

TSN Profiles for Service Provider Networks

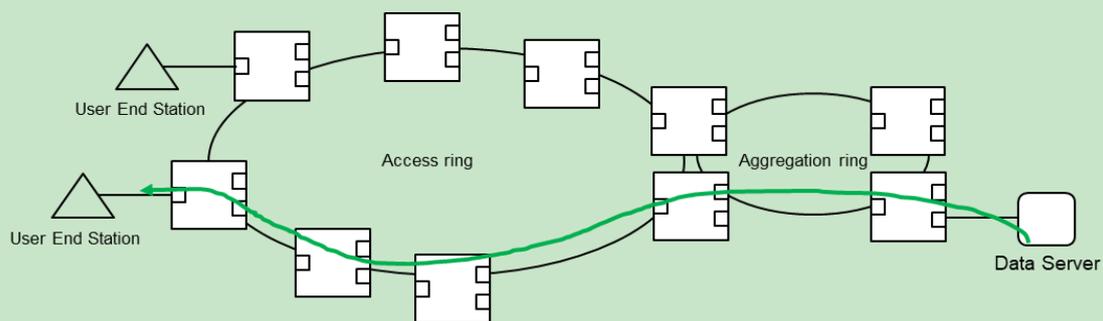
Tongtong Wang, April 2020

Huawei Technologies

1 Overview

Service provider networks, also called carrier networks, provide connectivity between access node and content sources (usually in data centers) for multiple users and applications. While 5G new technologies come into market, URLLC applications (e.g. vertical applications / utility networks) bring on stringent latency requirements over carrier networks.

As the following diagram shows, a typical service provider topology is like layered ring networks with sufficient redundant connections for better reliability and load balance. Usually user end stations are connected on access ring network, and multiple access rings could be linked to one aggregation ring where larger bandwidth links are shared by multiple users and services.



To specify and explain the selection of features and options, this document:

- Describe latency requirements for typical latency sensitive applications in service provider networks.
- Describes how the operation of bridges and bridged networks affects the quality of service provided by the carrier bridged network (Clause x), provides details in the calculation of latency. (Clause ...), and the potential impact of the use of flow control ;
- Specifies multiple profiles () that support the construction of bridged networks meeting latency requirements and jitter requirements.
- Defines service provider network profile conformance requirements () for bridges meeting either profile x requirements, for end stations and for synchronization.
- Provides a Profile Conformance Statement (PCS, Annex) to support clear detailed statements of equipment conformance to Service provider network profile requirements.
- Provide basic knowledge on Network Calculus to assist network latency evaluation.

2 Normative references

The following referenced documents are indispensable for the application of this document (i.e., they must be understood and used, so each referenced document is cited in the text and its relationship to this document is explained). For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments or corrigenda) applies.

IEEE Std 802 , IEEE Standard for Local and Metropolitan Area Networks—Overview and Architecture.4, 5

IEEE Std 802.1Q , IEEE Standard for Local and Metropolitan Area Networks—Bridges and Bridged Networks.

IEEE Std 802.3 , IEEE Standard for Ethernet.

IEEE Std 802.3br , IEEE Standard for Ethernet—Amendment 5: Specification and Management Parameters for Interspersing Express Traffic.

3GPP TS 23.501 G30 “System Architecture for the 5G System” V16.3.0

3 Definitions

For the purposes of this document, the following terms and definitions apply. The IEEE Standards Dictionary Online should be consulted for terms not defined in this clause.8

This standard makes use of the following terms defined in IEEE Std 802:

- bridge
- end station
- Ethernet
- forwarding
- frame
- Local Area Network (LAN)

This standard makes use of the following terms defined in IEEE Std 802.1Q:

- bridged network
- latency
- port
- traffic class

4 Acronyms and abbreviations

GBR – Guarantee Bit Rate

CBS – Credit Based Shaper

ATS – Asynchronous Traffic Shaping

TAS – Time Aware Shaper

5 Conformance

6 Service Provider Networks

<< Editor’s Note: This clause is a suggestion based on the presentation Suggestions for Service Provider Networks. <http://www.ieee802.org/1/files/public/docs2019/df-wangtt-SP-prof-outline-0519.pdf>

This clause will list a few representative use cases for service provider networks, and classify them from requirement perspective,

1. Bounded latency
2. Bounded jitter
 - ✓ Isolation
 - ✓ Slicing
3. Reliability

>>

Possible emerging applications on 5G carrier networks are discussed in 3GPP TS 23.501 [1], and summarized into three types of services shown in the following table. Bandwidth sensitive services have strict requirement on average bandwidth and loose constraint on latency, while connection services just require message delivery from time to time. A new type of service is the delay critical ones that will be the focus discussed in this draft.

Service Catalog	Examples	Packet delay budget	Packet loss rate	Default Max Data Burst
Bandwidth Sensitive Services (GBR)	Conservational Voice	100ms	10^{-2}	N/A
	Conversational Video (live streaming)	150ms	10^{-3}	N/A
	Real Time Gaming	50ms	10^{-3}	N/A
Connection Services (Non-GBR)	Buffered Streaming Video	300ms	10^{-6}	N/A

Latency Sensitive Services (Delay Critical)	Intelligent Transport Systems	30ms	10^{-5}	1354 bytes
	Smart Grid Tele-protection	5ms	10^{-5}	255 bytes

6.1 Bandwidth sensitive services

Bandwidth sensitive services like conversational voice usually have relaxed delay requirement over carrier network, and all packets traverse current IP DiffServ networks with legacy QoS methods like strict priority, weighted round robin, etc. Since carrier networks usually have large bandwidth and utilized as balanced as possible, thus traffic congestion rarely happens to high priority data streams. Bandwidth sensitive applications get satisfactory performance as long as adequate throughput and buffering capabilities are reserved and provided in time.

<<Editor Note: Consider to provide guideline on bandwidth analysis over carrier networks>>

6.2 Latency sensitive services

Latency sensitive services put more stringent requirement on end to end latency over carrier network, and any packet arrived later than deadline (tolerable deadline) is regarded as failure of packet delivery. A real case of latency sensitive service is smart grid tele-protection application, which requires end-to-end 2ms latency over carrier networks with 99.999% reliability. [1]

<<Editor Note: Consider to provide detailed latency evaluation method and compare multiple TSN techniques in Profiles section.>>

7 Profiles

7.1 Introduction

<< Editor's Note: This clause is a suggestion based on the presentation Suggestions for Service Provider Networks. <http://www.ieee802.org/1/files/public/docs2019/df-wangtt-SP-prof-outline-0519.pdf>

One or two profiles, for devices conformant to Clause 5, that will meet the needs of a significant market.

>>

7.2 Base Profile

<< Editor's Note: with existing TSN techniques, considering CBS, ATS, TAS, Strict Priority etc.

Compare latency and jitter, interference isolation etc.

>>

8 Interface with DetNet

8.1 Introduction

<< Editor's Note: This clause is a suggestion based on the presentation Suggestions for Service Provider Networks. <http://www.ieee802.org/1/files/public/docs2019/df-wangtt-SP-prof-outline-0519.pdf>

Control plane interface for resource reservation;

Data plane interface:

--Flow identification, flow aggregation; etc.

IETF DetNet has started working on the data plane;

>>

Annex A Concept for Network Calculus

<< Editor's Note: Basis of Network Calculus will be introduced briefly. Also considering re-visit latency evaluation for existing TSN techniques, like CBS, TAS, etc. Probably leads to maintenance for 802.1Q-2018 with update on latency analysis >>

Latency analysis based on Network Calculus (Informative)

<< Editor's Note: This clause may set an example on how to use profiles defined in this standard to setup a network to satisfy a certain use cases, such as smart grid or Cloud VR applications. >>

<< Editor's Note: briefly introduce Network Calculus methodology with examples. Illustrate how to use Network calculus to analyze delay on single node and cascaded networks.>>

Network calculus theory emerged during 1990s as a latency evaluation theory for quality of service analysis of packet switching networks, it is originally focus on performance analysis for

IntServ model over IP network. Data arrivals at a networked system are modelled by upper envelope functions. Minimum service guarantees that are provided by systems, such as a router, a scheduler, or a link, are characterized by service curves. Based on these concepts, network calculus offers convolution forms that enable worst case performance bounds evaluation including backlog and delay. Any number of bridged system in series can be transformed into a single equivalent system by convolution operation and obtain end-to-end performance.

A.1 Arrival curves

Flows can be described by arrival functions $F(t)$ that are given as the cumulated number of bits seen in an interval $[0,t]$. Arrival curves are defined to give an upper bound on the arrival functions, where

$$\alpha(t_2 - t_1) = F(t_2) - F(t_1);$$

Token bucket based arrival curve is usually featured like in equation, $\alpha(t) = b + rt$, where b is burst size, r is data rate;

<< Editor's Note: diagram of token bucket arrival curves will be helpful in this section. >>

A.2 Service Curves

The service offered by the scheduler on an output port can be characterized by a minimum service curve, denoted by $\beta(t)$. A common service curve is described as rate-latency service curve that includes a rate R and a latency T . Service curves of the rate-latency type can be implemented by Priority Queuing (PQ), Generalized Processor Sharing (GPS), Weighted Fair Queuing (WFQ), and further with TSN schedulers, where bandwidth resource R is assigned to selected traffic. However, in aggregated scheduling networks resources are provisioned on an aggregate basis.

<< Editor's Note: Consider a separate section to talk about aggregating mode. >>

Annex B Network Slicing

<<Editor's Note: Network slicing is essentially related to latency/jitter/reliability performance, TSN techniques provide different features and help implement network slicing over service provider networks.>>

Bibliography

[1] 3GPP TS 23.501 “System Architecture for the 5G System” V16.3.0;