mondal, N. Managalvedhe, and T. Thomas, "LTE-advanced: next-generation wireless broadband technology," *IEEE Wireless Communications*, vol. 17, no. 3, pp. 10-22, 2010.
- [4] 3GPP TR 23.882 V8.0.0: 3rd Generation Partnership Project; 3GPP System Architecture Evolution: Report on Technical Options and Conclusions.
- [5] www.3gpp.org
- [6] ITU Press Release, http://www.itu.int/net/pressoffice/press_releases/2011/01.aspx, retrieved Jan 6, 2012.
- [7] "Server Glitch Crashes T-Mobile Network", <http://www.cn-c114.net/576/a406617.html>, retrieved Sept 6, 2011.
- [8] M. Rahnema, "Overview of the GSM system and protocol architecture," *IEEE Commun. Mag.*, vol. 31, no. 4, pp. 92-100, 1993.
- [9] C. Bettstetter, H. Vogel and J. Eberspacher, "GSM phase 2+ general packet radio service GPRS: Architecture, protocols, and air interface," *IEEE, Communications Surveys and Tutorials*, vol. 2, no. 3, pp. 2-14, 1999.
- [10] Agilent Technologies, "Understanding General Packet Radio Service," <http://cp.literature.agilent.com/litweb/pdf/5988-2598EN.pdf>, 2001.
- [11] K. W. Richardson, "UMTS overview," *Electronics & Communication Engineering Journal*, vol. 12, no. 3, pp. 93-100, 2000.
- [12] Nokia White Paper, "HSDPA solution," http://www.nokia.com/NOKIA_COM_1/About_Nokia/Press/White_Papers/pdf_files/HSDPA_A4.pdf, 2003.
- [13] E. Dahlman et al., *3G Evolution: HSPA and LTE for Mobile Broadband*, 2nd ed., Academic Press, 2008.

- [14] Qualcomm White Paper, “HSPA+ for Enhanced Mobile Broadband”,
http://www.qualcomm.com/common/documents/white_papers/HSPAPlus_MobileBroadband_021309.pdf , 2009.
- [15] 3GPP TR 25.913 V8.0.0: 3rd Generation Partnership Project; Technical Specification Group Radio Access Network; Requirements for Evolved UTRA (E-UTRA) and Evolved UTRAN (E-UTRAN).
- [16] D. Astely et al., "LTE: The Evolution of Mobile Broadband", IEEE Communication Magazine, vol. 47, no. 4, pp. 44-51, 2009.
- [17] P. E. Mogensen et al., “LTE-Advanced: The Path towards Gigabit/s in Wireless Mobile Communications,” 1st International Conference on Wireless VITAE 2009, Aalborg, pp. 147-151, July 2009.
- [18] C. Makaya et al., “Service continuity support in self-organizing IMS networks,” 2nd International Conference on Wireless Communication, VITAE, pp. 1 –5, 2011.
- [19] 3GPP TS 32.500 V10.0.0: Telecommunication management; Self-organizing Networks (SON); Concepts and requirements.
- [20] 3GPP TS 32.541 V10.0.0: Telecommunication management; Self-organizing Networks (SON); Self-healing Concepts and Requirements.
- [21] D. Ghosh, R. Sharman, H. Raghav Rao, and S. Upadhyaya, "Self-healing systems - survey and synthesis", Decision Support Systems, vol. 42, no. 4, pp. 2164–2185, 2007.
- [22] A.G. Ganak and T.A. Corbi, “The Dawning of the Autonomic Computing Era, IBM Systems J., vol. 42, no. 1, pp. 5-18, 2003.
- [23] L. Kant, "Design and Performance Modeling & Simulation of Self-healing Mechanisms for Wireless Communication Networks", In Proceedings of 35th Annual Simulation Symposium, 2002.
- [24] A. Snow, U. Varshney, and A. Malloy, “Reliability and Survivability of Wireless and Mobile Networks,” IEEE Comp., vol. 33, no. 7, pp. 49–55, 2000.
- [25] D. Tipper, T. Dahlberg, H. Shin, C. Charnsripinyo, “Providing fault tolerance in wireless access networks”, IEEE Communications Magazine, vol. 40, no.1, pp. 58-64, 2002.
- [26] S. Tsao, “Scalable gateway GPRS support node for GPRS/UMTS networks”, 56th IEEE Vehicular Technology Conference (VTC Fall’02), Volume 4, pp. 2239, 2002.

- [27] S. Kustos, L. Bokor, G. Jeney, "Testbed Evaluation of Dynamic GGSN Load Balancing for High Bitrate 3G/UMTS Networks", Vehicular Technology Conference (VTC Spring), pp. 1-5, 2011.
- [28] L. Schmelz et al., "Framework for the development of self-organization methods", INFSO-ICT-216284 SOCRATES D2.4, 2008.
- [29] 3GPP TS 23.007 V9.7.0: 3rd Generation Partnership Project; Technical Specification Group Core Network and Terminals; Restoration Procedures.
- [30] 3GPP TR 23.857 V1.3.0: 3rd Generation Partnership Project; Technical Specification Group Core Network and Terminals; Study of EPC Nodes Restoration.
- [31] T. Taleb and K. Samdanis, "Ensuring Service Resilience in the EPS:MME Failure Restoration Case", IEEE GLOBECOM, Dec 5-9, 2011.
- [32] <http://www.opnet.com>
- [33] <http://www.mathworks.com/products/matlab>
- [34] M. Sauter, Beyond 3G - Bringing Networks, Terminals and the Web Together: LTE, WiMAX, IMS, 4G Devices and the Mobile Web 2.0, John Wiley & Sons, 2008.
- [35] M. Olsson et al., SAE and the Evolved Packet Core: Driving The Mobile Broadband Revolution, Academic Press, 2009.
- [36] The LTE Network Architecture, Alcatel-Lucent White Paper, 2009.
- [37] 3GPP TS 22.278 V8.10.0: 3rd Generation Partnership Project; Service requirements for the Evolved Packet System (EPS).
- [38] 3GPP TS 23.401 V10.2.0: General Packet Radio Service (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access.
- [39] R. Nossenson, "Long-Term Evolution Network Architecture", IEEE International Conference on COMCAS, pp. 1-4, 2009.
- [40] H. Holma et al., LTE for UMTS - OFDMA and SC-FDMA Based Radio Access, 1st ed., John Wiley & Sons, 2009.
- [41] 3GPP TS 24.301 V8.1.0: Non-Access-Stratum (NAS) protocol for Evolved Packet System (EPS).
- [42] Continuous Computing White Paper, "Unlocking Long Term Evolution - A Protocol Perspective", 2010.

- [43] 3GPP TS 36.413 V8.9.0; Evolved Universal Terrestrial Radio Access Network (E-UTRAN); S1 Application Protocol (S1AP).
- [44] 3GPP TS 36.423 V8.3.0 (2008-09); X2 Application Protocol (X2AP) specification.
- [45] P. Lescuyer and T. Lucidarme, Evolved Packet System (EPS): The LTE and SAE Information of 3G UMTS, J. Wiley & Sons, 2008.
- [46] H. Ekstrom et al., "QoS Control in the 3GPP Evolved Packet System," IEEE Communication Magazine, vol. 47, no. 2, pp. 76–83, 2009.
- [47] NEC White Paper, "Self Organizing Network," <http://www.nec.com/en/global/solutions/nsp/lte/pdf/son.pdf>, 2009.
- [48] 3G Americas: "The Benefits of SON in LTE", http://www.3gamericas.org/documents/2009_%203GA_LTE_SON_white_paper_12_15_09_Final.pdf, 2009.
- [49] INFSO-ICT-216284 SOCRATES D2.1: Use Cases for Self-Organising Networks, [http://www.fp7-socrates.eu/files/Deliverables/SOCRATES_D2.1 Use cases for self-organising networks.pdf](http://www.fp7-socrates.eu/files/Deliverables/SOCRATES_D2.1%20Use%20cases%20for%20self-organising%20networks.pdf),
- [50] 3GPP TR 36.902 V9.2.0: 3rd Generation Partnership Project; Self-configuring and self-optimizing network (SON) use cases and solutions.
- [51] Nomor Research, White Paper, "Self-Organizing Networks (SON) in 3GPP Long Term Evolution", 2008.
- [52] A. Bianco, J. Finochietto, L. Giraud, M. Modesti and F. Neri, "Network planning for disaster recovery ", 16th IEEE Workshop on LANMAN, September, 2008.
- [53] I. Widjaja, P. Bosch, H. La Roshe, "Comparison of MME Signaling Loads for Long-Term-Evolution Architectures", IEEE 70th VTC, Sept 20-23, 2009.
- [54] I. Widjaja and H. La Roshe, "Sizing X2 Bandwidth for Inter-connected eNBs", IEEE 70th VTC, Sept 20-23, 2009.
- [55] <http://3g4g.blogspot.ca/2009/06/lte-qci-and-end-to-end-bearer-qos-in.html>, retrieved Dec 12, 2011.
- [56] C. Cox, An Introduction to LTE: LTE, LTE-Advanced, SAE and 4G Mobile Communications, 1st ed., John Wiley & Sons, 2012.