

GTI IoT Core Network Architecture White Paper

GTI

<http://www.gtigroup.org>

GTI

| | |
|---------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Version: | V1.0 |
| Deliverable Type | <input checked="" type="checkbox"/> Procedural Document <input type="checkbox"/> Working Document |
| Confidential Level | <input checked="" type="checkbox"/> Open to GTI Operator Members <input checked="" type="checkbox"/> Open to GTI Partners <input checked="" type="checkbox"/> Open to Public |
| Working Group | IoT Program |
| Source members | China Mobile, Nokia, Ericsson, Huawei, ZTE |
| Support members | China Mobile, Nokia, Ericsson, Huawei, ZTE |
| Editor | China Mobile, Nokia, Ericsson, Huawei, ZTE |
| Last Edit Date | 08-04-2018 |
| Approval Date | DD-MM-YYYY |

Confidentiality: This document may contain information that is confidential and access to this document is restricted to the persons listed in the Confidential Level. This document may not be used, disclosed or reproduced, in whole or in part, without the prior written authorization of GTI, and those so authorized may only use this document for the purpose consistent with the authorization. GTI disclaims any liability for the accuracy or completeness or timeliness of the information contained in this document. The information contained in this document may be subject to change without prior notice.

Document History

| Date | Meeting # | Version # | Revision Contents |
|------------|-----------|-----------|-------------------------------------------|
| 10-30-2017 | | 0.1 | Initial Draft |
| 12-20-2018 | | 0.2 | Updates from Nokia, Huawei, Ericsson, ZTE |
| 03-14-2018 | | 0.3 | Updated TOC and formatting |
| 05-10-2018 | | 1.0 | Final Edit |
| | | | |
| | | | |

Content

| | |
|---------------------------------------------------------------------|-----------|
| EXECUTIVE SUMMARY / INTRODUCTION | 6 |
| 1. TERMINOLOGY | 7 |
| 2. ASSUMPTIONS AND SCOPE | 7 |
| 3. CELLULAR IOT SERVICES OVERVIEW | 7 |
| 3.1. CHALLENGES TO CELLULAR IOT | 7 |
| 3.2. PROGRESS IN 3GPP SPECIFICATIONS AND INDUSTRY DEVELOPMENT | 8 |
| 3.3. CELLULAR IOT SERVICES USE CASES | 9 |
| 3.3.1 <i>Use Case Overview</i> | 9 |
| 3.3.2 <i>User Case 1: Metering</i> | 10 |
| 3.3.3 <i>User Case 2: Asset monitoring</i> | 10 |
| 4. CELLULAR IOT CORE NETWORK ARCHITECTURAL REQUIREMENTS | 10 |
| 4.1. REQUIREMENTS FROM SERVICE ASPECT | 10 |
| 4.2. REQUIREMENTS FROM NETWORK OPERATION ASPECT | 10 |
| 4.2.1 <i>Network Flexibility</i> | 10 |
| 4.2.2 <i>Network Scalability</i> | 11 |
| 4.2.3 <i>Network Capability Exposure</i> | 11 |
| 4.2.4 <i>Automated network management</i> | 11 |
| 5. KEY TECHNOLOGIES IN CIOT CORE NETWORK ARCHITECTURE DESIGN | 11 |
| 5.1. OVERVIEW | 11 |
| 5.2. NETWORK FUNCTION VIRTUALIZATION | 11 |
| 5.3. CORE NETWORK ARCHITECTURE OPTIMIZATION | 12 |
| 5.4. SERVICE CAPABILITY EXPOSURE | 12 |
| 6. NFV BASED NETWORK ARCHITECTURE AND MANAGEMENT | 12 |
| 6.1. NFV OVERVIEW | 12 |
| 6.2. NFV ARCHITECTURE OVERVIEW | 13 |
| 6.3. NFV MANAGEMENT ARCHITECTURE OVERVIEW | 14 |
| 6.4. VNF LIFECYCLE MANAGEMENT | 14 |
| 6.5. NETWORK SERVICE MANAGEMENT AND ORCHESTRATION | 15 |
| 7. CIOT CORE NETWORK ARCHITECTURE AND SOLUTIONS | 15 |
| 7.1. CIOT OPTIMIZATION EPC ARCHITECTURE (NB-IOT AND EMTC) | 15 |
| 7.2.1 <i>CIOT Optimization EPC architecture</i> | 15 |
| 7.2.2 <i>Network Functions Enhancements</i> | 16 |
| 7.2.3 <i>Reference Points</i> | 16 |
| 7.2. KEY SOLUTION IN CIOT CORE NETWORK | 17 |
| 7.2.1 <i>Efficient data transmission</i> | 17 |
| 7.2.1.1 <i>Data over NAS (Control-Plane Optimization)</i> | 17 |
| 7.2.1.2 <i>User-Plane Optimization</i> | 18 |

| | |
|-----------------------------------------------------------|-----------|
| 7.2.1.3 Non-IP small data transmission | 18 |
| 7.2.2. <i>Low Power Consumption of Terminals</i> | 19 |
| 7.2.2.1 PSM | 19 |
| 7.2.2.2 Extended Periodic TAU Timer | 19 |
| 7.2.2.3 eDRX | 20 |
| 7.2.3. <i>Coverage Enhancement</i> | 20 |
| 7.2.4. <i>Massive Access Control</i> | 20 |
| 7.2.4.1 Congestion Control Management | 20 |
| 7.2.4.2 Rate Control | 21 |
| 7.2.5. <i>Inter RAT idle mode mobility to/from NB-IoT</i> | 22 |
| 7.2.6. <i>DECOR and eDECOR</i> | 22 |
| 7.2.7. <i>SMS</i> | 23 |
| 8. CIOT CORE NETWORK DEPLOYMENT ANALYSIS | 23 |
| 8.1. SHARING THE CORE NETWORK WITH CONVERGED MBB NETWORK | 23 |
| 8.2. CIOT DEDICATED CORE NETWORK FOR CIOT | 24 |
| 9. CIOT SERVICE CAPABILITY EXPOSURE | 24 |
| 9.1. OVERVIEW | 24 |
| 9.2. SERVICE CAPABILITY EXPOSURE REQUIREMENTS | 25 |
| 9.2.1. <i>Real-time user status awareness</i> | 25 |
| 9.2.2. <i>Communications</i> | 25 |
| 9.2.3. <i>Service parameter management</i> | 26 |
| 9.2.4. <i>Qos management</i> | 28 |
| 9.3. SERVICE CAPABILITY EXPOSURE ARCHITECTURE | 28 |
| 9.3.1. <i>Overview</i> | 28 |
| 9.3.2. <i>Interface</i> | 31 |
| 9.3.3. <i>Identifier</i> | 32 |
| 9.3.4. <i>Deployment and routing</i> | 33 |
| 9.4. SECURITY REQUIREMENTS OF SERVICE CAPABILITY EXPOSURE | 33 |
| 10. RESOURCES/REFERENCES | 34 |

Executive Summary / Introduction

The Internet of Things (IoT) is the next revolution in the mobile ecosystem. In 2025, an estimated 27 billion connected devices will be deployed, up from 5 billion in 2015.

Among them, cellular IoT modules are forecast to account for 2.2 billion, up from 0.334 billion in 2015. [Machina Research, May 2015].

This document gives a description on cellular IoT core network architectural requirement and existing architecture specification by SDOs. This whitepaper also investigates how to leverage NFV and other emerging network technologies to realize optimization of cellular IoT core network architecture.

1. Terminology

| Term | Description |
|--------|--------------------------------------|
| CIoT | Cellular IoT |
| CP | Control Plane |
| DCN | Dedicated Core Network |
| eDRX | extended Discontinuous Reception |
| eMTC | enhanced Machine Type Communication |
| EPC | Evolved Packet Core |
| EPS | Evolved Packet System |
| GW | Gateway |
| IoT | Internet of Things |
| LPWA | Low-Power Wide-Area |
| MME | Mobility Management Entity |
| MNO | Mobile Network Operator |
| MTC | Machine Type Communications |
| NAS | Non-Access Stratum |
| NB-IoT | Narrowband IoT |
| PO | paging opportunity |
| PSM | Power Saving Mode |
| SCEF | Service Capability Exposure Function |
| SMS | Short message service |
| UE | User Equipment |
| UP | User Plane |

2. Assumptions and Scope

This whitepaper provides an overview of mobile core network architecture which is optimized to support varied IoT services, with leveraging NFV technique. From the perspective of massive IoT services, this document describes IoT service requirements, architectural requirement, EPC optimized architecture to support NB-IoT or enhanced MTC devices, and service capability exposure architecture.

3. Cellular IoT Services Overview

3.1. Challenges to Cellular IoT

Internet of Things (IoT) is the network in which physical objects exchange information with the Internet. An IoT terminal can communicate with each other through the Internet. With its wide applications in industry and human life, IoT services require different transmission rates, as shown in the following figure.

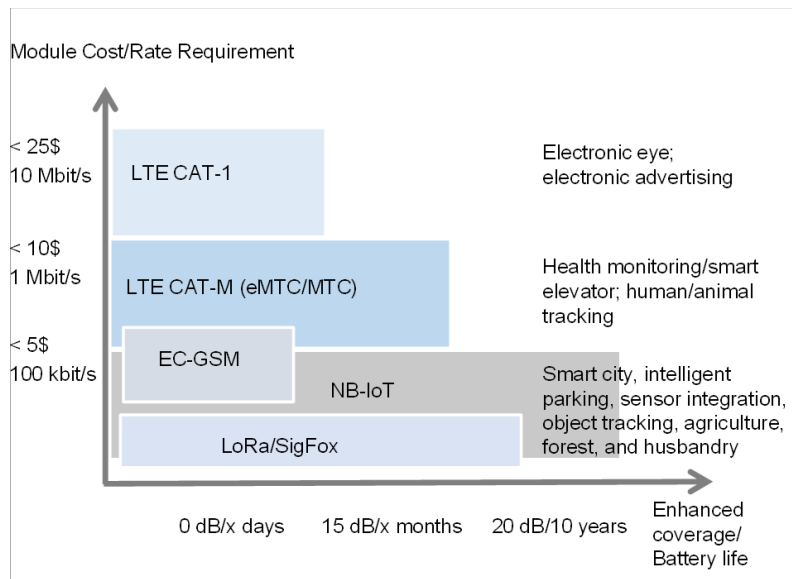


Figure 3.1-1 IoT technology comparison

- High-end applications with a high speed, such as video surveillance and electronic advertising. Currently, these applications can be implemented on the UMTS/LTE network of carriers, and there are not any non-3GPP technologies posing a threat to the 3GPP technology in these applications.
- Mid-range applications with a medium speed, such as smart household, point of sale (POS) machines. These applications consume low power and can be implemented on the enhanced Machine Type Communication (eMTC), which requires a data rate of less than 1 Mbit/s.
- Low-end applications with a low speed, such as intelligent meter reading and low-end vehicle-mounted devices, which are estimated to account for 70% of the IoT market. Non-3GPP technologies such as LoRa and SigFox have been deployed on unlicensed spectrum. To compete with these non-3GPP technologies, carriers and vendors jointly propose a brand-new air interface technology NB-IoT, which requires a data rate of lower than 100 kbit/s, meeting the requirements of low cost, low power consumption, and wide coverage.

In the future, there may be extra-low delay applications such as self-driving and LTE-V that is being researched by 3GPP. These applications are currently not included in this document.

3.2. Progress in 3GPP Specifications and Industry Development

In 3GPP RAN#69 Conference held in September 2015, the Narrowband Internet of Things (NB-IoT) project was successfully initiated, which was expected to compete in the low-speed IoT market. NB-IoT is a type of Low Power Wide Area (LPWA) IoT technology typically applied in intelligent meter reading, intelligent parking, smart environmental monitoring, and intelligent agriculture. The NB-IoT network is required to provide:

- Enhanced coverage: coverage with a 20 dB gain when compared with GPRS coverage
- Low power consumption: increased battery lifetime for more than 10 years
- Low cost: chip cost of less than 1 USD and module of less than 5 USD in mass production

- Massive connection: supporting up to 50,000 users per cell

Note: The battery life, connection quantity is estimated based on the traffic model from the 3GPP TR 45.820.

In 3GPP RAN#72 Conference held in Pusan, South Korea on June 16, 2016, NB-IoT was added to 3GPP R13 specifications as an important subject. The key NB-IoT standards widely supported by the wireless industry have been finalized over the 2-year period. NB-IoT adopts a brand-new air interface technology. IoT can also be supported with the evolution of LTE. To support IoT applications, LTE can be evolved to implement LTE-Machine-To-Machine communications. In 3GPP R12 specifications, machine type communication (MTC) is presented, with CAT-0 terminals whose uplink or downlink data rate is 1 Mbit/s. In 3GPP R13 specifications, eMTC is presented, with CAT-M1 terminals. The software features of the IoT network are designed according to eMTC technologies. The coverage is enhanced with a 15 dB gain, and the terminal cost is further reduced. This meets the wide-coverage and low-cost requirements of low-speed IoT.

Both NB-IoT and eMTC are implemented based on the S1 link. They are different in data rates as well as the following core network characteristics:

- In 3GPP specifications, an enumeration value "NB-IoT" is added to the RAT type defined for NB-IoT but not for eMTC. Therefore, UEs accessing the network through the NB-IoT RAN can be identified upon their access.
- NB-IoT does not support voice services. eMTC supports all LTE features as long as terminals supporting eMTC are used.
- Though enhanced Cellular Internet of Things (CIoT) features defined in 3GPP R13 specifications can be applied in both NB-IoT and eMTC, their applications are different. (The features include control-plane transmission optimization, user-plane transmission optimization, eDRX, PSM, non-IP data delivery, and service capability exposure.) For example, control-plane transmission optimization is mandatory in NB-IoT but is optional in eMTC.

3.3. Cellular IoT Services Use Cases

3.3.1 Use Case Overview

The LPWA market has existed for about 10 years; it's not a new thing. The current technologies (solutions) supporting this market are fragmented and non-standardized, therefore there are shortcomings like poor reliability, poor security, high operational and maintenance costs.

Furthermore, the new overlay network deployment is complex. CIoT overcomes the above defects, with all the advantages like wide area ubiquitous coverage, fast upgrade of existing network, low-power consumption guaranteeing 10 year battery life, high coupling, low cost terminal, plug and play, high reliability and high carrier-class network security. Initial network investment may be quite substantial and superimposed costs are very little. CIoT perfectly matches LPWA market requirements, enabling operators to enter this new field. CIoT enables operators to operate traditional businesses such as Smart Metering, Tracking, by virtue of ultra-low-cost (\$5 to \$10) modules and super connectivity (100K / Cell), also opens up more industry opportunities, for example, Smart City, eHealth. CIoT makes it possible for more things to be connected, but also managing the commercial value of the resulting Big Data is a big task,

operators can carry out cooperation with related industries, in addition to selling connections, they can also sell data.

3.3.2. User Case 1: Metering

Smart metering helps save manpower by remotely collecting electricity, water and gas meter data over the cellular network. This is gaining quite an amount of momentum with most of the top MNOs taking an interest in this topic mainly due to the market opportunity it presents. Smart metering will consequently help cut down cost generated from manual meter reading and changing of meter batteries, which seems to be the two major cost drivers for conventional metering. Smart metering includes smart meters for water, gas and electricity.

3.3.3. User Case 2: Asset monitoring

Asset tracking mainly deals with monitoring methods of physical assets made possible by a module on the asset broadcasting its location. Assets are usually tracked using GPS technology. This service is best leveraged in the logistics and transportation management industry, where using sensors in modules sending information over the cellular network it is possible to gather and manage data relating to the current geographical location of assets. Asset tracking helps the owners of the assets to detect and preemptively react to unexpected events.

4. Cellular IoT Core Network Architectural Requirements

4.1. Requirements from Service aspect

The energy-efficient and low-complexity IoT devices (e.g. smart energy distribution grid system or agriculture IoT devices) are expected being connected through cellular network with following capabilities expectations from the mobile network:

- Support power saving mode of IoT devices for less power consumption, if IoT device has equipped with power control functions.
- Support simplified connection establish procedure between IoT device and network.
- Support delay-tolerant capability with long term caching mechanism.

4.2. Requirements from Network Operation aspect

4.2.1. Network Flexibility

The CloT core network architecture is expected to be flexible enough to support variety of IoT devices and applications, to approach the efficiency of resource utilization in core network when handling various traffic models, it is expected CloT core network to support following capabilities:

- Network Function Virtualization (NFV) is a significant technology which can make the CloT network more flexible by realizing network functions as pure software components. Also, NFV is enabler to approach the programmability of network functions when provisioning of new emerging services.
- Separation of control planes and data planes can further enable the flexible deployment and dimensioning of control plane and user plane components independently.
- If possible, a core network can include multiple isolated logical networks tailored to support different service categories or subscriber groups. Furthermore, support creating the logical network instance on-demand.

4.2.2. Network Scalability

Comparing to traditional LTE UEs, the traffic generated by a very large number of connected low-complexity IoT devices typically will be a relatively low volume of non-delay-sensitive data and corresponding signaling. However, considering that one core network node serves a huge amount of such IoT services, MNO expects the network architecture to support:

- capacity scalable when meeting temporary intensive burst traffic
- prevent signaling and user data congestion raised from massive number of IoT devices.

4.2.3. Network Capability Exposure

To allow the 3rd party to access information regarding services provided by the cloT network (e.g. subscriber subscription information, connectivity information, mobility information, etc.), the CloT core network architecture is expected to provide open interfaces (APIs) for the 3rd party to access/exchange network information, that includes:

- Support the unified network capability exposure interfaces via centralized functions.
- Support acquiring information (e.g. network/connectivity information) and providing them to the 3rd party.
- Support to expose network capabilities to the 3rd party applications located within the operator domain close to the edge of the network, or outside the operator domain.

4.2.4. Automated network management

Accompanying with increasing number of network elements and complexity of the core network, automated network management becomes more and more important to reduce the OPEX and maintain the network in optimized operation mode. To realize this objective, it is expected that cloT network supports following related capabilities;

- Support a unified end-to-end network management to ensure compatibility and flexibility for the OAM of an CloT network.
- Support CloT dedicated performance management and configuration management.
- Support automated optimization mechanism to handle various traffic model shift.
- Support automated healing mechanism to avoid or mitigate service impact when fault happen.

5. Key technologies in CloT Core Network Architecture Design

5.1. Overview

According to architectural requirements identified in the chapter 4, following key technologies would be considered when designing a CloT core network architecture:

- Network function virtualization.
- Core network architecture optimization to support RAN IoT related features.
- Service capability exposure.

5.2. Network Function Virtualization

To realize network flexibility objective, the CloT is expected to support deployments in virtualized environments. NFV is a technology enabler which is relevant to a large amount of use cases including what were identified in the chapter 4. This key issue will determine the management

architecture impacts due to new characteristic of VNF management, for example, VNF lifecycle management operations (i.e. VNF instantiation, scaling, and termination).

5.3. Core Network Architecture Optimization

To fulfill various architectural requirement from both service and operation aspects, it is recommended to optimize the existing core network architecture to reflect the enhancement regarding service characteristic. This key issue will determine the potential solutions to meet identified service requirements, including (non-exclusive):

- Support power saving mode of IoT devices
- Optimize the signaling procedure to reduce the exchanged message amount or lead to lower the complexity of IoT devices
- Support logical network instance isolation if multiple logical network instances are deployed.

5.4. Service Capability Exposure

To offer on-demand or customized IoT services to IoT service provider, the CloT network is required to provide open interfaces to allow the 3rd party to access network information or to consume authorized network capabilities.

6. NFV based network architecture and management

6.1. NFV Overview

To address for CloT service characters in operator network and adapt to the TTM and TTC requirement for CloT deployment, the deployment of NFV and SDN technology in operator network is a key solution. It is well accepted in whole telecom industry that NFV based network frame is the key technology and enabler for future network evolution, and CloT will change information industry deeply and change the business model of operator and vendor. Tier1 operators in global has started the pace of NFV&SDN transformation to build up and to provide typical CloT UCs, it is expected based on LTE/EPC evolution, various CloT services will be deployed over NFV/SDN core network architecture.

NFV based network architecture has much benefits compared with legacy platform:

- Reduce operator CAPEX and OPEX by the decouple of network function for SW and HW as well as automatically deployment of operator service, it saves a lot of operator CAPEX and OPEX.
- New business model information industry including operator network interworks in much more flexibility in NFV network architecture, and in same time network exposure capability triggers more business model for CloT up-spring.
- Flexibility the network function is deployed in form of SW and chained together by the orchestration, the service network topology and capacity are deployed in big flexibility and scalability.
- Automation operator infrastructure network and service network works in way of automation that capability and scalable capacity be deployed by business centralized orchestrator.

6.2. NFV architecture Overview

NFV technical initiative and its specification cluster is specified by ETSI Industry Specification Group (ISG), according to specification of Network Functions Virtualization (NFV) by ETSI GS NFV 002, NFV architecture is described as Figure 6.2-1.

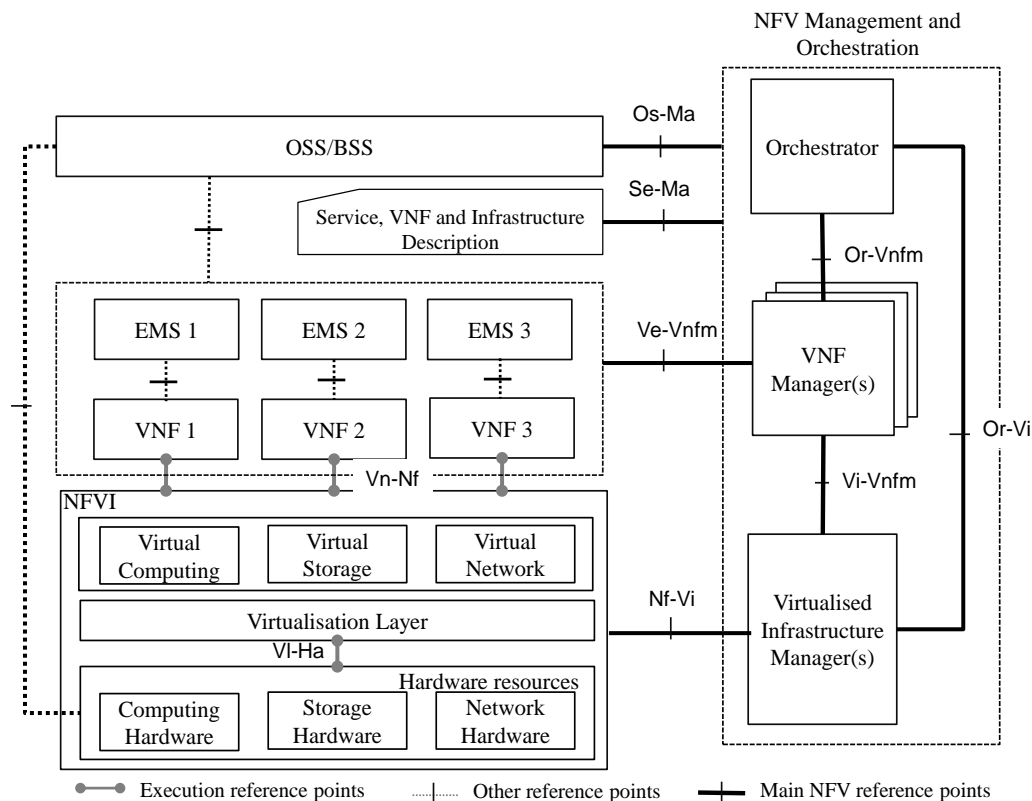


Figure 6.2-1 NFV architecture

In frame of NFV architecture SW/HW is decoupled and HW resources is represented as computing hardware、storage hardware and network hardware, which are virtualized by hypervisor. The use of hypervisors is one of the typical solutions for the deployment of VNF (Virtual Network Function), VM (Virtual Machine) is build up on hypervisor and operator's VNF is envisioned to be deployed on one or several VMs. VNF may also include software running on top of a non-virtualized server by means of an operating system.

From NFV's point of view, virtualized infrastructure management comprises the functions that are used to control and manage the interaction of a VNF with computing、storage and network resources, as well as their virtualization. According to hardware resources specified in the architecture, the Virtualized Infrastructure Manager performs the SW/HW resource inventory management, allocates and manages virtualized resource for VM/VNF and relevant network connectivity; it also fulfills alarm and performance management of SW/HW resource and other management functions.

Operator's network function is virtualized as VNF which is deployed in 3rd partner HW resources infrastructure. VNF is managed by VNF manager, which fulfills VNF life cycle management and coordinates with orchestrator to do network service deployment.

EMS manages VNF FCAPS, and operator's OSS/BSS system works together with NFV management and orchestration system implements service plan、 service deployment and service lifecycle management in way of orchestration and automation.

6.3. NFV Management Architecture Overview

NFV management system domain is composed of management and orchestration MANO in short. The MANO architecture is composed of virtualized infrastructure management VIM、 VNF management VNFM and Orchestrator.

ETSI GS NFV-MAN 001 specified NFV management architecture as the Figure 6.3-1

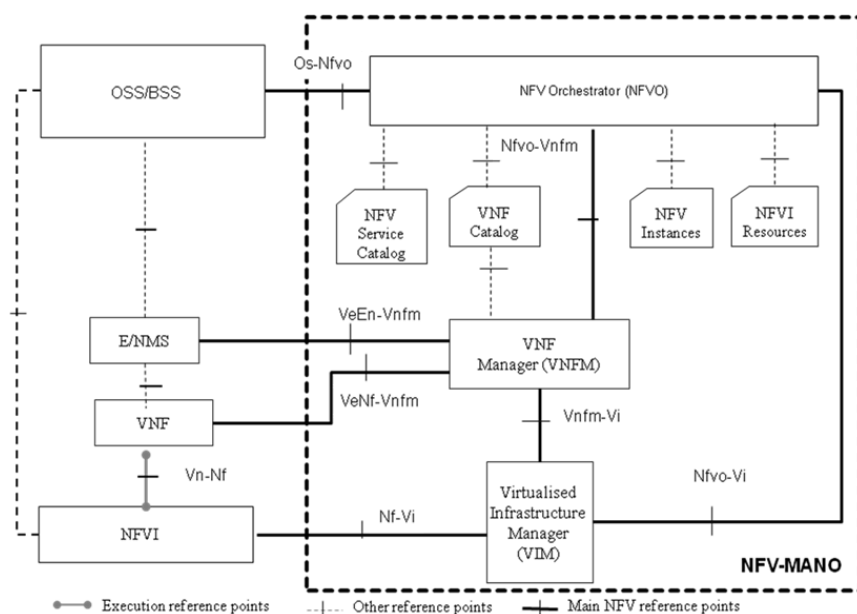


Figure 6.3-1 NFV management architecture

VIM、VNFM and Orchestrator interwork to manage NFVI resource and VNF lifecycle management. MANO works with operator's OSS/BSS system to fulfill network management、 analytics, it is responsible for deployment of network service and VNF instantiation.

Boarding NFV service catalog and VNF catalog describes operator services with NSD and VNFD; all network services and VNF services instantiated is managed in NFV repository. Orchestrator keeps the available NFVI resource in repository.

6.4. VNF Lifecycle Management

VNF works as the virtual network function instance which chained together to provide operator network service. The VNF is instantiated according to VNF catalog template by VNFM, which also fulfills all VNF instance lifecycle management includes:

- According to VNF catalog template to create a VNF instance
- To Scale a VNF instance capacity with scale-out and scale in operation
- To Upgrade VNF SW version and to update its configuration
- To terminate a VNF and to release its NFVI resource.

VNF lifecycle management interworks closely with Orchestrator and VIM, which manages NFVI resource management. VNFM has the capability to monitor the FACPS of VNF to trigger VNF lifecycle management actions. VNFM supplies service to orchestrator for VNF agility and automation management.

6.5. Network Service management and orchestration

NFV based network service management and orchestration works together with operator OSS/BSS system and works as entry to manage operator network agility and automation. VNFM and NFVI provide the agile capability and service to fulfill the network service deployment. network service has the capability to chain VNFs and PNFs together to support operator network migration

In general network service Orchestrator lifecycle management has below functions in operator NFV network architecture.

- To board a network service in form of catalog template.
- To instantiate a network service according to network service catalog template and to store it in NFV instance repository.
- To manage Network service capacity like agility
- To manage network service upgrades like SW upgrades and capacity expansion
- To manage network service logic with VNFFG
- VNFFG creation and deletion as well as update

Network service orchestrator interworks closely with OSS/BSS system in operator's network NFV evolution.

7. CloT Core Network Architecture and Solutions

7.1 CloT Optimization EPC architecture (NB-IoT and eMTC)

7.2.1. CloT Optimization EPC architecture

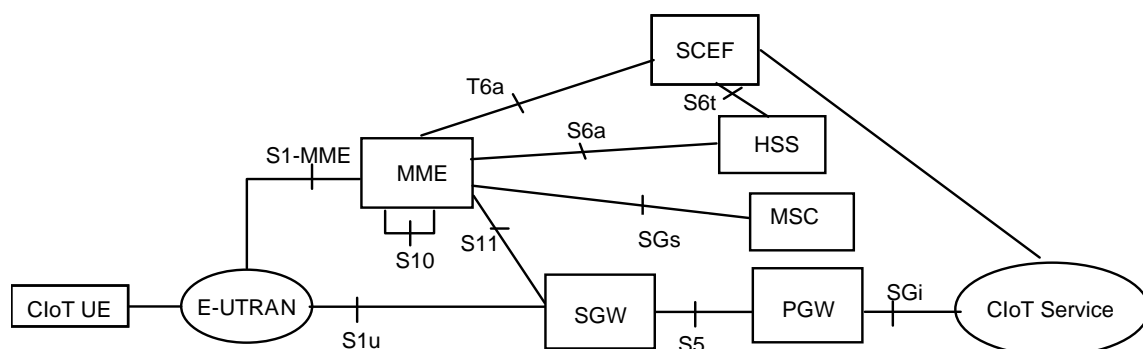


Figure 7.2-1: Optimized EPS architecture option for CloT - Non-roaming architecture

7.2.2. Network Functions Enhancements

CloT Optimization EPC including MME, SGW, PGW and SCEF supports necessary functionalities compared with the existing EPS core network elements and supports at least some of the following CloT optimizations:

- Control plane CloT EPS optimization for small data transmission.
- User plane CloT EPS optimization for small data transmission.
- Necessary security procedures for efficient small data transmission.
- SMS without combined attach for NB-IoT only UEs.
- Paging optimizations for coverage enhancements.
- Support for non-IP data transmission via SGi tunneling and/or SCEF.
- Support for Attach without PDN connectivity.

7.2.3. Reference Points

Modified interface:

- **S1-MME**: Reference point for the control plane protocol between E-UTRAN and MME.
- **S5**: It provides user plane tunnelling and tunnel management between Serving GW and PDN GW. It is used for Serving GW relocation due to UE mobility and if the Serving GW needs to connect to a non-collocated PDN GW for the required PDN connectivity.
- **S6a**: It enables transfer of subscription and authentication data for authenticating/authorizing user access to the evolved system between MME and HSS.
- **S10**: Reference point between MMEs for MME relocation and MME to MME information transfer.
- **S11-C**: Reference point between MME and Serving GW for signaling.

New added interface:

- **S11-U**: Reference point between MME and Serving GW for data transfer.
- **T6a**: Reference point between MME and SCEF (Service Capability Exposure Function).
- **S6t**: Reference point between HSS and SCEF.

Existing interface:

- **S1-U**: Reference point between E-UTRAN and Serving GW for the per bearer user plane tunnelling and inter eNodeB path switching during handover. S1-U does not apply to the Control Plane CloT EPS optimisation.
- **SGs**: Reference point between MME and MSC/VLR.
- **SGi**: It is the reference point between the PDN GW and the packet data network. Packet data network may be an operator external public or private packet data network or an intra operator packet data network

7.2. Key Solution in CloT Core Network

7.2.1. Efficient data transmission

7.2.1.1 Data over NAS (Control-Plane Optimization)

A CloT terminal accesses the PS core network over the S1 interface. When a terminal in the idle state has data to send, it needs to initiate a service request procedure to set up an S1-U connection and data radio bearer (DRB) before data transmission. For most CloT applications, only one data report is sent at a time. According to 3GPP CloT traffic model, 50% of the uplink reports have an ACK, indicating the other half does not have an ACK. The signaling overheads for setting up an S1-U connection and DRB are large. After a data report and associated ACK are exchanged, the S1-U connection and DRB need to be released, requiring more signaling overheads. Therefore, highly efficient small data transmission solutions are required.

To meet the service requirements, 3GPP proposes the following two data transmission solutions:

- Control-plane data transmission optimization: Service layer data can be transmitted in NAS messages without setting up an S1-U connection or DRB to reduce the total data transmission load (including signaling messages) for IoT terminals. In this way, the power consumption of IoT terminals is reduced, and the network load caused by signaling messages is also reduced.
- User-plane data transmission optimization: The UE, eNodeB, and MME save the data such as the S1AP association, UE context, and bearer context used for resuming a connection. When required later, the bearer can be quickly restored through a connection resume procedure, which reduces the network signaling load. However, the disadvantage of this solution is that the UE and eNodeB always must save the contexts.

The two optimization solutions can be used on the S1 link. When the NB-IoT technology is used, control-plane optimization is mandatory while user-plane optimization is optional. When MTC/eMTC technology is used, both control-plane optimization and user-plane optimization are optional.

Control-plane optimization requires the collaboration among the UE, RAN, and core network. The main optimization is described as follows and shown in the following figure:

- The UE transmits small data packets to the MME through NAS PDUs. In this process, NAS encryption and integrity protection are performed between the UE and MME.
- After receiving NAS PDUs, the MME converts the PDUs into GTP-U messages and sends them to the S-GW over the new S11-U interface. Then, the S-GW sends the messages to the P-GW and the P-GW sends the messages to the application server.

For network reconstruction using this solution, control-plane data transmission optimization needs to be selected on the MME. In addition, the S11-U interface needs to be configured between the MME and S-GW for data transmission.

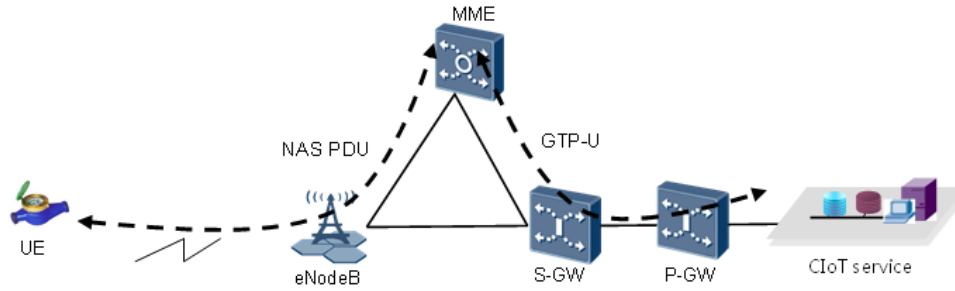


Figure 7.2.1.1-1: Control-Plane Optimization architecture

7.2.1.2 User-Plane Optimization

The traditional user-plane solution releases and reestablishes a connection through RRC connection release and establishment procedures. In the user-plane data transmission optimization solution, the connection suspend and resume procedures are used to suspend and resume a connection with the collaboration of the UE, RAN, and EPC. This is because the contexts are saved on the eNodeB, and the connection reestablishment requires only the copy of the information, such as the AS security context. This enables quick packet transmission.

This solution is implemented in the following procedure:

1. During the data transmission of the UE for the first time, contexts including the DRB and EPS bearer are set up.
2. After data transmission is complete, the air interface connection is released so that the UE enters the ECM-IDLE state. In a connection suspend procedure, the UE, eNodeB, and MME save the information such as the S1AP association, UE context, and bearer context used for resuming the connection.
3. When the UE needs to send packets again, a connection resume procedure can be performed to restore the contexts for data transmission.

7.2.1.3 Non-IP small data transmission

For most ClIoT applications, small data is reported infrequently. The usual data report size ranges from 20 to 200 bytes. The UDP/IP transport layer data is 28 bytes for IPv4-based transmission or 48 bytes for IPv6-based transmission, which is a high proportion in an entire data packet. Especially, when the valid load is lower than 20 bytes, the header size is larger than the data size. In this case, non-IP data transfer of UEs is a solution that significantly improves the data transmission efficiency on the radio network.

With the non-IP data transfer solution, non-IP contexts are created on MME/S-GW/P-GW for a UE and later the UE sends non-IP data involving only the application layer data without any IP and transport layer data. Non-IP data is transmitted on the EPC through the P-GW over the SGi interface or through the SCEF.

7.2.2. Low Power Consumption of Terminals

7.2.2.1 PSM

PSM is the main technology used for power consumption reduction of UEs. UEs in PSM disable their receiver and the access stratum (AS) functions over the air interface are stopped. Therefore, UEs in PSM do not monitor the paging channel over the air interface. PSM is suitable for low-frequency, periodic, and delay-insensitive small packets and can be applied to applications with extremely high power consumption reduction requirements. For example, data reporting is required once a day for water and gas meters with only batteries, and PSM which requires collaboration of UEs and the MME can be used.

In a service procedure related to PSM, the UE and MME negotiate through NAS messages to determine the length of an active timer. The timer is started when the UE becomes idle. When the active timer expires, the MME determines that the UE enters PSM and rejects downlink services and paging services for the UE, and the UE disables the AS functions such as cell selection, which reduces power consumption. In this way, the UE and network have synchronized information for PSM of the UE. The following figure shows the state transition of UEs supporting PSM.

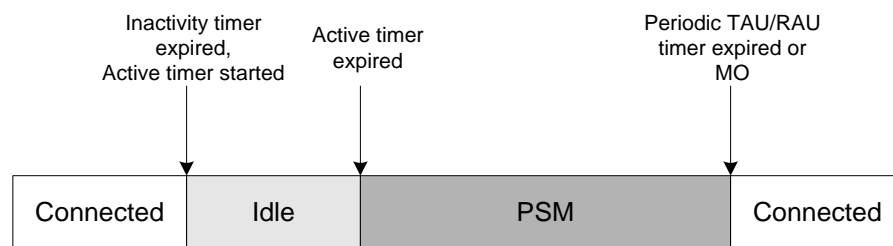


Figure 7.2.2.1-1: State of UE supporting PSM

After the UE enters PSM, the UE exits from the PSM only when it has MO data to send or periodic TAU/RAU needs to be performed upon periodic TAU/RAU timer expiry. The periodic TAU/RAU timer length can be specified on the network. To save power of the UE and reduce signaling load on the network, the periodic TAU/RAU timer length is generally not set to a small value.

7.2.2.2 Extended Periodic TAU Timer

When a UE enters PSM, it can be wakened upon an MO service or TAU/RAU timer expiry. The length of the periodic TAU/RAU timer on traditional networks is 186 minutes, which cannot meet the requirements of long-time dormancy for power saving purposes. Therefore, 3GPP extends the periodic TAU/RAU timer length to a maximum of 310 hours. In addition, associated subscription information is added on the HSS.

7.2.2.3 eDRX

DRX is used to reduce the power consumption of terminals in traditional networks, as shown in the left part of the above figure. A terminal monitors the page indicator channel (PICH) only at the paging opportunity within each DRX period, which helps reduce power consumption of this terminal. To achieve a compromise between low power consumption of terminals and prompt service provision, traditional networks require that the maximum DRX period be 2.56s.

To meet the requirements of more delay-insensitive services with lower power consumption, 3GPP proposes eDRX, whose period is extended to several minutes or even hours when compared with the DRX period, as shown in the right part of the above figure. The eDRX period ranges from 10.24s to about 2.913h (10.24s x 210), achieving a compromise between power saving and low delay for interactive applications such as smart bands, tracking, and street lamp control.

7.2.3. Coverage Enhancement

Enhanced coverage is one of the key objectives for the LPWA solution. Compared with the GPRS coverage, the NB-IoT coverage is enhanced with a 20 dB gain, and the eMTC coverage is enhanced with a 15 dB gain. However, IoT terminals are not always located in the area with poor coverage. Therefore, dividing the coverage into different levels and adopting different MCSs and repetition and spreading schemes for these levels will increase the system capacity, reduce power consumption of terminals, and decrease the delay without sacrificing the coverage. The scheme of Coverage Level based Paging Enhancement requires collaboration among the UE, eNodeB, and MME.

This scheme is implemented in the following procedure:

1. The UE and eNodeB negotiate and determine the coverage level of the UE. Then, the eNodeB includes the coverage level of the UE and the lists of recommended eNodeBs and cells in the UE Context Release Complete message sent to the MME for storage.
2. The MME determines whether to enable precise paging according to the local policy and decides the paging method according to the recommended information sent from the eNodeB and the list of learned adjacent eNodeBs. In the paging message, the MME includes the coverage level and assistant parameters including the number of paging attempts calculated for this paging decision and total number of planned paging attempts. Based on the information, the eNodeB determines the paging method.

7.2.4. Massive Access Control

7.2.4.1 Congestion Control Management

IoT terminals may frequently initiate attach or other procedures when, for example, periodic handshaking fails between IoT terminals and the IoT application server due to a server fault. This may lead to an abrupt increase of request messages, which triggers network-initiated flow

control. In this case, all UEs including common UEs may be affected.

There are a large number of IoT applications in different industries. Carriers currently have weak control over these applications. All IoT terminals in an industry may initiate services within a very short period. For example, all electric meters in an electric power company report data within a very short period. The access of a large number of IoT terminals (much more than common UEs) to the network within a short period easily results in network congestion, which affects services of common UEs. To guarantee user experience of common UEs and high-priority IoT terminals in the preceding scenario, congestion control solutions based on the MME, P-GW, and peripheral NEs need to be provided.

The MME congestion control solution includes the following two policies:

- **Signaling Congestion Control Based on the Back-off Timer**

A list of low-priority IoT APNs for congestion control is configured on the MME. The MME autonomously adjusts the access rates of terminals with low-priority IoT APNs according to the CPU usage. When the MME is slightly congested due to bursts of low-priority IoT service messages, the MME rejects excessive messages according to the low-priority access control threshold and includes a back-off timer in the rejection messages.

This policy aims to prevent congestion on the MME due to bursts of service messages sent from IoT terminals. If the traffic caused by bursts of IoT services is excessively heavy, which triggers queue-based flow control and self-protection flow control, the two flow control functions are performed in sequence before the signaling congestion control based on the back-off timer is performed.
- **Low-Priority Access Control**

3GPP specifications require that IoT terminals include a low-priority indicator in the RRC connection setup cause during their low-priority access. Then, the MME can perform low-priority access control based on the indicator. The MME autonomously adjusts the access rates of IoT terminals requesting low-priority services according to the CPU usage. When the MME is congested due to bursts of low-priority IoT services, the MME discards the low-priority service messages exceeding the allowed access rate. This prevents and mitigates system congestion.

When detecting a P-GW overload according to the CPU usage, the P-GW includes a P-GW back-off timer in the Create Session Response message to reject services with a low-priority IoT APN. This triggers the signaling congestion control based on the back-off timer policy. After receiving the P-GW back-off timer, the MME restricts the rate of the messages sent to the P-GW based on the low-priority IoT APN. This mitigates the P-GW congestion and improves user experience of high-priority services.

7.2.4.2 Rate Control

In normal scenarios, UEs send data packets infrequently for CIoT services. However, in abnormal scenarios (such as data is retransmitted due to timeout or malicious software is implanted in UEs), UEs send data packets frequently during a short period. This may cause network congestion and affect services. Air-interface resources of CIoT are precious. Therefore, 3GPP specifications

introduce the following two mechanisms for rate control when UEs send data packets:

- **Serving PLMN Based Clot Access Rate Control**

A rate threshold is configured on the MME and sent to the UE and S-GW/P-GW for controlling the numbers of uplink and downlink data packets transmitted between the UE and MME/S-GW/P-GW. This feature applies to control-plane data transmission optimization of Clot.

- **APN Based Clot Access Rate Control**

A rate threshold is configured on the P-GW for controlling the numbers of uplink and downlink data packets transmitted within a certain PDN connection of the UE. This feature applies to both control-plane and user-plane data transmission optimization of Clot.

7.2.5. Inter RAT idle mode mobility to/from NB-IoT

Considering that NB-IoT terminal has extreme power saving requirement and Low-frequency small packet, 3GPP R13 does not support connected/ idle mode mobility to and from the NB-IoT RAT, which the UE need detach/reattach every time when it moves from/to NB-IoT for dual-mode terminal.

However, 3GPP R14 considers that a UE which supports the NB-IoT RAT might benefit from supporting other RATs (and vice-versa) and then introduced the support for Inter-RAT idle mode mobility (connected mode mobility to and from NB-IoT RAT is still remained unsupported).

For example, an asset tracking application might support both an NB-IoT RAT and a 2G RAT to increase the number of locations where the UE will be able to gain connectivity. Moreover, in future, NB-IoT capability might be implemented on a smartphone. With this capability, the smartphone can be reachable when outside of WB-E-UTRAN coverage in some situation.

7.2.6. DECOR and eDECOR

Operator may require deploying an independent core network for IoT services. In such networking, how terminals are routed to the independent IoT PS core network and to the existing traditional PS core network must be solved.

As a new RAT, the NB-IoT requires independent spectrum and cell resources on the RAN side. Therefore, NB-IoT terminals can be routed to the dedicated IoT core network through RAN configurations.

If IoT terminals access the network in MTC/eMTC or LTE mode, discrimination between IoT terminals and traditional H2H terminals must be performed because the two RAT modes are considered the same and the RAN side cannot route them based on the RAT. To enable the discriminated routing, the dedicated core network (DECOR) and enhanced DECOR (eDECOR) are introduced to 3GPP specifications.

- According to 3GPP Release 13 specifications, the selection of the dedicated core network is achieved by the DECOR procedure. A new parameter "UE Usage Type" is added to the subscription data on the HSS. If the MME cannot provide services for the UE type indicated by "UE Usage Type", the MME selects the target MME for this specific UE type. The original MME sends a reroute command to instruct the RAN side to forward the initial UE message to the target MME. In this way, the selection of the dedicated core network is completed. The following figure shows the details.
- While enabling the selection of the dedicated core network, the DECOR solution increases

the signaling between the RAN and MME in the situation of massive IoT connections because a reroute procedure is required each time a terminal request to access the network. To solve this problem, eDECOR is introduced in 3GPP Release 14 specifications. When a terminal access the network for the first time, the MME performs a reroute procedure if a wrong core network is selected. The MME also sends the correct DCN ID to the UE. In subsequent access attempts, this terminal always sends this correct DCN ID to the RAN side. Therefore, the correct dedicated core network can be selected for the terminal without initiating a reroute procedure. As an enhancement to the DECOR solution, the eDECOR solution requires UE collaboration.

7.2.7. SMS

In addition to IP services, some IoT terminals need to support the SMS service to, for example, carry Device Trigger and initiate eSIM Profile download for UEs. Some IoT terminals do not support the IP protocol so that IoT service data needs to be transmitted over SMS. These scenarios require that the EPS support SMS. The EPS supports SMS in two ways. The SMS service can be implemented over the SGs interface between the MME and MSC/VLR or over the SGd interface between the MME and SMS-GMSC/IWMSC/SMS Router. The SGs interface applies when the carrier has deployed the 2G/3G CS domain. The SGd interface applies when the carrier does not have the 2G/3G CS domain or wishes to simplify CloT network architecture.

When a UE that only supports NB-IoT request SMS service but set the Attach Type to EPS attach instead of combined EPS/IMSI attach, a supporting MME provides SMS transfer without the UE performing the combined EPS attach.

After an NB-IoT UE initiates an EPS attach procedure through an Attach Request message containing the "Additional update type" IE with the value "SMS only," the MME exchanges messages between the UE and MSC as an intermediate node. This helps implement the SMS procedure.

8. CloT Core Network Deployment Analysis

8.1. Sharing the Core Network with Converged MBB Network

The CloT network shares the macro network. Based on the legacy or cloud platform on the existing network, the MME and SAE-GW are upgraded to seamlessly support NB-IoT and eMTC access. The software upgrade has minor impacts on peripheral networking. Only the S11-U interface is introduced between the MME and S-GW. The external networking of the core network remains unchanged. Sharing the Core Network with existing MBB Network has several advantages.

- Existing-network resources are fully utilized or shared. IoT services are processed by the existing EPC, saving the investment cost.
- The network topology is unchanged after the existing EPC is upgraded. Configurations are not required for massive eNodeBs/transport networks.
- The O&M is simple with no NEs added to the existing EMS/NMS.
- NB-IoT supports evolution according to subsequent 3GPP specifications. Inter-RAT-based project initiation is added in 3GPP Release 14, oriented to multi-RAT services established in the future and ensuring optimal experience.

8.2. CloT Dedicated core network for CloT

The CloT applications are supported by an independent core network. The cloud platform is recommended, using virtualization technologies and cloud architecture feature (such as elastic scaling) to meet the IoT requirements for massive connections. Dedicated Core Network has several advantages.

- The resources of common services and CloT services are separated and have independent O&M, preventing interference between them.
- Virtualization techniques such as flexible scaling meets the requirement of massive connections of IoT terminals.
- Independent software upgrade can be done on separated network if needed.
- Independent parameters can be configured.

9. CloT Service Capability Exposure

9.1. Overview

In recent years, a wealth of Internet applications brings the huge changes of the business form and operating mode, which gradually threatens the dominant position of operators in the mobile Internet ecological chain. Operators only rely on pipeline resources and network advantages, which cannot meet the mobile Internet business development need. Therefore, operators begin to explore how to deeply support the actual business needs of business providers, and give full play to the advantages of operators, to build a unified network of open capacity to provide a platform for operators and business providers to win a new situation of cooperation and win-win situation.

If application can easily obtain and use the mobile network capability information, good interaction is implemented between the application and the network. Operators can provide network capabilities based on the needs of business applications, optimize network resource allocation and traffic management, provide the right services for the upper business, thus new pipeline values are created, and new business models are formed for mobile networks, the new opportunity is brought to operators.

The mobile core network has service capabilities, such as Communications, Context, Subscription and Service Control that may be valuable to application providers. Communications refers to functions like voice calling, SMS, MMS. Subscription includes Subscription identity, feature sets and preference. Context covers real-time user information such as location, presence, profile, device capabilities and data connection type. Service Control addresses functions like Quality of Service, policy and security.

Mobile Network Operators (MNO) can offer value added services by exposing these service capabilities to external application providers, businesses and partners using web based API. In addition, mobile network operators can combine other internal or external services with their network capabilities to provide richer, composite API services to their partners. This brings mobile network intelligence to applications, allowing new, profitable business relationships to be created between MNOs and a wide range of external providers of enterprise/business solutions and web-based services or content.

Service Capability Exposure scenarios are widely used, for example, MTC terminal power saving mode setting, monitoring event for device management and fault detection, location service for

low power consumption and no GPS signal coverage, QoS management for establishing the session of a specific QoS for specific service, etc.

9.2. Service Capability Exposure Requirements

9.2.1. Real-time user status awareness

The Real-time user status awareness feature is intended for monitoring of specific events in 3GPP system and making such monitoring events information available for Service Capability Server (SCS) or Application Server (AS) via Service Capability Exposure Function (SCEF).

The events exposed via SCEF to SCS/AS are supported either by HSS, MME/SGSN or by PCRF. The following event types can be supported by HSS or MME/SGSN:

- Monitoring the association of the UE and UICC and/or new IMSI-IMEI-SV association;
- UE reachability;
- Location of the UE, and change in location of the UE;
- Loss of connectivity;
- Communication failure;
- Roaming status (i.e. Roaming or No Roaming) of the UE, and change in roaming status of the UE;
- Number of UEs present in a geographical area;
- Availability after DDN failure.

The following event types can be supported by PCRF:

- Location of the UE;
- Communication failure;

A User Equipment (UE) may adopt the Power Saving Mode (PSM) for reducing its power consumption. This mode like UE power-off, but the UE still stays registered with the network and there is no need to re-attach or re-establish PDN connections. A UE in PSM is not immediately reachable for mobile terminating data delivery. A UE using PSM is available for mobile terminating data delivery when the UE is in connected mode and during the period of an Active Time that is after the connected mode. SCS/AS can know the UE Power Saving status via monitoring the events of UE reachability and Loss of connectivity. SCS/AS can also know the mobile terminating data delivery status via monitoring the event of Availability after DDN failure. SCEF can also expose User subscription activation status to HSS, the status information is stored and provided from HSS.

9.2.2. Communications

Small data transmission service exposure is intended to provide a flexible method for Applications (SCS/AS) to communicate with UE. Service data will be contained in the API interface between SCEF and Applications (SCS/AS). These services include NIDD, SMS, MBMS, and so on.

- Non-IP data

Non-IP data is used for specific UE and Applications (SCS/AS) to save IP address resources for the scenarios of massive access. Because of its low-throughput, control plane transmission is

commonly used to reduce the radio resource consumption. Furthermore, NIDD via SCEF builds a short data transmission path between UE and Applications (SCS/AS), only eNB, MME and SCEF are involved in the path. Applications (SCS/AS) can deliver the service payload data directly using SCEF northbound API interface.

Between the SCEF and MME, a connection is established for each UE who wants to NIDD. Applications (SCS/AS) can also requests this connection establishment control by subscription data when this subscription data is exposed through SCEF. Applications (SCS/AS) should request the SCEF with configuration of NIDD such as Duration and Destination Address, and then SCEF will use the value as the default preference from the Applications (SCS/AS) when handling all non-IP packets associated with the NIDD connection. Through the API interface, Applications (SCS/AS) can deliver non-IP data of a given user as identified by External Identifier or MSISDN. SCEF is also required to execute authorization and load control for NIDD.

- SMS

SMS is the traditional way to realize the small data transmission and other additional value-added services, especially MTC device triggering. SCEF will acts as SME and provide SMS service to Applications (SCS/AS).

If MME supports SMS transmission via SGd, the applications data between UE and Applications (SCS/AS) which is carried by SMS payload will transport through the path of eNB-MME-SCEF, it is more like NIDD via SCEF.

For multiband device, Applications (SCS/AS) sends SMS to the UE to trigger the UE to initiate communication with the SCS when an IP address for the UE is not available or reachable by the Applications (SCS/AS).

9.2.3. Service parameter management

Service parameter should be exposed via SCEF to Applications for various service requirement. These parameters can be set for individual UE or a group of UE. Each parameter set normally has an associated validity time.

- power saving parameters

Although UE can negotiate the power saving parameters with core network, the core network should also have a flexible policy to control the request of UEs. On the other side, the Applications are really cleared with the requirement of service, so it is necessary for Applications to set the power saving parameters of core network. SCEF should expose the configuration of these parameters to the Applications by API. In the APIs, the parameters in the network side such

as Active Time, Periodic TAU Timer and eDRX cycle length, may be replaced by the description of transmission requirement such as Maximum Latency, Maximum Response Time and Maximum Connectivity Time, as power saving will go with high latency communications. SCEF will understand this requirement and modify relevant subscription in HSS, and then HSS insert subscription to MME. When MME negotiate the power saving parameters with UE, MME will consider the subscription as higher priority than UE's requested.

- communication parameters

3GPP has already defined Communication Pattern (CP) parameters for individual UEs or groups of UEs to describe predictable communication behaviors, including periodic time, duration time and scheduled time. This CP parameters is provided by Applications, provisioned through SCEF, and stored into subscription in HSS.

The CP parameters may be extended to traffic volume, location and mobility. Different parameters may be used by different network elements for different purpose, such as inactive timer value tuning in RAN side, mobility restriction, transmission mode decision and power saving parameters tuning in core network side, and so on.

SCEF may be able to map the CP parameters to corresponding network parameters and distribute them to different network elements.

- mobility parameters

Capability exposure for mobility management make it is possible to apply flexible mobility management and paging policy. For example, Applications can inform core network about UE's low-mobility character and even fixed location via SCEF, and then MME can restrict UE's access area and detect if terminal is stolen or misused.

For paging, efficiency and latency should be both taken in account. Low-mobility or fixed terminals should be paged in a small area and try to use coverage enhancement to save the resource of RAN and Core side, and high-mobility terminals should extend paging area when paging retry to improve paging successful rate, while low latency service need a quick paging and restrict using coverage enhancement. Applications are in charge of these communication requirement and mobility character, and indicate to SCEF through API interface. SCEF connect to corresponding network element and perform the configuration.

- Control Plane and User Plane preference

Control Plane CIoT EPS Optimization provides improved support of small data transfer via MME

by encapsulating them in NAS messages. It is the solution to support highly efficient handling of infrequent small data transmissions for IoT devices, and will reduce the total number of control plane messages when handling a short data transaction. Meanwhile, UE and Applications may change their transmission mode from control plane to user plane because of the change of transmission requirement from small infrequent data to continuous media stream. Since network is unable to predict the future transmission, SCEF should provide the Control Plane and User Plane preference and switching capability exposure to Applications. When Application is planning to perform continuous and large packet transmission, Applications is required to use API interface to change the Control Plane and User Plane preference data in subscription, or directly indicate MME to establish the S1-u bearer through SCEF.

9.2.4. Qos management

Qos management application scenarios are extensive. For example: when the MTC equipment and the server establish a reliable session connection for emergency services, a high Qos is required; when wireless network congestion occurs, Qos management for congestion cells is implemented; Applications can specify user's data flow to provide service acceleration or speed limit etc. through network capability exposure interface.

- Qos and bearers management

Applications (SCS/AS) may request that a data session to a UE that is served by the 3rd party service provider is set up with a specific QoS (e.g. low latency or jitter) and priority handling. This functionality is exposed via the SCEF towards the Applications (SCS/AS).

Applications (SCS/AS) can request the network to provide QoS for the AS session based on the application and service requirements with the help of a QoS reference parameter which refers to pre-defined QoS information. SCEF receives the request from the Applications (SCS/AS) to provide QoS for an AS session, the SCEF acts as an AF and transfers the request to provide QoS for an AS session to the PCRF via the Rx interface.

9.3. Service Capability Exposure Architecture

9.3.1. Overview

The core Network allows for exposure of network service capabilities to external 3rd party Application providers via Service Capability Exposure Function (SCEF). The SCEF provides a means to securely expose the services and capabilities provided by network. The SCEF provides a means for the discovery of the exposed services and capabilities. The SCEF provides access to network

capabilities through homogenous network APIs. The SCEF abstracts the services from the underlying core network interfaces and protocols.

The functionality of the SCEF may include the following:

- Authentication and Authorization:
- Identification of the API consumer,
- Profile management,
- ACL (access control list) management.
- Ability for the external entities to discover the exposed service capabilities
- Policy enforcement:
 - Infrastructural Policy: policies to protect platforms and network. An example of which maybe ensuring that a service node such as SMS-SC is not overloaded.
 - Business Policy: policies related to the specific functionalities exposed. An example may be number portability, service routing, subscriber consent etc.
 - Application Layer Policy: policies that are primarily focused on message payload or throughput provided by an application. An example may be throttling.
- Assurance:
 - Integration with O&M systems,
 - Assurance process related to usage of APIs.
- Accounting for inter operator settlements.
- Access: issues related to external interconnection and point of contact
- Abstraction: hides the underlying network interfaces and protocols to allow full network integration. The following functions are among those that may be supported:
 - Underlying protocol connectivity, routing and traffic control,
 - Mapping specific APIs onto appropriate network interfaces,
 - Protocol translation.

The SCEF shall protect the other PLMN entities (e.g. HSS, MME) from requests exceeding the permission arranged in the SLA with the third-party service provider.

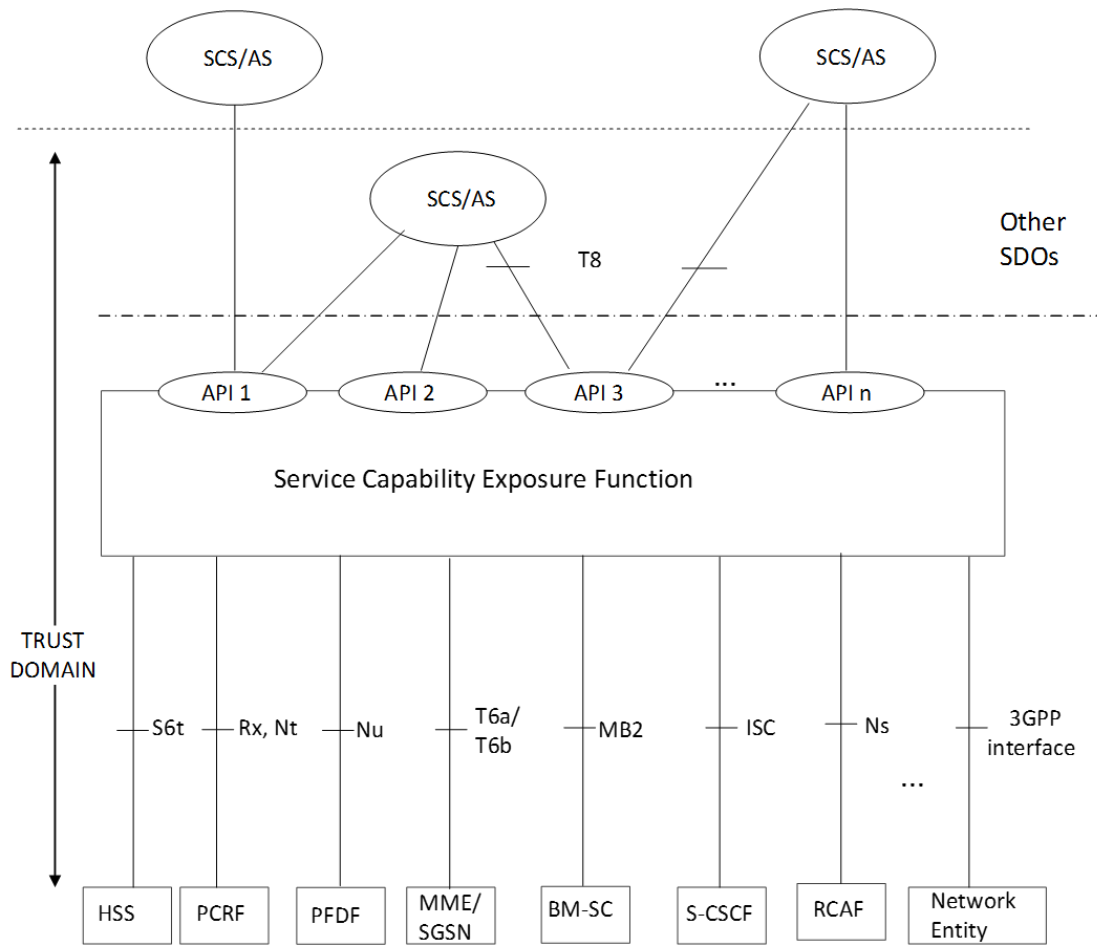


Figure 9.3.1-1 Architecture for Service Capability Exposure

The trust domain (see above figure) cover network elements that are protected by adequate network domain security, e.g. TLS, Oauth2.0. The elements and interfaces within the trust domain may all be within one operator's control, or some may be controlled by a trusted business partner which has a trust relationship with the operator e.g. another operator or a 3rd party. The APIs exposed from SCEF to trusted or non-trusted SCS/AS are defined as T8 interface in TS 29.122.

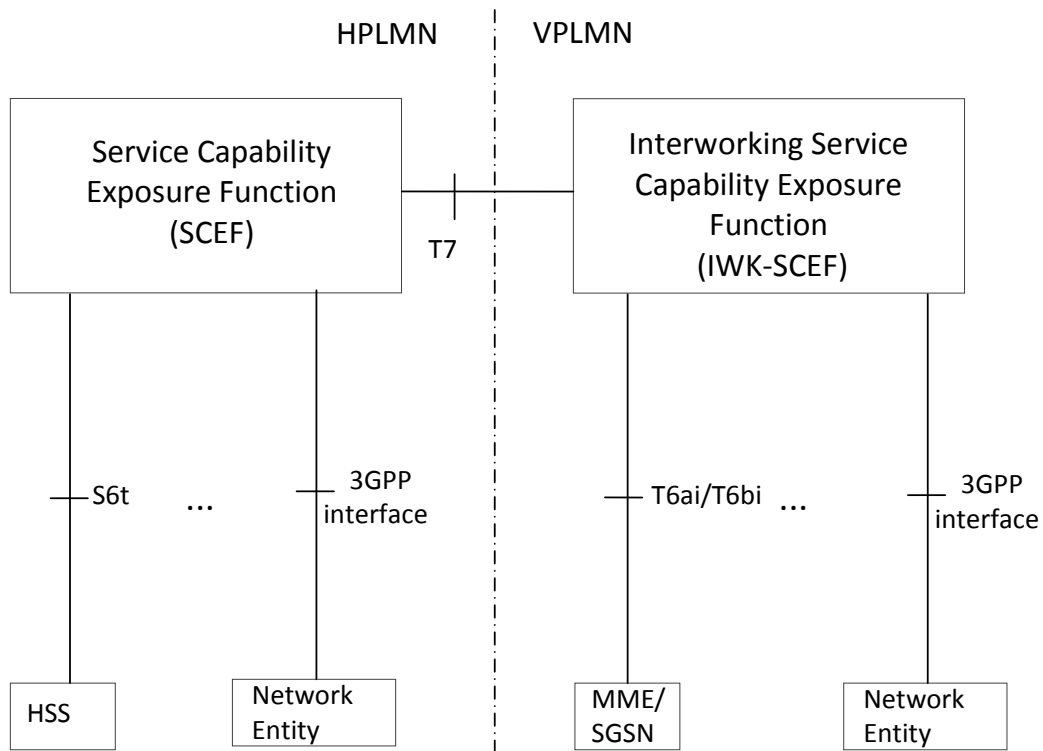


Figure 9.3.1-2 Roaming Architecture for Service Capability Exposure

For the roaming scenario, the Interworking SCEF (IWK-SCEF) shall have the connection with SCEF within the home network via T7 interface and with serving MME/SGSN in the visited network via T6ai/T6bi.

9.3.2. Interface

A set of interfaces for Network Service Capability Exposure are defined as below:

- Tsp: Reference interface used by a SCS to communicate with the MTC-IWF related control plane signalling.
- T4: Reference interface used between MTC-IWF and the SMS-SC in the HPLMN.
- T6a: Reference interface used between SCEF and serving MME.
- T6b: Reference interface used between SCEF and serving SGSN.
- T6ai: Reference interface used between IWK-SCEF and serving MME.
- T6bi: Reference interface used between IWK-SCEF and serving SGSN.
- T7: Reference interface used between IWK-SCEF and SCEF.

- T8: Reference interface used between the SCEF and the SCS/AS.
- S6m: Reference interface used by MTC-IWF to interrogate HSS/HLR.
- S6n: Reference interface used by MTC-AAA to interrogate HSS/HLR.
- S6t: Reference interface used between SCEF and HSS.
- MB2: Reference interface used by SCEF and BM-SC.
- Rx: Reference interface used by SCEF and PCRF.
- Ns: Reference interface used between SCEF and RCAF.
- Nt: Reference interface used by SCEF and PCRF.
- Nu: Reference interface used by SCEF to interact with the PPDF.

9.3.3. Identifier

External Identifier

Every UE may with one IMSI may have one or several External Identifier(s), the mapping of IMSI and External Identifier(s) are stored in the HSS. A standard External Identifier shall have follow components:

- Domain Identifier: identifies a domain that is under the control of a Mobile Network Operator (MNO). The Domain Identifier is used to identify where services or capabilities provided by the SCEF can be accessed. An operator may use different domain identifiers to provide access to different services and capabilities.
- Local Identifier: Identifier used to derive or obtain the IMSI. The Local Identifier shall be unique within the applicable domain. It is managed by the Mobile Network Operator.

External Group Identifier

Multiple External Identifiers can be grouped and identified by a global unique External Group Identifier. The mapping of External Group Identifiers and IMSI-Group Identifiers are stored in HSS. The External Group Identifier shall be formatted the same as the External Identifier.

The External Group Identifier is used on the interface between the SCS/AS and the SCEF and on the interface between the SCEF and the HSS.

9.3.4. Deployment and routing

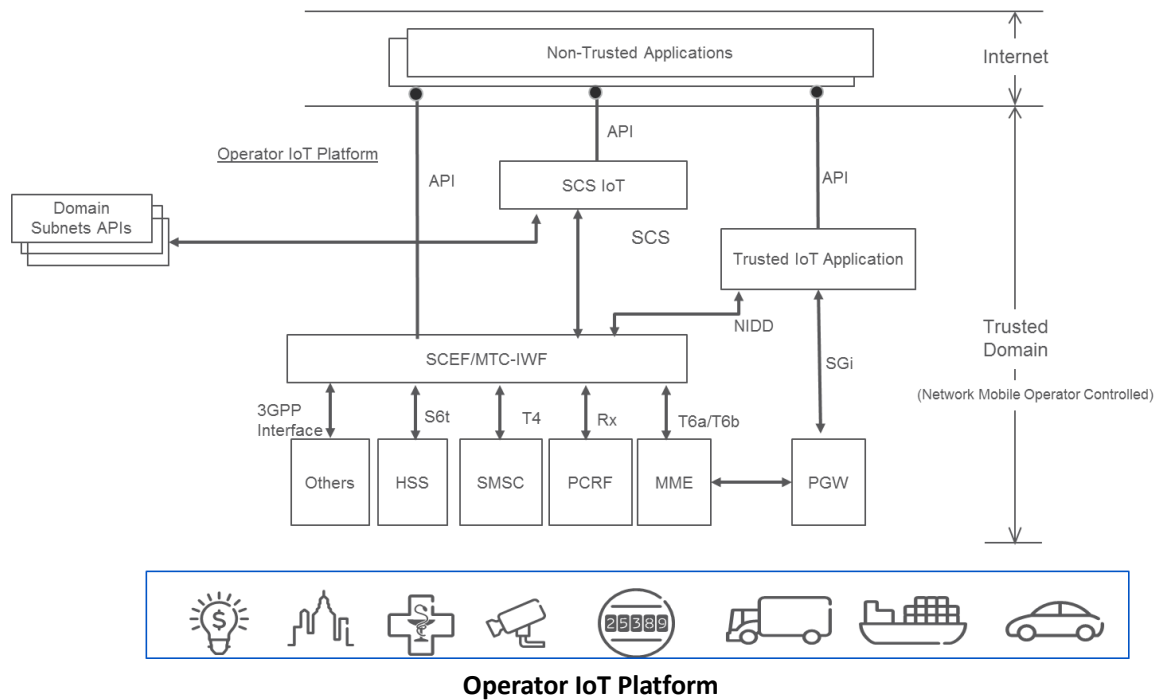


Figure 9.3.4

In certain deployments, the MTC-IWF may be co-located with the SCEF in which case MTC-IWF functionality is exposed to the SCS/AS via T8 interface. Operator from trusted domain can provide SCEF and SCS collocated deployment to expose not only pure network capabilities but also IoT Capabilities to non-trusted IoT application.

9.4. Security requirements of Service Capability Exposure

As the deployment of SCEF, it will be always within the trust domain, while applications belong to the trust domain or and outside the trust domain are both exist. So, the security requirement of SCEF northbound side is firstly using APP ID to identify different API consumers, and then do basic authentication and authorization by user name and password. Profile management and access control list management may will be also needed. Furthermore, API parameter validity checking and access resource authorization. For the Transmission security, HTTPS is used for the security of data transmission. For the deployment option of IoT platform lays over the SCEF, thus security of northbound can be simplified.

The security requirement of SCEF for the southbound side will be based on mutual authentication between 3GPP Network Entity and SCEF. For the privacy of user, the numbering information e.g. IMSI/IMPI should be only used in the operator's domain, and MSISDN or external identity used in northbound.

One more important things about security is to build effective policies to protect the whole system including IoT platforms, core network and SCEF itself keeping away from overloaded. It is required for SCEF to do the message payload or throughput control and throttling. API combined and user grouped is also the way to reduce this payload.

10. Resources/References

- [1] *3GPP TS 23.401: "General Packet Radio Service (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access".*
- [2] *Machina Research Report: "Service Provider Opportunities & Strategies in the Internet of Things".*
- [3] *GSMA INDUSTRY Paper: "Mobile Internet of Things Low Power Wide Area Connectivity".*
- [4] *3GPP TS 23.682: "Architecture enhancements to facilitate communications with packet data networks and applications".*
- [5] *3GPP TR 45.820.: "Cellular system support for ultra-low complexity and low throughput Internet of Things (CIoT)".*