

Network Working Group
Internet Draft
Intended status: Standard
Expires: September 10, 2021

L. Dunbar
H. Chen
Futurewei
Aijun Wang
China Telecom
March 10, 2021

**OSPF extension for 5G Edge Computing Service
draft-dunbar-lsr-5g-edge-compute-ospf-ext-04**

Abstract

This draft describes an OSPF extension for routers to advertise the running status and environment of the directly attached 5G Edge Computing servers. The AppMetaData can be used by the routers in the 5G Local Data Network to make intelligent decisions to optimize the forwarding of flows from UEs. The goal is to improve latency and performance for 5G Edge Computing services.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#). This document may not be modified, and derivative works of it may not be created, except to publish it as an RFC and to translate it into languages other than English.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>

This Internet-Draft will expire on April 7, 2021.

Copyright Notice

Copyright (c) 2021 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in [Section 4.e](#) of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction.....	3
1.1.	5G Edge Computing Background.....	3
1.2.	Problem#1: ANYCAST in 5G EC Environment.....	4
1.3.	Problem #2: Unbalanced Anycast Distribution due to UE Mobility.....	5
1.4.	Problem 3: Application Server Relocation.....	5
2.	Conventions used in this document.....	5
3.	Solution Overview.....	7
3.1.	Flow Affinity to an ANYCAST server.....	8
3.2.	IP Layer Metrics to Gauge App Server Running Status	8
3.3.	To Equalize traffic among Multiple ANYCAST Locations.....	9
3.4.	Reason for using IGP Based Solution.....	10
4.	Aggregated Cost Computed by Egress routers.....	11
4.1.	OSPFv3 LSA to carry the Aggregated Cost.....	11
4.2.	OSPFv2 LSA to carry the Aggregated Cost.....	11
5.	IP Layer App-Metrics Advertisements.....	11
5.1.	OSPFv3 Extension to carry the App-Metrics.....	12

5.2. OSPFv2 Extension to advertise the IP Layer App-Metrics.....	13
5.3. IP Layer App-Metrics Sub-TLVs.....	14
6. Soft Anchoring of an ANYCAST Flow.... Error! Bookmark not defined.	
7. Manageability Considerations.....	16
8. Security Considerations.....	16
9. IANA Considerations.....	16
10. References.....	16
10.1. Normative References.....	17
10.2. Informative References.....	17
11. Acknowledgments.....	18

[1. Introduction](#)

This document describes an OSPF extension to distribute the 5G Edge Computing App running status and environment so that other routers in the 5G Local Data Network (LDN) can make intelligent decisions to optimize the forwarding of flows from UEs. The goal is to improve latency and performance for 5G Edge Computing services.

1.1. 5G Edge Computing Background

As described in [[3GPP-EdgeComputing](#)], it is desirable for a mission critical Application to have multiple Application Servers hosted in multiple Edge Computing data centers to minimize the latency and to optimize the user experience. Those Edge Computing data centers are usually very close to or co-located with 5G base stations.

When a UE (User Equipment) initiates application packets using the destination address from a DNS reply or its cache, the packets from the UE are carried in a PDU session through 5G Core [5GC] to the 5G UPF-PSA (User Plan Function - PDU Session Anchor). The UPF-PSA decapsulates the 5G GTP outer header and forwards the packets from the UEs to the Ingress router of the Edge Computing (EC) Local Data Network (LDN) which is responsible for forwarding the packets to the intended destinations.

When the UE moves out of coverage of its current gNB (next-generation Node B) (gNB1), the handover procedure is initiated which includes the 5G SMF (Session Management

Function) selecting a new UPF-PSA [3GPP TS 23.501 and TS 23.502]. When the handover process is complete, the UE has a new IP address and the IP point of attachment is to the new UPF-PSA. 5GC may maintain a path from the old UPF to new the UPF for a short time for SSC [Session and Service Continuity] mode 3 to make the handover process more seamless.

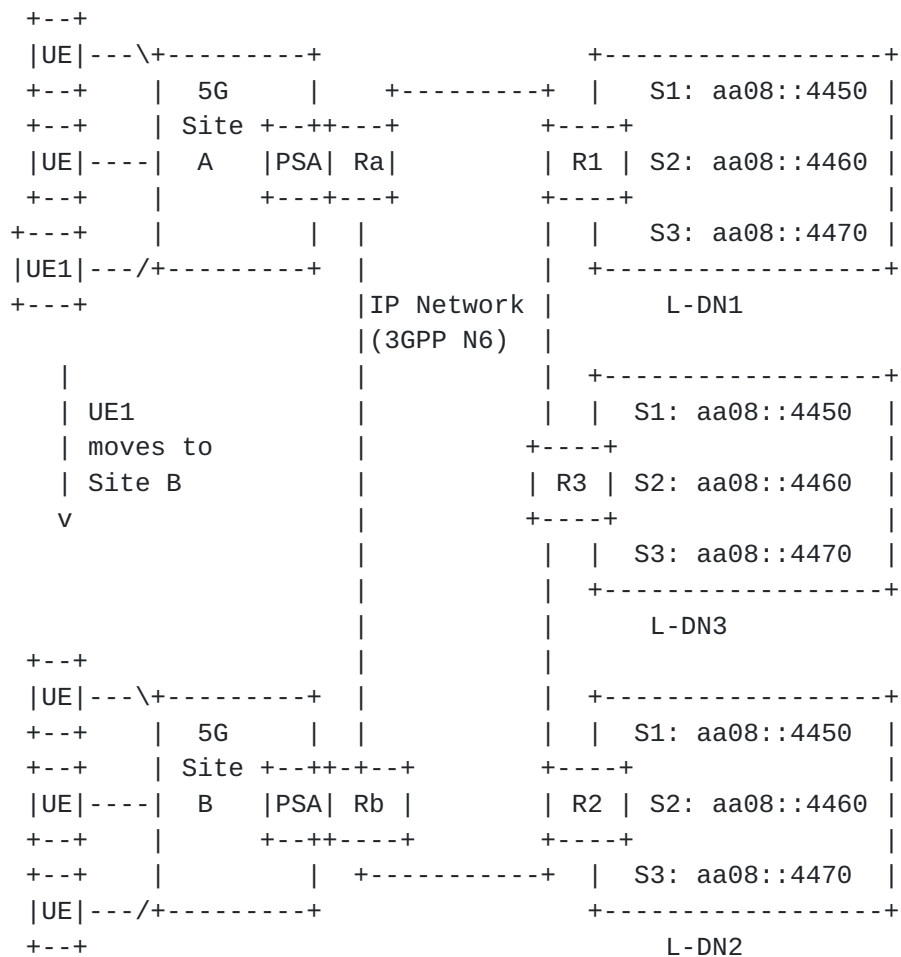


Figure 1: App Servers in different edge DCs

1.2. Problem#1: ANYCAST in 5G EC Environment

Increasingly, ANYCAST is used extensively by various application providers and CDNs because ANYCAST makes it possible to dynamically load balance across server locations based on network conditions. With multiple

servers having the same ANYCAST address, it eliminates the single point of failure and bottleneck at the application layer load balancers. Another benefit of using ANYCAST address is removing the dependency on how UEs get the IP addresses for their Applications. Some UEs (or clients) might use stale cached IP addresses for an extended period.

But, having multiple locations of the same ANYCAST address in 5G Edge Computing environment can be problematic because all those edge computing Data Centers can be close in proximity. There might be very little difference in the routing cost to reach the Application Servers in different Edge DCs, which can cause packets from one flow to be forwarded to different locations, resulting in service glitches.

1.3. Problem #2: Unbalanced Anycast Distribution due to UE Mobility

UEs' frequent moving from one 5G site to another can make it difficult to plan where the App Servers should be hosted. When one App server is heavily utilized, other App servers of the same address close-by can be very under-utilized. Since the condition can be short-lived, it is difficult for the application controller to anticipate the move and adjust.

1.4. Problem 3: Application Server Relocation

When an Application Server is added to, moved, or deleted from a 5G Edge Computing Data Center, not only the reachability changes but also the utilization and capacity for the Data Center might change.

Note: for the ease of description, the Edge Computing server, Application server, App server are used interchangeably throughout this document.

2. Conventions used in this document

A-ER: Egress Router to an Application Server, [A-ER] is used to describe the last router that the Application Server is attached. For 5G EC

environment, the A-ER can be the gateway router to a (mini) Edge Computing Data Center.

Application Server: An application server is a physical or virtual server that hosts the software system for the application.

Application Server Location: Represent a cluster of servers at one location serving the same Application. One application may have a Layer 7 Load balancer, whose address(es) are reachable from an external IP network, in front of a set of application servers. From IP network perspective, this whole group of servers is considered as the Application server at the location.

Edge Application Server: used interchangeably with Application Server throughout this document.

EC: Edge Computing

Edge Hosting Environment: An environment providing the support required for Edge Application Server's execution.

NOTE: The above terminologies are the same as those used in 3GPP TR 23.758

Edge DC: Edge Data Center, which provides the Edge Computing Hosting Environment. It might be co-located with 5G Base Station and not only host 5G core functions, but also host frequently used Edge server instances.

gNB next generation Node B

LDN: Local Data Network

PSA: PDU Session Anchor (UPF)

SSC: Session and Service Continuity

UE: User Equipment

UPF: User Plane Function

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [BCP 14](#) [[RFC2119](#)] [[RFC8174](#)] when, and only when, they appear in all capitals, as shown here.

3. Solution Overview

From IP Layer, the Application Servers are identified by their IP (ANYCAST) addresses. To a router, having multiple servers with the same (ANYCAST) address attached to different egress routers (A-ER) is same as having multiple paths to reach the (ANYCAST) address.

There are many tools available to influence the path section on a router, such as the routing distance, TE metrics, policies, etc. This draft describes a solution to add "Site-Cost" to influence the path selection. The "Site-Cost", which is derived from "site-capacity + load measurement + Preference + xxx", can be raw measurements collected by the egress routers based on the instructions from a controller or can be informed by the App Controller periodically.

The proposed solution is for the egress router (A-ER) that have a direct connection to the Application Servers to collect desired measurements about the Servers' running status and advertise the metrics to other routers in 5G EC LDN.

The solution assumes that the 5G Edge Computing controller or management system is aware of the ANYCAST addresses that need optimized forwarding. To minimize the processing on routers, only the application flows that match with the ACLs configured by the 5G Edge Computing controller will collect and advertise the desired measurements.

3.1. Flow Affinity to an ANYCAST server

Having multiple Edge Computing Servers or App Layer Load Balancers with the same ANYCAST address attached to multiple A-ERs, Flow Affinity means routers sending the packets of the same flow to the same A-ER even if the cost towards the A-ER is no longer optimal.

Many commercial routers today support some forms of flow affinity to ensure packets belonging to one flow be forwarded along the same path.

Editor's note: for IPv6 traffic, Flow Affinity can be supported by the routers of the Local Data Network (LDN) forwarding the packets with the same Flow Label in the packets' IPv6 Header along the same path towards the same egress router.

3.2. IP Layer Metrics to Gauge App Server Running Status

Most applications do not expose their internal logic to the network. Their communications are generally encrypted. Most of them do not even respond to PING or ICMP messages initiated by routers or network gears.

[5G-EC-Metrics] describes the IP Layer Metrics that can gauge the application servers running status and environment:

- IP-Layer Metric for App Server Load Measurement:
The Load Measurement to an App Server is a weighted combination of the number of packets/bytes to the App Server and the number of packets/bytes from the App Server which are collected by the A-ER that has the direct connection to the App Server.
The A-ER is configured with an ACL that can filter out the packets for the Application Server.
- Capacity Index:
Capacity Index is used to differentiate the running environment of the attached application server. Some data centers can have hundreds, or thousands, of servers behind an application server's App Layer Load Balancer. Other data centers can have a very small number of servers for the application. "Capacity

Index", which is a numeric number, is used to represent the capacity of the application server attached to an A-ER.

- Site preference index:
[IPv6-StickyService] describes a scenario that some sites are more preferred for handling an application than others for flows from a specific UE.

For ease of description, those metrics, more may be added later, are called IP Layer App-Metrics throughout the document.

3.3. To Equalize traffic among Multiple ANYCAST Locations

The main benefit of using ANYCAST is to leverage the network layer information to balance the traffic among multiple Application Server locations.

For 5G Edge Computing environment, the routers in the LDN need to be notified of various measurements of the App Servers attached to each A-ER to make the intelligent decision on where to forward the traffic for the application from UEs.

[5G-EC-Metrics] describes the algorithms that can be used by the routers in LDN to compare the cost to reach the App Servers between the Site-i or Site-j:

$$\text{Cost-i} = \min(w * \frac{\text{Load-i} * \text{CP-j}}{\text{Load-j} * \text{CP-i}} + (1-w) * \frac{\text{Pref-j} * \text{Network-Delay-i}}{\text{Pref-i} * \text{Network-Delay-j}})$$

Load-i: Load Index at Site-i, it is the weighted combination of the total packets or/and bytes sent to and received from the Application Server at Site-i during a fixed time period.

CP-i: capacity index at site I, a higher value means higher capacity.

Network Delay-i: Network latency measurement (RTT) to the A-ER that has the Application Server attached at the site-i.

Noted: Ingress nodes can easily measure RTT to all the egress nodes by existing IPPM metrics. But it is not so easy for ingress nodes to measure RTT to all the App Servers. Therefore, "Network-Delay-i", a.k.a. Network latency measurement (RTT), is between the Ingress nodes and egress nodes. The link cost between the egress nodes to their attached servers are embedded in the "capacity index".

Pref-i: Preference index for site-i, a higher value means higher preference.

w: Weight for load and site information, which is a value between 0 and 1. If smaller than 0.5, Network latency and the site Preference have more influence; otherwise, Server load and its capacity have more influence.

3.4. Reason for using IGP Based Solution

Here are some benefits of using IGP to propagate the IP Layer App-Metrics:

- Intermediate routers can derive the aggregated cost to reach the Application Servers attached to different egress nodes, especially:
 - The path to the optimal egress node can be more accurate or shorter
 - Convergence is shorter when there is any failure along the way towards the optimal ANYCAST server.
 - When there is any failure at the intended ANYCAST server, all the transient packets can be optimally forwarded to another App Server attached to a different egress router.
- Doesn't need the ingress nodes to establish tunnels with egress nodes.

There are limitations of using IGP too, such as:

- The IGP approach might not suit well to 5G EC LDN operated by multiple ISPs networks.
For LDN operated by multiple ISPs, BGP should be used. AppMetaData NLRI Path Attribute [[5G-AppMetaData](#)] describes the BGP UPDATE message to propagate IP Layer App-Metrics crossing multiple ISPs.

4. Aggregated Cost Computed by Egress Routers

If all egress routers that have a direct connection to the App Servers can get a periodic update of the aggregated cost to the App Servers or can be configured with a consistent algorithm to compute an aggregated cost that takes into consideration the Load Measurement, Capacity value, and Preference value, this aggregated cost can be considered as the Metric of the link to the App Server.

In this scenario, there is no protocol extension needed.

4.1. OSPFv3 LSA to carry the Aggregated Cost

If the App Servers use IPv6 ANYCAST address, the aggregated cost computed by the egress routers can be encoded in the Metric field [the interface cost] of Intra-Area-Prefix-LSA specified by [Section 3.7](#) of the [[RFC5340](#)].

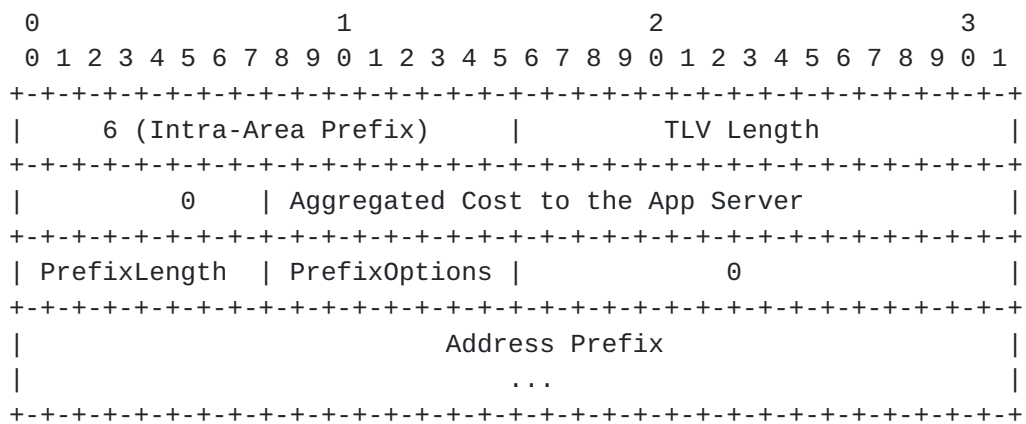


Figure 2: Aggregated Cost to App Server

4.2. OSPFv2 LSA to carry the Aggregated Cost

For App Servers in IPv4 address, the Aggregated Cost can be encoded in the "Metric" field of the Stub Link LSA [Link type =3] specified by [Section 12.4](#) of the [[RFC2328](#)].

5. IP Layer App-Metrics Advertisements

This section describes the OSPF extension that can carry the detailed IP Layer Metrics when it is not possible for all the egress routers to have a consistent algorithm to compute the aggregated cost or some routers need all the detailed IP Layer metrics for the App Servers for other purposes.

Since only a subset of routers within an IGP domain need to know those detailed metrics, it makes sense to use the OSPFv2 Extended Prefix Opaque LSA for IPv4 and OSPFv3 Extended LSA with Intra-Area-Prefix TLV to carry the detailed sub-TLVs. For routers that don't care about those metrics, they can ignore them very easily.

It worth noting that not all hosts (prefix) attached to an A-ER are ANYCAST servers that need network optimization. An A-ER only needs to advertise the App-Metrics for the ANYCAST addresses that match with the configured ACLs.

Draft [[draft-wang-lsr-passive-interface-attribute](#)] introduces the Stub-Link TLV for OSPFv2/v3 and ISIS protocol respectively. Considering the interfaces on an edge router that connects to the App servers are normally configured as passive interfaces, these IP-layer App-metrics can also be advertised as the attributes of the passive/stub link. The associated prefixes can then be advertised in the "Stub-Link Prefix Sub-TLV" that is defined in [[draft-wang-lsr-passive-interface-attribute](#)]. All the associated prefixes share the same characteristic of the link. Other link related sub-TLVs defined in [[RFC8920](#)] can also be attached and applied to the calculation of path to the associated prefixes.

5.1. OSPFv3 Extension to carry the App-Metrics

For App Servers using IPv6, the OSPFv3 Extended LSA with the Intra-Area-Prefix Address TLV specified by the [Section 3.7 of RFC8362](#) can be used to carry the App-Metrics for the attached App Servers.

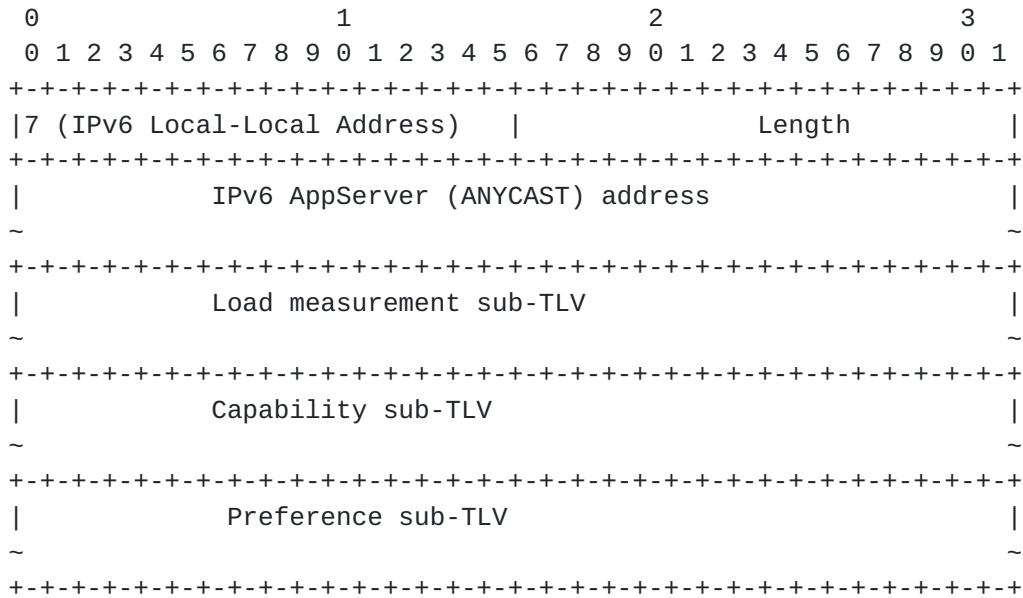


Figure 3: IPv6 App Server App-Metrics Encoding

5.2. OSPFv2 Extension to advertise the IP Layer App-Metrics

For App Servers using IPv4 addresses, the OSPFv2 Extended Prefix Opaque LSA with the extended Prefix TLV can be used to carry the App Metrics sub-TLVs, as specified by the [Section 2.1 \[RFC7684\]](#).

Here is the proposed encoding:

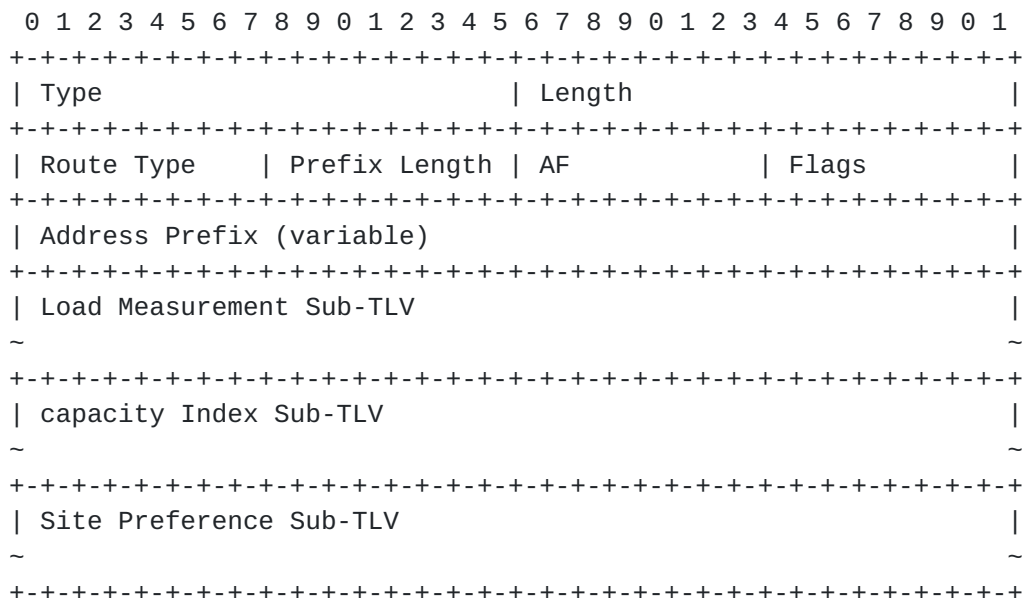


Figure 4: App-Metrix Sub-TLVs in OSPFv2 Extended Prefix TLV

5.3. IP Layer App-Metrics Sub-TLVs

Two types of Load Measurement Sub-TLVs are specified:

- a) The Aggregated Load Index based on a weighted combination of the collected measurements;
- b) The raw measurements of packets/bytes to/from the App Server address. The raw measurement is useful when the egress routers cannot be configured with a consistent algorithm to compute the aggregated load index or the raw measurements are needed by a central analytic system.

The Aggregated Load Index Sub-TLV has the following format:

```

0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           Type (TBD2)           |           Length           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           Measurement Period           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           Aggregated Load Index to reach the App Server           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Figure 5: Aggregated Load Index Sub-TLV

Type=TBD2 (to be assigned by IANA) indicates that the sub-TLV carries the Aggregated Load Measurement Index derived from the Weighted combination of bytes/packets sent to/received from the App server:

$$\text{Index} = w1 * \text{ToPackets} + w2 * \text{FromPackets} + w3 * \text{ToBytes} + w4 * \text{FromBytes}$$

Where w_i is a value between 0 and 1; $w1 + w2 + w3 + w4 = 1$.

The Raw Load Measurement sub-TLV has the following format:

```

0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|          Type (TBD3)          |          Length          |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|          Measurement Period          |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| total number of packets to the AppServer |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| total number of packets from the AppServer |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| total number of bytes to the AppServer |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| total number of bytes from the AppServer |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Figure 6: Raw Load Measurement Sub-TLV

Type= TBD3 (to be assigned by IANA) indicates that the sub-TLV carries the Raw measurements of packets/bytes to/from the App Server ANYCAST address.

Measurement Period: A user-specified period in seconds, default is 3600 seconds.

The Capacity Index sub-TLV has the following format:

```

0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|          Type (TBD3)          |          Length          |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|          Capacity Index          |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Figure 7: Capacity Index Sub-TLV

The Preference Index sub-TLV has the following format:

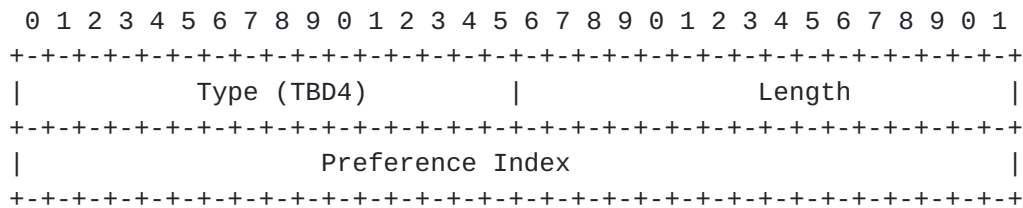


Figure 8: Preference Index Sub-TLV

Note: "Capacity Index" and "Site preference" can be more stable for each site. If those values are configured to nodes, they might not need to be included in every OSPF LSA.

6. Manageability Considerations

To be added.

7. Security Considerations

To be added.

8. IANA Considerations

The following Sub-TLV types need to be added by IANA to OSPFv4 Extended-LSA Sub-TLVs and OSPFv2 Extended Link Opaque LSA TLVs Registry.

- Aggregated Load Index Sub-TLV type
- Raw Load Measurement Sub-TLV type
- Capacity Index Sub-TLV type
- Preference Index Sub-TLV type

9. References

9.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2328] J. Moy, "OSPF Version 2", [RFC 2328](#), April 1998.
- [RFC7684] P. Psenak, et al, "OSPFv2 Prefix/Link Attribute Advertisement", [RFC 7684](#), Nov. 2015.
- [RFC8200] S. Deering R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", July 2017.
- [RFC8326] A. Lindem, et al, "OSPFv3 Link State advertisement (LSA0 Extensibility", [RFC 8362](#), April 2018.

9.2. Informative References

- [3GPP-EdgeComputing] 3GPP TR 23.748, "3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; Study on enhancement of support for Edge Computing in 5G Core network (5GC)", Release 17 work in progress, Aug 2020.
- [5G-AppMetaData] L. Dunbar, K. Majumdar, H. Wang, "BGP NLRI App Meta Data for 5G Edge Computing Service", [draft-dunbar-idr-5g-edge-compute-app-meta-data-01](#), work-in-progress, Nov 2020.
- [5G-EC-Metrics] L. Dunbar, H. Song, J. Kaippallimalil, "IP Layer Metrics for 5G Edge Computing Service", [draft-dunbar-ippm-5g-edge-compute-ip-layer-metrics-01](#), work-in-progress, Nov 2020.
- [5G-StickyService] L. Dunbar, J. Kaippallimalil, "IPv6 Solution for 5G Edge Computing Sticky Service", [draft-dunbar-6man-5g-ec-sticky-service-00](#), work-in-progress, Oct 2020.

- [RFC5521] P. Mohapatra, E. Rosen, "The BGP Encapsulation Subsequent Address Family Identifier (SAFI) and the BGP Tunnel Encapsulation Attribute", April 2009.
- [BGP-SDWAN-Port] L. Dunbar, H. Wang, W. Hao, "BGP Extension for SDWAN Overlay Networks", [draft-dunbar-idr-bgp-sdwan-overlay-ext-03](#), work-in-progress, Nov 2018.
- [SDWAN-EDGE-Discovery] L. Dunbar, S. Hares, R. Raszuk, K. Majumdar, "BGP UPDATE for SDWAN Edge Discovery", [draft-dunbar-idr-sdwan-edge-discovery-00](#), work-in-progress, July 2020.
- [Tunnel-Encap] E. Rosen, et al "The BGP Tunnel Encapsulation Attribute", [draft-ietf-idr-tunnel-encaps-10](#), Aug 2018.

10. Acknowledgments

Acknowledgements to Acee Lindem, Gyan Mishra, Jeff Tantsura, and Donald Eastlake for their review and suggestions.

This document was prepared using 2-Word-v2.0.template.dot.

Authors' Addresses

Linda Dunbar
Futurewei
Email: ldunbar@futurewei.com

Huaimo Chen
Futurewei
Email: huaimo.chen@futurewei.com

Aijun Wang
China Telecom
Email: wangaj3@chinatelecom.cn