### Architecture for delivering low latency service
### draft-ding-low-latency-delivery-arch-00

Abstract

   Demand for Ultra high-Reliability and Low-latency Communication
   (URLLC) and Broadband Assured IP Services (BAS) will grow as new
   service scenarios like 5G, IoT, AR/VR, Cloud are deployed.  As these
   new service scenarios will typically rely on shared packet
   infrastructure like Internet, methods to ensure URLLC and BAS
   performance across the underlying network resources will be required.
   This document outlines the motivation and key requirements for URLLC
   or BAS connectivity across heterogeneous network domains.  It also
   outlines the corresponding models and architecture required for
   providing orchestrated URLLC or BAS communication.

Status of This Memo

   This Internet-Draft is submitted in full conformance with the
   provisions of BCP 78 and BCP 79.

   Internet-Drafts are working documents of the Internet Engineering
   Task Force (IETF).  Note that other groups may also distribute
   working documents as Internet-Drafts.  The list of current Internet-
   Drafts is at http://datatracker.ietf.org/drafts/current/.

   Internet-Drafts are draft documents valid for a maximum of six months
   and may be updated, replaced, or obsoleted by other documents at any
   time.  It is inappropriate to use Internet-Drafts as reference
   material or to cite them other than as "work in progress."

   This Internet-Draft will expire on January 20, 2018.

Table of Contents

## 1.  Introduction

Low latency communications have been recently received much interest
such as those in Ultra high-Reliability and Low-latency Communication
(URLLC) and Broadband Assured IP Services (BAS) [BAS-Architecture].
Further investigation and requirements gathering is required.  Such
investigation should also build on existing IETF work, including
transport, security, and web technology and protocols effort [I-
D.arkko-arch-low-latency].  In parallel to the IETF efforts, relevant
discussion is ongoing including Time-Sensitive Networking Task Group
[TSN8021] in IEEE 802.1, 5G requirements for next-generation access
technology [TS38913]in 3GPP and BAS in BBF.

We may further scope the URLLC and BAS application requirements by explicitly involving end-to-end (E2E) service characteristics and capability requirements.  E2E service usually traverses multiple domains and involves multiple layers.  Yet, existing standards and current discussion typically focuses on a specific layer, protocol or link layer technology.  This myopic view lacks a holistic approach or system view on solving the URLLC and BAS problem space.

This draft identifies common URLLC and BAS application requirements in heterogeneous networks and key challenges for delivering suitable application and user Quality of Experience (QoE).  It analyses the applicability of existing technologies, and where necessary documents the gaps between URLLC/BAS requirements and network implementations.

Furthermore, the document proposes models and architecture to provide orchestrated URLLC or BAS communication.

## 1.1.  Requirements Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119.

## 2.  Requirements for delivery low latency services in heterogeneous networks

Emerging URLLC and BAS applications, such as self-driving cars, industrial control, real-time gaming, AR/VR, and Cloud based applications, introduce new requirements such as high reliability and low latency on data transmission.  For instance, in 5GPPP, the most stringent requirements on latency and reliability that we have identified relate to self-driving cars, where E2E latencies down to 1ms must be provided with a reliability of $1-10^{-9}$.  In other words, only one message in $10^9$ data transfers may be lost or delayed by more than 1ms when the latency budget is set to 1ms.

Since these applications would force the development for high reliability and low latency networking, monitoring and storing the latency performance across the network, and latency guarantee in each network segment or even node.  Unlike current packet network which is typically best-effort, making such latency guarantee is very difficult to achieve.

Multiple methods are being proposed to solve such latency guarantee issue, including cooperative hierarchical caching and routing, hardware acceleration, high-data throughput in the aggregation network, fog computing and mobile edge computing facilitating the

placement of compute and applications as close to the consumer as possible.

From the Internet stack perspective, improvement for communication latency may be achieved at multiple layers.  More recent technologies are being developed to reduce the communication latency, such as [L4S], [DETNET], [FlexE].  With such technologies, different network operators can build their own low latency networks.  For instance, each technology can be modelled as a network service for latency improvement, but often restricted to a specific domain or layer.

Typically, heterogeneous networks are composed of a wide-range of network segments traversing multiple domains and involving multiple layers.  Existing low latency technologies typically focus on specific layer, protocol or link.  With such diverse networks, it becomes very challenging to deliver low latency for E2E services.

Multiple technical proposals have described similar requirements discussed in this document, such as [I-D.dunbar-e2e-latency-arch-view-and-gaps] and [I-D.arkko-arch-low-latency].  BAS has discussed performance assurance of E2E services [BAS-Architecture].  From industrial automation perspective, 3GPP specification 22.282 has also defined latency requirement for robot control applications.

## [3](#). Application requirements and network performance

Application requirements can be modeled as Quality of Experience (QoE), and qualified by various service KQIs.  From users' perspective, QoE is the overall performance perception of the service.

Network performance can be evaluated by network KPIs such as delay, jitter, and packet loss.  As mentioned, URLLC and BAS applications require the capabilities of high reliability and low latency networking, which is unlike the current best-effort packet network.  Hence, it is important to identify and manage network KPIs to quantify and achieve the corresponding service KQI, as shown in below figure.  The KQI for a given service can be expressed as a function of a set of KPIs, expressed as KQI=f(KPI1, KPI2, ..., KPIn).

```
                          +-------------+
                          |Requirements |
                   +-------------+------+---------+
                   |             |                |
                   |             |                |
            +-----+-----+   +----+---+     +------+----+
Service KQI |   KQI1    |   | KQI2  |     |   KQI3    |
            +-----+-----+   +--------+     +----------+
                  |
         +-----------------------+--------------+-----------+
         |         |             |              |           |
Network+---v-+   +---v--+    +--------v--------+  +-v---+  +----v------+
KPI    |KPI1 |   | KPI2 |    |      KPI3       |  | ... |  |    ...    |
       +-----+   +------+    +-----------------+  +-----+  +----------+
```

   However, how to map URLLC and BAS application KQI to the
   corresponding set of network KPIs could be challenging.  Furthermore,
   there could be potential need of defining new network KPI (e.g.
   latency down to 1ms must be provided with a reliability of 1-10^-9)
   to reflect some new application requirements.

## 4.  Low latency delivery models

### 4.1.  Application service and network service model

   Application service model has information about application level
   policies and requirements, such as end user information, application
   service attributes.  Such model is constructed based on service KQIs.
   Below figure shows an example of application service model.

```
   +------------------+----------------------------------+
   | Service Name     | KQI Value                        |
   +------------------+----------------------------------+
   | 4K/8K Video      | quality/zap time/response time   |
   +------------------+----------------------------------+
   | Bank Transaction | transaction Rate/Locking/Idle Time|
   +------------------+----------------------------------+
   | Driving assistant| map updated time/map accuracy    |
   |------------------+----------------------------------+
   | VR/AR            | data rate/delay                  |
   |------------------+----------------------------------+
   | Factory Automation| real-time control/automation    |
   |------------------+----------------------------------+
```

   Network service model is used to describe the configuration, state
   data, operations and notifications of abstract representations of

services.  Take L2VPN service model as example [L2SM], it provides an
abstracted view of the L2VPN service configuration components, which
contains L2VPN domain relevant information, as well as network QoS or
KPI information in the L2VPN domain.

As mentioned in previous section, the latency sensitive applications
might traverse multiple domains and need E2E latency guarantee across
multiple domains.  Assuming the maximum latency is guaranteed and
cannot exceed a predefined value called MAX-LATENCY, the MAX-LATENCY
should be divided into multiple latency values and mapped to multiple
domains.  In each domain, the transmission latency must be guaranteed
less than the latency value allocated to it.

Some network KPI metrics of the network service model are listed in
below figure.  Note that the KPIs of latency bound and reliability
could be new element, compared to existing network service model, in
order to support the aforementioned new URLLC and BAS applications.

```
+--------------------+-----------------+
| KPI Name           | KPI Value       |
+--------------------+-----------------+
| Service type       | 4K/8K/VR etc    |
+--------------------+-----------------+
| User Information   | Triple-5/User ID |
+--------------------+-----------------+
| Service Profile    | Platinum/Gold/  |
+--------------------+-----------------+
| Latency bound      | MAX-LATENCY     |
|--------------------+-----------------+
| Reliability        | MAX-RELIABILITY |
|--------------------+-----------------+
| throughput         | MAX-THROUGHPUT  |
|--------------------+-----------------+
| packet loss rate   | MAX-PKTLOSSRATE |
|--------------------+-----------------+
| jitter             | MAX-JITTER      |
|--------------------+-----------------+
|bandwidth           | MAX-BANDWIDTH   |
|--------------------+-----------------+
```

The network service model shown in Figure 3 can be generic in the
sense that it has no assumption on the underlying network
technologies.  It is up to the network provider to translate this
network service model to specific network service models based on the

underlying network implementation, such as L2/L3VNF service model,
Detnet service model, FlexE/MPLS configurations, etc.
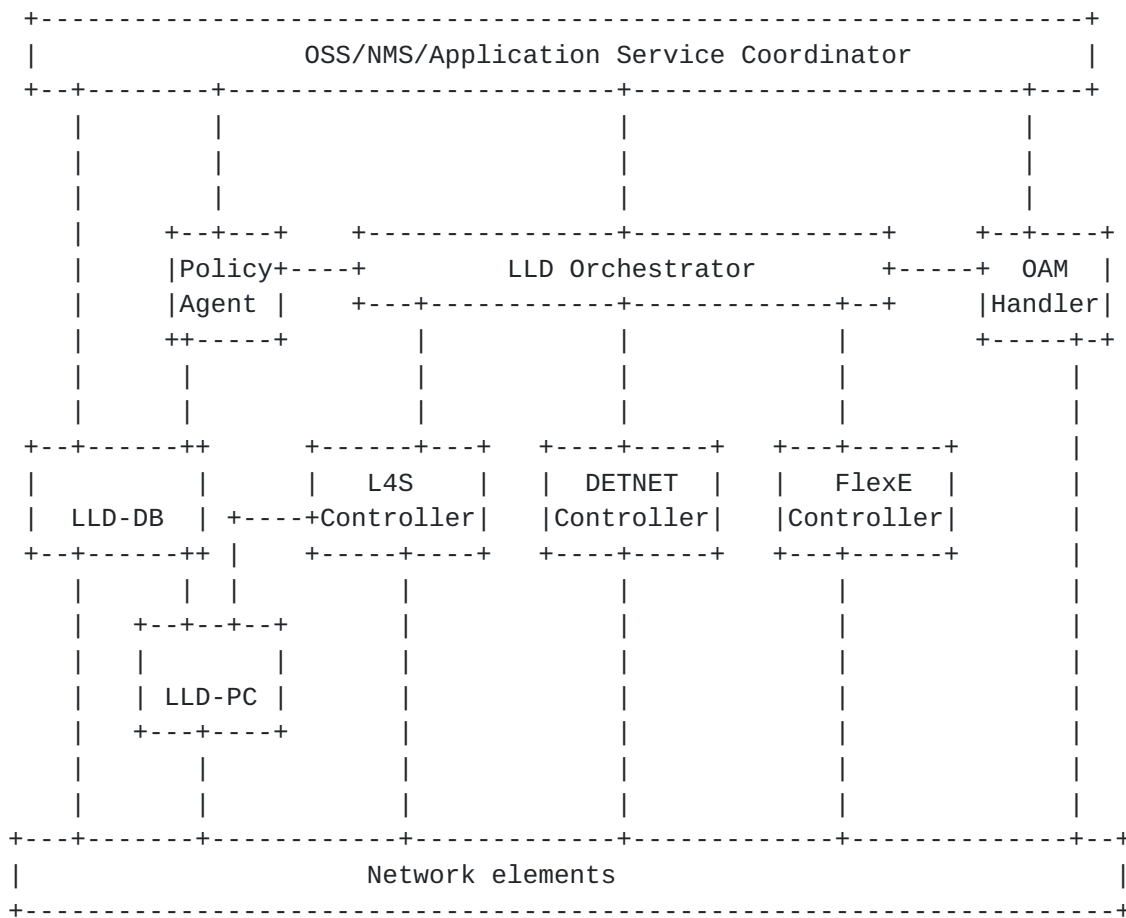
**4.2.  OAM model**

During each latency performance measurement period, latency metric is
sent to the OAM model ready to be analyzed.  Periodically, OAM model
retrieves aggregated monitored data and applies data classification
techniques to filter the data.  OAM model is responsible for
monitoring the reliability of the filtered data, and performs trouble
shooting based on the preconfigured reliability requirement.  If the
analyzed reliability of traffic data is lower than the preconfigured
reliability, OAM model issues a problem report.  Some parameters in
OAM model are listed in below figure.

```
+--------------------+------------------+
| Name               | Elements         |
+--------------------+------------------+
|traffic data        |                  |
+--------------------+------------------+
|minimum latency     |                  |
+--------------------+------------------+
|maximum latency     |                  |
|--------------------+------------------+
|average latency     |                  |
|--------------------+------------------+
|percentile latency  |                  |
|--------------------+------------------+
|queue length/size   |                  |
|--------------------+------------------+
```

**5.  Low latency delivery architecture**

**5.1.  Architecture Overview**

Below figure shows an architecture for low latency service delivery
(LLD).  It could be beneficial to define low latency delivery
architecture (or cook book) to coordinate and orchestrate multiple
low latency tools, in order to support low latency requirements from
user's perspective.  Note that LLD architecture has referred to ABNO
architecture [ABNO] especially the layer design, components
definition, etc.

```
 +----------------------------------------------------------------------+
 |                 OSS/NMS/Application Service Coordinator               |
 +--+--------+--------------------------+------------------------+---+
    |        |                          |                        |
    |        |                          |                        |
    |        |                          |                        |
    |     +--+---+    +---------------+----------------+    +--+----+
    |     |Policy+----+        LLD Orchestrator    +-----+  OAM  |
    |     |Agent |    +---+-----------+------------+--+    |Handler|
    |     ++-----+        |           |            |       +-----+-+
    |      |  |           |           |            |             |
    |      |  |           |           |            |             |
 +--+------++     +------+---+    +----+-----+    +---+------+    |
 |          |     |  L4S     |    | DETNET   |    |  FlexE   |    |
 |  LLD-DB  | +----+Controller|   |Controller|   |Controller|    |
 +--+------++ |    +-----+----+    +----+-----+    +---+------+    |
    |  |  | |        |           |            |             |
    |  +--+--+--+    |           |            |             |
    |  |     |       |           |            |             |
    |  | LLD-PC |    |           |            |             |
    |  +---+----+    |           |            |             |
    |      |         |           |            |             |
    |      |         |           |            |             |
 +---+-------+-----------+------------+------------+-------------+--+
 |                      Network elements                            |
 +----------------------------------------------------------------------+
```

## 5.2.  Components

### 5.2.1.  LLD orchestrator

   The LLD orchestrator is responsible to translate the generic network
   service model into the specific network service models, such as data
   model for L2VPN service delivery [L2SM], data model for L3VPN service
   delivery [RFC8049], and data model for EVPN [draft-ietf-bess-evpn-
   yang], in corresponding domains.

   Each domain has a separate controller that is responsible for
   receiving the network configuration from LLD orchestrator.  Based on
   the network configuration, the controller learns how to control the
   network elements.  One representative example of controller is PCE
   controller.

### 5.2.2.  OAM Handler

Latency measurement is also very crucial to make sure the latency
bound is not violated and useful for E2E latency aware OAM mechanism.
There is a need to support the measurement of latency inside of a
network device.

Existing technologies such as OWAMP [RFC4656] and TWAMP [RFC5357] is
focused on providing one way and two-way IP performance metrics.
Latency is one of metrics that can be used for E2E deterministic
latency provisioning.  Use OWAMP/TWAMP protocols or extension on that
to support measurement of flow latency performance is feasible.

The OAM Handler is responsible for monitoring the network elements,
collecting the measurement results and receiving notifications from
the network elements.  The OAM Handler also reports network
performance and problems to NMS/OSS/application service coordinator.

### 5.2.3.  Policy agent

Policy agent is configured by the NMS/OSS, and it is connected to
some components where the corresponding policy can be applied to.

### 5.3.  Functional interfaces

### 5.3.1.  Low latency path computation

Low latency path computation is a critical and fundamental feature
because individual controller in each domain is only able to share
abstracted information that is local to their domain.

Via the interface between LLD orchestrator and controller, the
controller gets the network service configuration and learns the
latency upper bound value in its domain.  After that, the controller
computes the optimized path to cover the latency upper bound, and
reserves and activate corresponding network resource for the path.

### 5.3.2.  OAM and report

OAM Handler interacts with the network to perform several actions:

Enabling OAM function within the network.

Performing proactive OAM operations in the network.

Receiving notifications of network events.

For low latency service, OAM handler correlates events reported from network and reports them onward to the LLD orchestrator and to the NMS/OSS/Application service coordinator.

## 6. Use cases

### 6.1. Network slicing

TBD

### 6.2. Provisioning E2E low latency path

TBD

## 7. Security Considerations

TBD

## 8. IANA Considerations

TBD

## 9. References

### 9.1. Normative References

TBD

### 9.2. Informative References

[I-D.dunbar-e2e-latency-arch-view-and-gaps] Dunbar, L., "Architectural View of E2E Latency and Gaps", draft-dunbar-e2e-latency-arch-view-and-gaps-01 (work in progress), March 2017.
[I-D.arkko-arch-low-latency] J. Arkko, "Low Latency Applications and the Internet Architecture", draft-arkko-arch-low-latency-00 (work in progress), November 2017.
[TS38913] "3rd Generation Partnership Project; Technical Specification Group Radio Access Network; Study on Scenarios and Requirements for Next Generation Access Technologies; (Release 14)", 3GPP Technical Report TR 38.913 V14.2.0, March 2017 (https://portal.3gpp.org/desktopmodules/Specifications/ SpecificationDetails.aspx?specificationId=2996).
[BAS-Architecture] Y.L. Jiang, "Broadband Assured IP Services Architecture", draft WT-387-00, broadband forum (BBF), July, 2016.
[IMTC-SDN] C. Lauwers, P. Menezes, M. Arndt,etc, International Multimedia Telecommunications Consortium (IMTC). http://lp.imtc.org/IMTC-SDN/.
[L2SM] B. Wen, G. Fioccola, C. Xie, L. Jalil, "A YANG Data Model for L2VPN Service Delivery", draft-ietf-l2sm-l2vpn-service-model-01 (work in progress), May 2017.

[RFC7491] D. King and A. Farrel, "A PCE-Based Architecture for Application-Based Network Operations ", RFC 7491, March 2015,

[DETNET] "Deterministic Networking (DETNET) Working Group", March2016 (https://tools.ietf.org/wg/detnet/charters).

[L4S] "Low Latency Low Loss Scalable throughput (L4S) Birds-of-Feather Session", July 2016 (https://datatracker.ietf.org/wg/l4s/charter/).

[FlexE] Stephen, J. and David. R. Stauffer, "FlexE Implementation Agreement", 2016.

[I-D.ietf-netmod-yang-model-classification] Bogdanovic, D., Claise, B., and C. Moberg, "YANG Module Classification", draft-ietf-netmod-yang-model-classification-07 (work in progress), May 2017.

Authors' Addresses

    Xiaojian Ding
    Huawei Technologies
    101 Software Avenue, Yuhua District
    Nanjing
    China

    EMail: dingxiaojian1@huawei.com


    Jinwei Xia
    Huawei Technologies
    101 Software Avenue, Yuhua District
    Nanjing
    China

    EMail: xiajinwei@huawei.com