

Networking Working Group
Internet-Draft
Intended status: Standards Track
Expires: September 7, 2019

P. Psenak, Ed.
C. Filsfils
Cisco Systems
A. Bashandy
Individual
B. Decraene
Orange
Z. Hu
Huawei Technologies
March 6, 2019

IS-IS Extensions to Support Routing over IPv6 Dataplane
draft-bashandy-isis-srv6-extensions-05.txt

Abstract

Segment Routing (SR) allows for a flexible definition of end-to-end paths by encoding paths as sequences of topological sub-paths, called "segments". Segment routing architecture can be implemented over an MPLS data plane as well as an IPv6 data plane. This draft describes the IS-IS extensions required to support Segment Routing over an IPv6 data plane.

Requirements Language

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.

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 <https://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 September 7, 2019.

Copyright Notice

Copyright (c) 2019 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 (<https://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
2. SRv6 Capabilities sub-TLV	3
3. Advertising Supported Algorithms	4
4. Advertising Maximum SRv6 SID Depths	4
4.1. Maximum Segments Left MSD Type	5
4.2. Maximum End Pop MSD Type	5
4.3. Maximum T.Insert MSD Type	5
4.4. Maximum T.Encaps MSD Type	5
4.5. Maximum End D MSD Type	6
5. SRv6 SIDs and Reachability	6
6. Advertising Locators and End SIDs	7
6.1. SRv6 Locator TLV Format	8
6.2. SRv6 End SID sub-TLV	9
7. Advertising SRv6 End.X SIDs	11
7.1. SRv6 End.X SID sub-TLV	11
7.2. SRv6 LAN End.X SID sub-TLV	13
8. Advertising Endpoint Behaviors	14
9. IANA Considerations	15
9.1. SRv6 Locator TLV	15
9.1.1. SRv6 End SID sub-TLV	15
9.1.2. Revised sub-TLV table	16
9.2. SRv6 Capabilities sub-TLV	16
9.3. SRv6 End.X SID and SRv6 LAN End.X SID sub-TLVs	17
9.4. MSD Types	17
10. Security Considerations	17
11. Contributors	17
12. References	18
12.1. Normative References	18
12.2. Informative References	20
Authors' Addresses	21

1. Introduction

With Segment Routing (SR) [I-D.ietf-spring-segment-routing], a node steers a packet through an ordered list of instructions, called segments.

Segments are identified through Segment Identifiers (SIDs).

Segment Routing can be directly instantiated on the IPv6 data plane through the use of the Segment Routing Header defined in [I-D.ietf-6man-segment-routing-header]. SRv6 refers to this SR instantiation on the IPv6 dataplane.

The network programming paradigm [I-D.filsfils-spring-srv6-network-programming] is central to SRv6. It describes how any function can be bound to a SID and how any network program can be expressed as a combination of SID's.

This document specifies IS-IS extensions that allow the IS-IS protocol to encode some of these functions.

Familiarity with the network programming paradigm [I-D.filsfils-spring-srv6-network-programming] is necessary to understand the extensions specified in this document.

This document defines one new top level IS-IS TLV and several new IS-IS sub-TLVs.

The SRv6 Capabilities sub-TLV announces the ability to support SRv6 and some Endpoint functions listed in Section 7 as well as advertising limitations when applying such Endpoint functions.

The SRv6 Locator top level TLV announces SRv6 locators - a form of summary address for the set of topology/algorithm specific SIDs associated with a node.

The SRv6 End SID sub-TLV, the SRv6 End.X SID sub-TLV, and the SRv6 LAN End.X SID sub-TLV are used to advertise which SIDs are instantiated at a node and what Endpoint function is bound to each instantiated SID.

2. SRv6 Capabilities sub-TLV

A node indicates that it has support for SRv6 by advertising a new SRv6- capabilities sub-TLV of the router capabilities TLV [RFC7981].

The SRv6 Capabilities sub-TLV may contain optional sub-sub-TLVs. No sub-sub-TLVs are currently defined.

The SRv6 Capabilities sub-TLV has the following format:

```

0                               1                               2                               3
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      |      Length      |      Flags      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| optional sub-sub-TLVs... |

```

Type: Suggested value 25, to be assigned by IANA

Length: 2 + length of sub-sub-TLVs

Flags: 2 octets The following flags are defined:

```

0                               1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| |O| | | | | | | | | | | | | | | | | | | | | | | | | | | |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

O-flag: If set, the router supports use of the O-bit in the Segment Routing Header(SRH) as defined in [I-D.ali-spring-srv6-oam].

3. Advertising Supported Algorithms

SRv6 capable router indicates supported algorithm(s) by advertising the SR Algorithm TLV as defined in [I-D.ietf-isis-segment-routing-extensions].

4. Advertising Maximum SRv6 SID Depths

[I-D.ietf-isis-segment-routing-msd] defines the means to advertise node/link specific values for Maximum SID Depths (MSD) of various types. Node MSDs are advertised in a sub-TLV of the Router Capabilities TLV [RFC7981]. Link MSDs are advertised in a sub-TLV of TLVs 22, 23, 141, 222, and 223.

This document defines the relevant SRv6 MSDs and requests MSD type assignments in the MSD Types registry created by [I-D.ietf-isis-segment-routing-msd].

4.1. Maximum Segments Left MSD Type

The Maximum Segments Left MSD Type specifies the maximum value of the "SL" field [I-D.ietf-6man-segment-routing-header] in the SRH of a received packet before applying the Endpoint function associated with a SID.

SRH Max SL Type: 41 (Suggested value - to be assigned by IANA)

If no value is advertised the supported value is assumed to be 0.

4.2. Maximum End Pop MSD Type

The Maximum End Pop MSD Type specifies the maximum number of SIDs in the top SRH in an SRH stack to which the router can apply "PSP" or USP" as defined in [I-D.filsfils-spring-srv6-network-programming] flavors.

SRH Max End Pop Type: 42 (Suggested value - to be assigned by IANA)

If the advertised value is zero or no value is advertised then it is assumed that the router cannot apply PSP or USP flavors.

4.3. Maximum T.Insert MSD Type

The Maximum T.Insert MSD Type specifies the maximum number of SIDs that can be inserted as part of the "T.insert" behavior as defined in [I-D.filsfils-spring-srv6-network-programming].

SRH Max T.insert Type: 43 (Suggested value - to be assigned by IANA)

If the advertised value is zero or no value is advertised then the router is assumed not to support any variation of the "T.insert" behavior.

4.4. Maximum T.Encaps MSD Type

The Maximum T.Encaps MSD Type specifies the maximum number of SIDs that can be included as part of the "T.Encaps" behavior as defined in [I-D.filsfils-spring-srv6-network-programming] .

SRH Max T.encaps Type: 44 (Suggested value - to be assigned by IANA)

If the advertised value is zero then the router can apply T.Encaps only by encapsulating the incoming packet in another IPv6 header without SRH the same way IPinIP encapsulation is performed.

If the advertised value is non-zero then the router supports both IPinIP and SRH encapsulation subject to the SID limitation specified by the advertised value.

4.5. Maximum End D MSD Type

The Maximum End D MSD Type specifies the maximum number of SIDs in an SRH when performing decapsulation associated with "End.Dx" functions (e.g., "End.DX6" and "End.DT6") as defined in [I-D.filsfils-spring-srv6-network-programming].

SRH Max End D Type: 45 (Suggested value - to be assigned by IANA)

If the advertised value is zero or no value is advertised then it is assumed that the router cannot apply "End.DX6" or "End.DT6" functions if the extension header right underneath the outer IPv6 header is an SRH.

5. SRv6 SIDs and Reachability

As discussed in [I-D.filsfils-spring-srv6-network-programming], an SRv6 Segment Identifier (SID) is 128 bits and represented as

LOC:FUNCT

where LOC (the locator portion) is the L most significant bits and FUNCT is the 128-L least significant bits. L is called the locator length and is flexible. Each operator is free to use the locator length it chooses.

A node is provisioned with topology/algorithm specific locators for each of the topology/algorithm pairs supported by that node. Each locator is a covering prefix for all SIDs provisioned on that node which have the matching topology/algorithm.

Locators MUST be advertised in the SRv6 Locator TLV (see Section 6.1). Forwarding entries for the locators advertised in the SRv6 Locator TLV MUST be installed in the forwarding plane of receiving SRv6 capable routers when the associated topology/algorithm is supported by the receiving node.

Locators are routable and MAY also be advertised in Prefix Reachability TLVs (236 or 237).

Locators associated with algorithm 0 (for all supported topologies) SHOULD be advertised in a Prefix Reachability TLV (236 or 237) so that legacy routers (i.e., routers which do NOT support SRv6) will install a forwarding entry for algorithm 0 SRv6 traffic.

In cases where a locator advertisement is received in both in a Prefix Reachability TLV and an SRv6 Locator TLV, the Prefix Reachability advertisement MUST be preferred when installing entries in the forwarding plane. This is to prevent inconsistent forwarding entries on SRv6 capable/SRv6 incapable routers.

SRv6 SIDs are advertised as sub-TLVs in the SRv6 Locator TLV except for SRv6 End.X SIDs/LAN End.X SIDs which are associated with a specific Neighbor/Link and are therefore advertised as sub-TLVs in TLVs 22, 23, 222, 223, and 141.

SRv6 SIDs are not directly routable and MUST NOT be installed in the forwarding plane. Reachability to SRv6 SIDs depends upon the existence of a covering locator.

Adherence to the rules defined in this section will assure that SRv6 SIDs associated with a supported topology/algorithm pair will be forwarded correctly, while SRv6 SIDs associated with an unsupported topology/algorithm pair will be dropped. NOTE: The drop behavior depends on the absence of a default/summary route covering a given locator.

In order for forwarding to work correctly, the locator associated with SRv6 SID advertisements MUST be the longest match prefix installed in the forwarding plane for those SIDs. There are a number of ways in which this requirement could be compromised

- o Another locator associated with a different topology/algorithm is the longest match
- o A prefix advertisement (i.e., from TLV 236 or 237) is the longest match

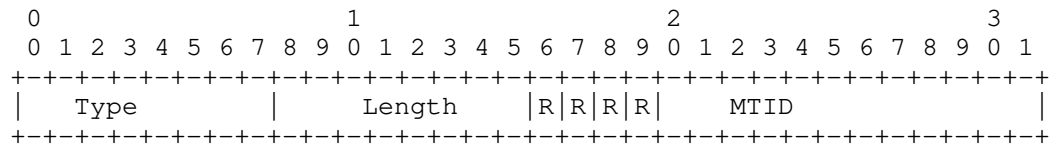
6. Advertising Locators and End SIDs

The SRv6 Locator TLV is introduced to advertise SRv6 Locators and End SIDs associated with each locator.

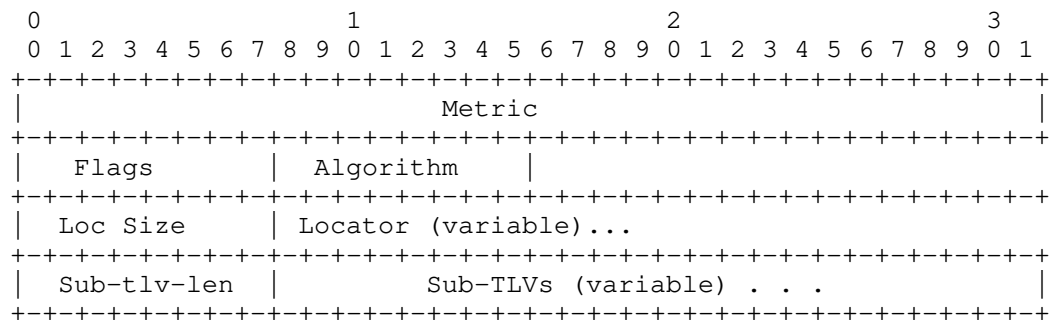
This new TLV shares the sub-TLV space defined for TLVs 135, 235, 236 and 237.

6.1. SRv6 Locator TLV Format

The SRv6 Locator TLV has the following format:



Followed by one or more locator entries of the form:



Type: 27 (Suggested value to be assigned by IANA)

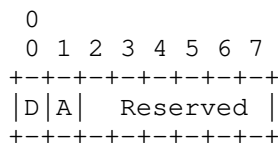
Length: variable.

MTID: Multitopology Identifier as defined in [RFC5120].
Note that the value 0 is legal.

Locator entry:

Metric: 4 octets. As described in [RFC5305].

Flags: 1 octet. The following flags are defined



where:

D bit: When the Locator is leaked from level-2 to level-1, the D bit MUST be set. Otherwise, this bit MUST be clear. Locators with the D bit set MUST NOT be leaked from level-1 to level-2.

This is to prevent looping.

A bit: When the Locator is configured as anycast, the A bit SHOULD be set. Otherwise, this bit MUST be clear.

The remaining bits are reserved for future use. They SHOULD be set to zero on transmission and MUST be ignored on receipt.

Algorithm: 1 octet. Associated algorithm. Algorithm values are defined in the IGP Algorithm Type registry.

Loc-Size: 1 octet. Number of bits in the Locator field.
(1 - 128)

Locator: 1-16 octets. This field encodes the advertised SRv6 Locator. The Locator is encoded in the minimal number of octets for the given number of bits.

Sub-TLV-length: 1 octet. Number of octets used by sub-TLVs

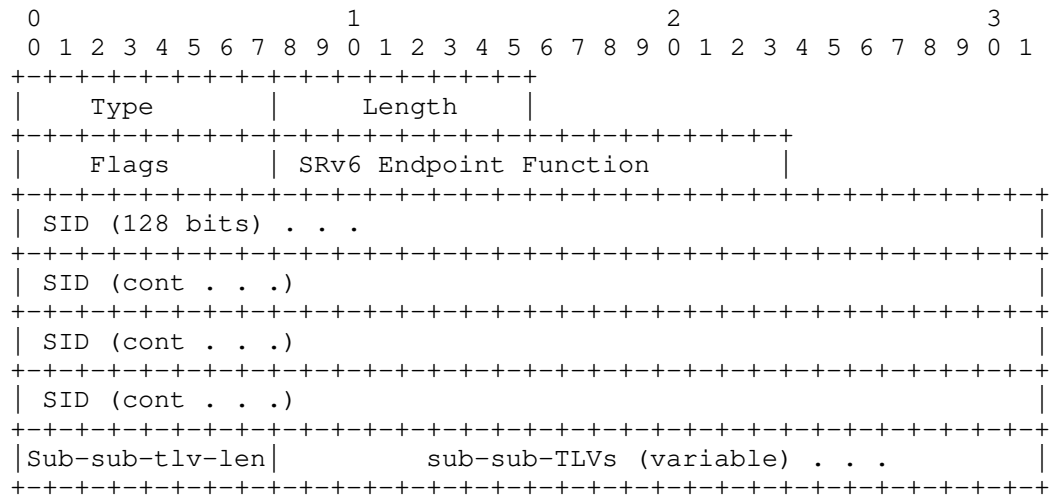
Optional sub-TLVs.

6.2. SRv6 End SID sub-TLV

The SRv6 End SID sub-TLV is introduced to advertise SRv6 Segment Identifiers (SID) with Endpoint functions which do not require a particular neighbor in order to be correctly applied [I-D.filsfils-spring-srv6-network-programming]. SRv6 SIDs associated with a neighbor are advertised using the sub-TLVs defined in Section 6.

This new sub-TLV is advertised in the SRv6 Locator TLV defined in the previous section. SRv6 End SIDs inherit the topology/algorithm from the parent locator.

The SRv6 End SID sub-TLV has the following format:



Type: 5 (Suggested value to be assigned by IANA)

Length: variable.

Flags: 1 octet. No flags are currently defined.

SRv6 Endpoint Function: 2 octets. As defined in
 [I-D.filsfils-spring-srv6-network-programming]
 Legal function values for this sub-TLV are defined in Section 7.

SID: 16 octets. This field encodes the advertised SRv6 SID.

Sub-sub-TLV-length: 1 octet. Number of octets used by sub-sub-TLVs

Optional sub-sub-TLVs

The SRv6 End SID MUST be a subnet of the associated Locator. SRv6 End SIDs which are NOT a subnet of the associated locator MUST be ignored.

Multiple SRv6 End SIDs MAY be associated with the same locator. In cases where the number of SRv6 End SID sub-TLVs exceeds the capacity of a single TLV, multiple Locator TLVs for the same locator MAY be advertised. For a given MTID/Locator the algorithm MUST be the same in all TLVs. If this restriction is not met all TLVs for that MTID/Locator MUST be ignored.

7. Advertising SRv6 End.X SIDs

Certain SRv6 Endpoint functions

[I-D.filsfils-spring-srv6-network-programming] must be associated with a particular neighbor, and in case of multiple layer 3 links to the same neighbor, with a particular link in order to be correctly applied.

This document defines two new sub-TLVs of TLV 22, 23, 222, 223, and 141 - namely "SRv6 End.X SID" and "SRv6 LAN End.X SID".

IS-IS Neighbor advertisements are topology specific - but not algorithm specific. End.X SIDs therefore inherit the topology from the associated neighbor advertisement, but the algorithm is specified in the individual SID.

All End.X SIDs MUST be a subnet of a Locator with matching topology and algorithm which is advertised by the same node in an SRv6 Locator TLV. End.X SIDs which do not meet this requirement MUST be ignored.

7.1. SRv6 End.X SID sub-TLV

This sub-TLV is used to advertise an SRv6 SID associated with a point to point adjacency. Multiple SRv6 End.X SID sub-TLVs MAY be associated with the same adjacency.

The SRv6 End.X SID sub-TLV has the following format:

```

      0                               1                               2                               3
      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      |      Length      |                                     |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  Flags     |      Algorithm   |      Weight      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  SRv6 Endpoint Function  |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  SID (128 bits) . . .   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  SID (cont . . .)       |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  SID (cont . . .)       |
+-----+-----+-----+-----+-----+-----+-----+-----+
|  SID (cont . . .)       |
+-----+-----+-----+-----+-----+-----+-----+-----+
|Sub-sub-tlv-len|      Sub-sub-TLVs (variable) . . .   |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Type: 43 (Suggested value to be assigned by IANA)

Length: variable.

Flags: 1 octet.

```

    0 1 2 3 4 5 6 7
    +-+-+-+-+-+-+-+-+
    |B|S|P|Reserved |
    +-+-+-+-+-+-+-+-+

```

where:

B-Flag: Backup flag. If set, the End.X SID is eligible for protection (e.g., using IPFRR) as described in [RFC8355].

S-Flag. Set flag. When set, the S-Flag indicates that the End.X SID refers to a set of adjacencies (and therefore MAY be assigned to other adjacencies as well).

P-Flag. Persistent flag. When set, the P-Flag indicates that the End.X SID is persistently allocated, i.e., the End.X SID value remains consistent across router restart and/or interface flap.

Other bits: MUST be zero when originated and ignored when received.

Algorithm: 1 octet. Associated algorithm. Algorithm values are defined in the IGP Algorithm Type registry.

Weight: 1 octet. The value represents the weight of the End.X SID for the purpose of load balancing. The use of the weight is defined in [I-D.ietf-spring-segment-routing].

SRv6 Endpoint Function: 2 octets. As defined in [I-D.filsfils-spring-srv6-network-programming]
Legal function values for this sub-TLV are defined in Section 7.

SID: 16 octets. This field encodes the advertised SRv6 SID.

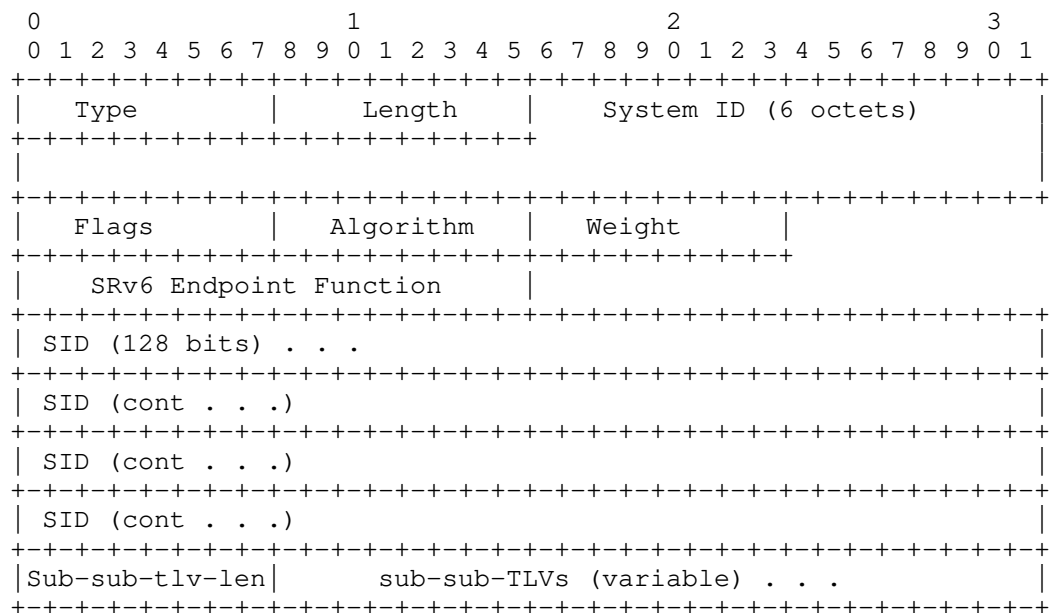
Sub-sub-TLV-length: 1 octet. Number of octets used by sub-sub-TLVs

Note that multiple TLVs for the same neighbor may be required in order to advertise all of the SRv6 End.X SIDs associated with that neighbor.

7.2. SRv6 LAN End.X SID sub-TLV

This sub-TLV is used to advertise an SRv6 SID associated with a LAN adjacency. Since the parent TLV is advertising an adjacency to the Designated Intermediate System(DIS) for the LAN, it is necessary to include the System ID of the physical neighbor on the LAN with which the SRv6 SID is associated. Given that a large number of neighbors may exist on a given LAN a large number of SRv6 LAN END.X SID sub-TLVs may be associated with the same LAN. Note that multiple TLVs for the same DIS neighbor may be required in order to advertise all of the SRv6 End.X SIDs associated with that neighbor.

The SRv6 LAN End.X SID sub-TLV has the following format:

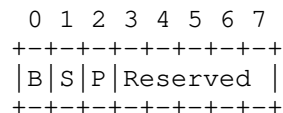


Type: 44 (Suggested value to be assigned by IANA)

Length: variable.

System-ID: 6 octets of IS-IS System-ID of length "ID Length" as defined in [ISO10589].

Flags: 1 octet.



where B,S, and P flags are as described in Section 6.1.
Other bits: MUST be zero when originated and ignored when received.

Algorithm: 1 octet. Associated algorithm. Algorithm values are defined in the IGP Algorithm Type registry.

Weight: 1 octet. The value represents the weight of the End.X SID for the purpose of load balancing. The use of the weight is defined in [I-D.ietf-spring-segment-routing].

SRv6 Endpoint Function: 2 octets. As defined in [I-D.filsfils-spring-srv6-network-programming]
Legal function values for this sub-TLV are defined in Section 7.

SID: 16 octets. This field encodes the advertised SRv6 SID.

Sub-sub-TLV-length: 1 octet. Number of octets used by sub-sub-TLVs.

8. Advertising Endpoint Behaviors

Endpoint behaviors are defined in [I-D.filsfils-spring-srv6-network-programming] and [I-D.ali-spring-srv6-oam]. The numerical identifiers for the Endpoint behaviors are defined in the "SRv6 Endpoint Behaviors" registry defined in [I-D.filsfils-spring-srv6-network-programming]. This section lists the Endpoint behaviors and their identifiers, which MAY be advertised by IS-IS and the SID sub-TLVs in which each type MAY appear.

Endpoint Behavior	Endpoint Behavior Identifier	End SID	End.X SID	Lan End.X SID
End (PSP, USP, USD)	1-4, 28-31	Y	N	N
End.X (PSP, USP, USD)	5-8, 32-35	N	Y	Y
End.T (PSP, USP, USD)	9-12, 36-39	Y	N	N
End.DX6	16	N	Y	Y
End.DX4	17	N	Y	Y
End.DT6	18	Y	N	N
End.DT4	19	Y	N	N
End.DT64	20	Y	N	N
End.OP	40	Y	N	N
End.OTP	41	Y	N	N

9. IANA Considerations

This document requests allocation for the following TLVs, sub-TLVs, and sub-sub-TLVs as well updating the ISIS TLV registry and defining a new registry.

9.1. SRv6 Locator TLV

This document adds one new TLV to the IS-IS TLV Codepoints registry.

Value: 27 (suggested - to be assigned by IANA)

Name: SRv6 Locator

This TLV shares sub-TLV space with existing "Sub-TLVs for TLVs 135, 235, 236 and 237 registry". The name of this registry needs to be changed to "Sub-TLVs for TLVs 27, 135, 235, 236 and 237 registry".

9.1.1. SRv6 End SID sub-TLV

This document adds the following new sub-TLV to the (renamed) "Sub-TLVs for TLVs 27, 135, 235, 236 and 237 registry".

Value: 5 (suggested - to be assigned by IANA)

Name: SRv6 End SID

This document requests the creation of a new IANA managed registry for sub-sub-TLVs of the SRv6 End SID sub-TLV. The registration procedure is "Expert Review" as defined in [RFC7370]. Suggested registry name is "sub-sub-TLVs for SRv6 End SID sub-TLV". No sub-sub-TLVs are defined by this document except for the reserved value.

0: Reserved

1-255: Unassigned

9.1.2. Revised sub-TLV table

The revised table of sub-TLVs for the (renamed) "Sub-TLVs for TLVs 27, 135, 235, 236 and 237 registry" is shown below:

Type	27	135	235	236	237
1	n	y	y	y	y
2	n	y	y	y	y
3	n	y	y	y	y
4	y	y	y	y	y
5	y	n	n	n	n
11	y	y	y	y	y
12	y	y	y	y	y

9.2. SRv6 Capabilities sub-TLV

This document adds the definition of a new sub-TLV in the "Sub- TLVs for TLV 242 registry".

Type: 25 (Suggested - to be assigned by IANA)

Description: SRv6 Capabilities

This document requests the creation of a new IANA managed registry for sub-sub-TLVs of the SRv6 Capability sub-TLV. The registration procedure is "Expert Review" as defined in [RFC7370]. Suggested registry name is "sub-sub-TLVs for SRv6 Capability sub-TLV". No sub-sub-TLVs are defined by this document except for the reserved value.

0: Reserved

1-255: Unassigned

9.3. SRv6 End.X SID and SRv6 LAN End.X SID sub-TLVs

This document adds the definition of two new sub-TLVs in the "sub-TLVs for TLV 22, 23, 25, 141, 222 and 223 registry".

Type: 43 (suggested - to be assigned by IANA)

Description: SRv6 End.X SID

Type: 44 (suggested - to be assigned by IANA)

Description: SRv6 LAN End.X SID

Type	22	23	25	141	222	223
43	Y	Y	Y	Y	Y	Y
44	Y	Y	Y	Y	Y	Y

9.4. MSD Types

This document defines the following new MSD types. These types are to be defined in the IGP MSD Types registry defined in [I-D.ietf-isis-segment-routing-msd] .

All values are suggested values to be assigned by IANA.

Type	Description
41	SRH Max SL
42	SRH Max End Pop
43	SRH Max T.insert
44	SRH Max T.encaps
45	SRH Max End D

10. Security Considerations

Security concerns for IS-IS are addressed in [ISO10589], [RFC5304], and [RFC5310].

11. Contributors

The following people gave a substantial contribution to the content of this document and should be considered as co-authors:

Stefano Previdi
Huawei Technologies
Email: stefano@previdi.net

Paul Wells
Cisco Systems
Saint Paul,
Minnesota
United States
Email: pauwells@cisco.com

Daniel Voyer
Email: daniel.voyer@bell.ca

Satoru Matsushima
Email: satoru.matsushima@g.softbank.co.jp

Bart Peirens
Email: bart.peirens@proximus.com

Hani Elmalky
Email: hani.elmalky@ericsson.com

Prem Jonnalagadda
Email: prem@barefootnetworks.com

Milad Sharif
Email: msharif@barefootnetworks.com>

Robert Hanzl
Cisco Systems
Millenium Plaza Building, V Celnici 10, Prague 1,
Prague, Czech Republic
Email rhanzl@cisco.com

Ketan Talaulikar
Cisco Systems, Inc.
Email: ketant@cisco.com

12. References

12.1. Normative References

[I-D.ali-spring-srv6-oam]

Ali, Z., Filsfils, C., Kumar, N., Pignataro, C.,
faiqbal@cisco.com, f., Gandhi, R., Leddy, J., Matsushima,
S., Raszuk, R., daniel.voyer@bell.ca, d., Dawra, G.,
Peirens, B., Chen, M., and G. Naik, "Operations,
Administration, and Maintenance (OAM) in Segment Routing
Networks with IPv6 Data plane (SRv6)", draft-ali-spring-
srv6-oam-02 (work in progress), October 2018.

[I-D.filsfils-spring-srv6-network-programming]

Filsfils, C., Camarillo, P., Leddy, J.,
daniel.voyer@bell.ca, d., Matsushima, S., and Z. Li, "SRv6
Network Programming", draft-filsfils-spring-srv6-network-
programming-07 (work in progress), February 2019.

[I-D.ietf-6man-segment-routing-header]

Filsfils, C., Previdi, S., Leddy, J., Matsushima, S., and
d. daniel.voyer@bell.ca, "IPv6 Segment Routing Header
(SRH)", draft-ietf-6man-segment-routing-header-16 (work in
progress), February 2019.

[I-D.ietf-isis-segment-routing-extensions]

Previdi, S., Ginsberg, L., Filsfils, C., Bashandy, A.,
Gredler, H., and B. Decraene, "IS-IS Extensions for
Segment Routing", draft-ietf-isis-segment-routing-
extensions-22 (work in progress), December 2018.

[I-D.ietf-isis-segment-routing-msd]

Tantsura, J., Chunduri, U., Aldrin, S., and L. Ginsberg,
"Signaling MSD (Maximum SID Depth) using IS-IS", draft-
ietf-isis-segment-routing-msd-19 (work in progress),
October 2018.

[ISO10589]

Standardization", I. "O. F., "Intermediate system to
Intermediate system intra-domain routing information
exchange protocol for use in conjunction with the protocol
for providing the connectionless-mode Network Service (ISO
8473), ISO/IEC 10589:2002, Second Edition.", Nov 2002.

[RFC2119]

Bradner, S., "Key words for use in RFCs to Indicate
Requirement Levels", BCP 14, RFC 2119,
DOI 10.17487/RFC2119, March 1997,
<<http://www.rfc-editor.org/info/rfc2119>>.

- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC7370] Ginsberg, L., "Updates to the IS-IS TLV Codepoints Registry", RFC 7370, DOI 10.17487/RFC7370, September 2014, <<https://www.rfc-editor.org/info/rfc7370>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

12.2. Informative References

- [I-D.ietf-spring-segment-routing]
Filsfils, C., Previdi, S., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", draft-ietf-spring-segment-routing-15 (work in progress), January 2018.
- [RFC8355] Filsfils, C., Ed., Previdi, S., Ed., Decraene, B., and R. Shakir, "Resiliency Use Cases in Source Packet Routing in Networking (SPRING) Networks", RFC 8355, DOI 10.17487/RFC8355, March 2018, <<https://www.rfc-editor.org/info/rfc8355>>.

Authors' Addresses

Peter Psenak (editor)
Cisco Systems
Pribinova Street 10
Bratislava 81109
Slovakia

Email: ppsenak@cisco.com

Clarence Filsfils
Cisco Systems
Brussels
Belgium

Email: cfilsfil@cisco.com

Ahmed Bashandy
Individual

Email: abashandy.ietf@gmail.com

Bruno Decraene
Orange
Issy-les-Moulineaux
France

Email: bruno.decraene@orange.com

Zhibo Hu
Huawei Technologies

Email: huzhibo@huawei.com

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: March 24, 2019

H. Chen
D. Cheng
Huawei Technologies
M. Toy
Verizon
Y. Yang
IBM
September 20, 2018

LS Flooding Reduction
draft-cc-ospf-flooding-reduction-04

Abstract

This document proposes an approach to flood link states on a topology that is a subgraph of the complete topology per underline physical network, so that the amount of flooding traffic in the network is greatly reduced, and it would reduce convergence time with a more stable and optimized routing environment. The approach can be applied to any network topology in a single area.

Requirements Language

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 [RFC2119].

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 <https://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 March 24, 2019.

Copyright Notice

Copyright (c) 2018 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 (<https://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
2. Terminology	3
3. Conventions Used in This Document	4
4. Problem Statement	4
5. Flooding Topology	5
5.1. Construct Flooding Topology	5
5.2. Backup for Flooding Topology Split	7
6. Extensions to OSPF	7
6.1. Extensions for Operations	8
6.2. Extensions for Centralized Mode	9
6.2.1. Message for Flooding Topology	9
6.2.2. Encodings for Backup Paths	16
6.2.3. Message for Incremental Changes	24
6.2.4. Leaders Selection	25
7. Extensions to IS-IS	27
7.1. Extensions for Operations	27
7.2. Extensions for Centralized Mode	27
7.2.1. TLV for Flooding Topology	27
7.2.2. Encodings for Backup Paths	28
7.2.3. TLVs for Incremental Changes	29
7.2.4. Leaders Selection	30
8. Flooding Behavior	30
8.1. Nodes Perform Flooding Reduction without Failure	30
8.1.1. Receiving an LS	30
8.1.2. Originating an LS	31
8.1.3. Establishing Adjacencies	31
8.2. An Exception Case	32
8.2.1. A Critical Failure	32
8.2.2. Multiple Failures	32
9. Security Considerations	33
10. IANA Considerations	33

10.1.	OSPFv2	33
10.2.	OSPFv3	35
10.3.	IS-IS	36
11.	Acknowledgements	36
12.	References	36
12.1.	Normative References	36
12.2.	Informative References	37
Appendix A.	Algorithms to Build Flooding Topology	37
A.1.	Algorithms to Build Tree without Considering Others	37
A.2.	Algorithms to Build Tree Considering Others	39
A.3.	Connecting Leaves	41
Authors' Addresses	42

1. Introduction

For some networks such as dense Data Center (DC) networks, the existing Link State (LS) flooding mechanism is not efficient and may have some issues. The extra LS flooding consumes network bandwidth. Processing the extra LS flooding, including receiving, buffering and decoding the extra LSs, wastes memory space and processor time. This may cause scalability issues and affect the network convergence negatively.

This document proposes an approach to minimize the amount of flooding traffic in the network. Thus the workload for processing the extra LS flooding is decreased significantly. This would improve the scalability, speed up the network convergence, stable and optimize the routing environment.

This approach is also flexible. It has multiple modes for computation of flooding topology. Users can select a mode they prefer, and smoothly switch from one mode to another. The approach is applicable to any network topology in a single area. It is backward compatible.

2. Terminology

Flooding Topology:

A sub-graph or sub-network of a given (physical) network topology that has the same reachability to every node as the given network topology, through which link states are flooded.

critical link or interface on a flooding topology:

A only link or interface among some nodes on the flooding topology. When this link or interface goes down, the flooding topology will be split.

critical node on a flooding topology:

A only node connecting some nodes on the flooding topology. When this node goes down, the flooding topology will be split.

backup path:

A path or a sequence of links, when a critical link or node goes down, providing a connection to connect two parts of a split flooding topology. When a critical node goes down, the flooding topology may be split into more than two parts. In this case, two or more backup paths are needed to connect all the split parts into one.

Remaining Flooding Topology:

A topology from a flooding topology by removing the failed links and nodes from the flooding topology.

LSA:

A Link State Advertisement in OSPF.

LSP:

A Link State Protocol Data Unit (PDU) in IS-IS.

LS:

A Link State, which is an LSA or LSP.

3. Conventions Used in This Document

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 [RFC2119].

4. Problem Statement

OSPF and IS-IS deploy a so-called reliable flooding mechanism, where a node must transmit a received or self-originated LS to all its interfaces (except the interface where an LS is received). While this mechanism assures each LS being distributed to every node in an area or domain, the side-effect is that the mechanism often causes redundant LS, which in turn forces nodes to process identical LS more than once. This results in the waste of link bandwidth and nodes' computing resources, and the delay of topology convergence.

This becomes more serious in networks with large number of nodes and links, and in particular, higher degree of interconnection (e.g., meshed topology, spine-leaf topology, etc.). In some environments such as in data centers, the drawback of the existing flooding mechanism has already caused operational issues, including repeated and waves of flooding storms, chock of computing resources, slow

convergence, oscillating topology changes, instability of routing environment.

One example is as shown in Figure 1, where Node 1, Node 2 and Node 3 are interconnected in a mesh. When Node 1 receives a new or updated LS on its interface I11, it by default would forward the LS to its interface I12 and I13 towards Node 2 and Node 3, respectively, after processing. Node 2 and Node 3 upon reception of the LS and after processing, would potentially flood the same LS over their respective interface I23 and I32 toward each other, which is obviously not necessary and at the cost of link bandwidth as well as both nodes' computing resource.

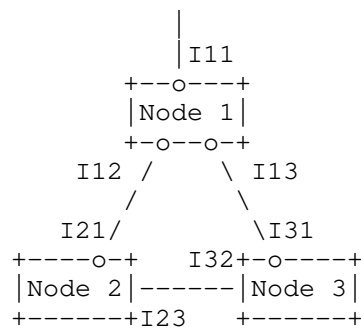


Figure 1

5. Flooding Topology

For a given network topology, a flooding topology is a sub-graph or sub-network of the given network topology that has the same reachability to every node as the given network topology. Thus all the nodes in the given network topology MUST be in the flooding topology. All the nodes MUST be inter-connected directly or indirectly. As a result, LS flooding will in most cases occur only on the flooding topology, that includes all nodes but a subset of links. Note even though the flooding topology is a sub-graph of the original topology, any single LS MUST still be disseminated in the entire network.

5.1. Construct Flooding Topology

Many different flooding topologies can be constructed for a given network topology. A chain connecting all the nodes in the given network topology is a flooding topology. A circle connecting all the nodes is another flooding topology. A tree connecting all the nodes is a flooding topology. In addition, the tree plus the connections

between some leaves of the tree and branch nodes of the tree is a flooding topology.

The following parameters need to be considered for constructing a flooding topology:

- o Number of links: The number of links on the flooding topology is a key factor for reducing the amount of LS flooding. In general, the smaller the number of links, the less the amount of LS flooding.
- o Diameter: The shortest distance between the two most distant nodes on the flooding topology is a key factor for reducing the network convergence time. The smaller the diameter, the less the convergence time.
- o Redundancy: The redundancy of the flooding topology means a tolerance to the failures of some links and nodes on the flooding topology. If the flooding topology is split by some failures, it is not tolerant to these failures. In general, the larger the number of links on the flooding topology is, the more tolerant the flooding topology to failures.

There are many different ways to construct a flooding topology for a given network topology. A few of them are listed below:

- o Central Mode: One node in the network builds a flooding topology and floods the flooding topology to all the other nodes in the network (This seems not good. Flooding the flooding topology may increase the flooding. The amount of traffic for flooding the flooding topology should be minimized.);
- o Distributed Mode: Each node in the network automatically calculates a flooding topology by using the same algorithm (No flooding for flooding topology);
- o Static Mode: Links on the flooding topology are configured statically.

Note that the flooding topology constructed by a node is dynamic in nature, that means when the base topology (the entire topology graph) changes, the flooding topology (the sub-graph) MUST be re-computed/re-constructed to ensure that any node that is reachable on the base topology MUST also be reachable on the flooding topology.

For reference purpose, some algorithms that allow nodes to automatically compute flooding topology are elaborated in Appendix A.

However, this document does not attempt to standardize how a flooding topology is established.

5.2. Backup for Flooding Topology Split

It is hard to construct a flooding topology that reduces the amount of LS flooding greatly and is tolerant to multiple failures. To get around this, we can compute and use backup paths for a critical link and node on the flooding topology. Using backup paths may also speed up convergence when the link and node fail.

When a critical link on the flooding topology fails, the flooding topology without the critical link (i.e., the remaining flooding topology) is split into two parts. A backup path for the critical link connects the two parts into one. Through the backup path and the remaining flooding topology, an LS can be flooded to every node in the network. The combination of the backup path and the flooding topology is tolerant to the failure of the critical link.

When a critical node on the flooding topology goes down, the flooding topology without the critical node and the links attached to the node (i.e., the remaining flooding topology) is split into two or more parts. One or more backup paths for the critical node connects the split parts into one. Through the backup paths and the remaining flooding topology, an LS can be flooded to every live node in the network. The combination of the backup paths and the flooding topology is tolerant to the failure of the critical node.

In addition to the backup paths for a critical link and node, backup paths for every non critical link and node on the flooding topology can be computed. When the failures of multiple links and nodes on the flooding topology happen, through the remaining flooding topology and the backup paths for these links and nodes, an LS can be flooded to every live node in the network. The combination of the backup paths and the flooding topology is tolerant to the failures of these links and nodes. If there are other failures that break the backup paths, an LS can be flooded to every live node by the traditional flooding procedure.

In a centralized mode, the leader computes the backup paths and floods them to all the other nodes. In a distributed mode, every node computes the backup paths.

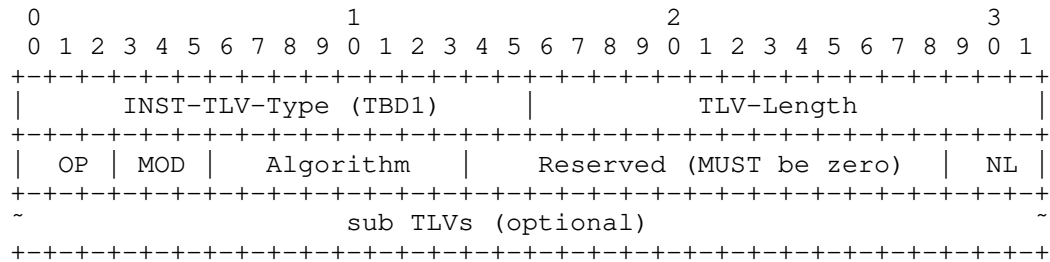
6. Extensions to OSPF

The extensions to OSPF comprises two parts: one part is for operations on flooding reduction, the other is specially for centralized mode flooding reduction.

6.1. Extensions for Operations

A new TLV is defined in OSPF RI LSA [RFC7770]. It contains instructions about flooding reduction, which is called Flooding Reduction Instruction TLV or Instruction TLV for short. This TLV is originated from only one node at any time.

The format of a Flooding Reduction Instruction TLV is as follows.



Flooding Reduction Instruction TLV

A OP field of three bits is defined in the TLV. It may have a value of the followings.

- o 0x001 (R): Perform flooding Reduction, which instructs the nodes in a network to perform flooding reduction.
- o 0x010 (N): Roll back to Normal flooding, which instructs the nodes in a network to roll back to perform normal flooding.

When any of the other values is received, it is ignored.

A MOD field of three bits is defined in the TLV and may have a value of the followings.

- o 0x001 (C): Central Mode, which instructs 1) the nodes in a network to select leaders (primary/designated leader, secondary/backup leader, and so on); 2) the leaders in a network to compute a flooding topology and the primary leader to flood the flooding topology to all the other nodes in the network; 3) every node in the network to receive and use the flooding topology originated by the primary leader.
- o 0x010 (D): Distributed Mode, which instructs every node in a network to compute and use its own flooding topology.

- o 0x011 (S): Static Mode, which instructs every node in a network to use the flooding topology statically configured on the node.

When any of the other values is received, it is ignored.

An Algorithm field of eight bits is defined in the TLV to instruct the leader node in central mode or every node in distributed mode to use the algorithm indicated in this field for computing a flooding topology.

A NL field of three bits is defined in the TLV, which indicates the number of leaders to be selected when Central Mode is used. NL set to 2 means two leaders (a designated/primary leader and a backup/secondary leader) to be selected for an area, and NL set to 3 means three leaders to be selected. When Central Mode is not used, The NL field is not valid.

Some optional sub TLVs may be defined in the future, but none is defined now.

6.2. Extensions for Centralized Mode

6.2.1. Message for Flooding Topology

A flooding topology can be represented by the links in the flooding topology. For the links between a local node and a number of its adjacent (or remote) nodes, we can encode the local node in a way, and encode its adjacent nodes in the same way or another way. After all the links in the flooding topology are encoded, the encoded links can be flooded to every node in the network. After receiving the encoded links, every node decodes the links and creates and/or updates the flooding topology.

For every node in an area, we may use an index to represent it. Every node in an area may order the nodes in a rule, which generates the same sequence of the nodes on every node in the area. The sequence of nodes have the index 0, 1, 2, and so on respectively. For example, every node orders the nodes by their router IDs in ascending order.

6.2.1.1. Links Encoding

A local node can be encoded in two parts: encoded node index size indication (ENSI) and compact node index (CNI). ENSI value plus a number (e.g., 9) gives the size of compact node index. For example, ENSI = 0 indicates that the size of CNIs is 9 bits. In the figure below, Local node LN1 is encoded as ENSI=0 using 3 bits and CNI=LN1's Index using 9 bits. LN1 is encoded in 12 bits in total.

```

  0 1 2 3 4 5 6 7 8
+---+---+---+---+---+---+
| 0 0 0 |           ENSI (3 bits) [9 bits CNI]
+---+---+---+---+---+---+
| LN1 Index Value | CNI (9 bits)
+---+---+---+---+---+---+

```

An Example of Local Node Encoding

The adjacent nodes can be encoded in two parts: Number of Nodes (NN) and compact node indexes (CNIs). The size of CNIs is the same as the local node. For example, three adjacent nodes RN1, RN2 and RN3 are encoded below in 30 bits (i.e., 3.75 bytes).

```

  0 1 2 3 4 5 6 7 8
+---+---+---+---+---+---+
| 0 1 1 |           NN (3 bits) [3 adjacent nodes]
+---+---+---+---+---+---+
| RN1's Index | CNI (9 bits) for RN1
+---+---+---+---+---+---+
| RN2's Index | CNI (9 bits) for RN2
+---+---+---+---+---+---+
| RN3's Index | CNI (9 bits) for RN3
+---+---+---+---+---+---+

```

An Example of Adjacent Nodes Encoding

The links between a local node and a number of its adjacent (or remote) nodes can be encoded as the local node followed by the adjacent nodes. For example, three links between local node LN1 and its three adjacent nodes RN1, RN2 and RN3 are encoded below in 42 bits (i.e., 5.25 bytes).

0	1	2	3	4	5	6	7	8	
+--+--+--+--+--+--+--+--+									
0 0 0	ENSI (3 bits) [9 bits CNI]						-		
+--+--+--+--+--+--+--+--+							}	Encoding for	
LN1 Index Value	CNI (9 bits) for LN1						-	Local Node LN1	
+--+--+--+--+--+--+--+--+							-		
0 1 1	NN (3 bits) [3 nodes]						-		
+--+--+--+--+--+--+--+--+								Encoding for	
RN1's Index	CNI (9 bits) for RN1							3 adjacent nodes	
+--+--+--+--+--+--+--+--+							}	RN1, RN2, RN3	
RN2's Index	CNI (9 bits) for RN2							of LN1	
+--+--+--+--+--+--+--+--+							-		
RN3's Index	CNI (9 bits) for RN3						-		
+--+--+--+--+--+--+--+--+									

An Example of Links Encoding

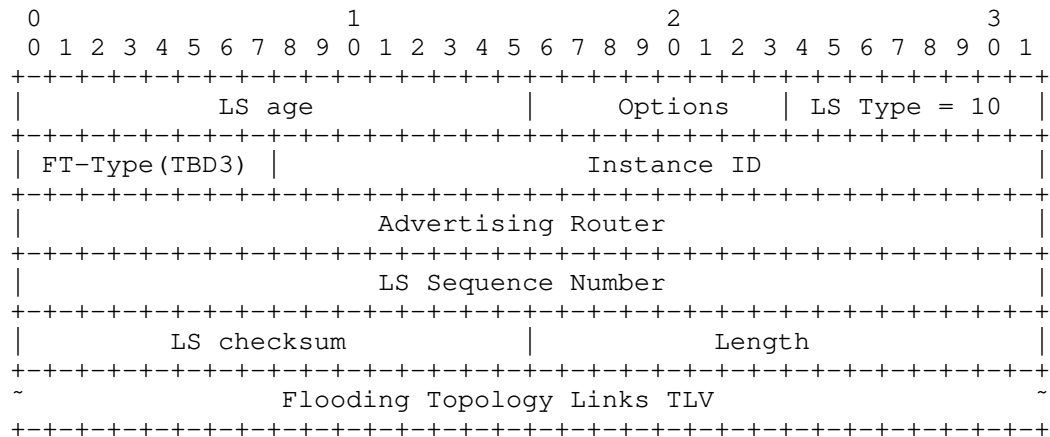
For a flooding topology computed by a leader of an area, it may be represented by all the links on the flooding topology. A Type-Length-Value (TLV) of the following format for the links encodings can be included in an LSA to represent the flooding topology (FT) and flood the FT to every node in the area.

[illegible]

Flooding Topology Links TLV

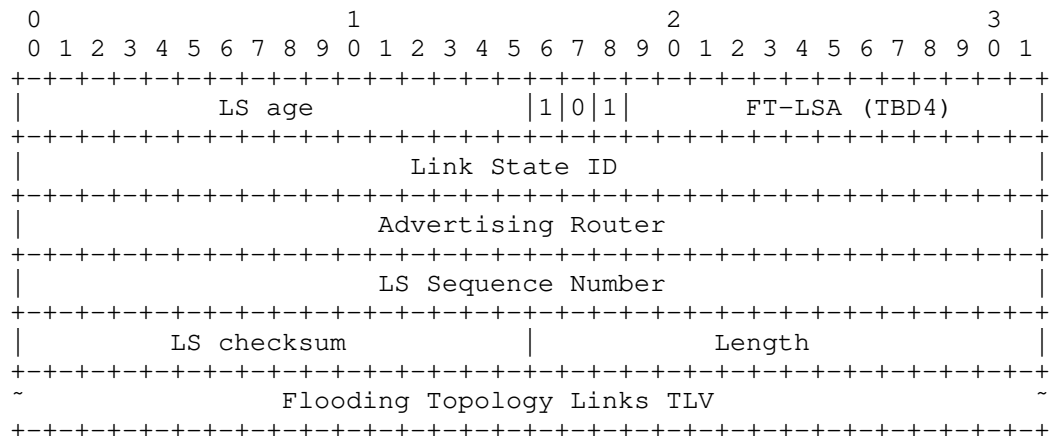
Note that a link between a local node LN and its adjacent node RN can be encoded once and as a bi-directional link. That is that if it is encoded in a Links Encoding from LN to RN, then the link from RN to LN is implied or assumed.

For OSPFv2, an Opaque LSA of a new opaque type (TBD3) containing a Flooding Topology Links TLV is used to flood the flooding topology from the leader of an area to all the other nodes in the area.



OSPFv2: Flooding Topology Opaque LSA

For OSPFv3, an area scope LSA of a new LSA function code (TBD4) containing a Flooding Topology Links TLV is used to flood the flooding topology from the leader of an area to all the other nodes in the area.



OSPFv3: Flooding Topology LSA

The U-bit is set to 1, and the scope is set to 01 for area-scoping.

6.2.1.2. Block Encoding

Block encoding uses a single structure to encode a block (or part) of topology, which can be a block of links in a flooding topology. It can also be all the links in the flooding topology. It starts with a local node LN and its adjacent (or remote) nodes RN_i ($i = 1, 2, \dots, n$), and can be considered as an extension to the links encoding.

The encoding of links between a local node and its adjacent nodes described in Section 6.2.1.1 is extended to include the links attached to the adjacent nodes.

The encoding for the adjacent nodes is extended to include Extending Flags (E Flags for short) between the NN (Number of Nodes) field and the CNIs (Compact Node Indexes) for the adjacent nodes. The length of the E Flags field is NN bits. The following is an example encoding of the adjacent nodes with E Flags of 3 bits, which is the value of the NN (the number of adjacent nodes).

```

  0 1 2 3 4 5 6 7 8
+---+---+---+---+---+---+
| 0 1 1 |          NN (3 bits)   [3 adjacent nodes]
+---+---+
| 1 0 1 |          E Flags [NN=3 bits]
+---+---+---+---+---+---+
|  RN1's Index  |  CNI (9 bits) for RN1
+---+---+---+---+---+---+
|  RN2's Index  |  CNI (9 bits) for RN2
+---+---+---+---+---+---+
|  RN3's Index  |  CNI (9 bits) for RN3
+---+---+---+---+---+---+

```

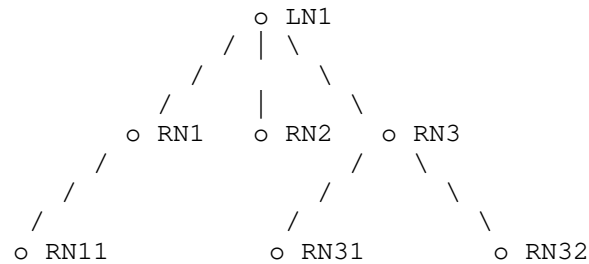
An Example of Adjacent Nodes with E Flags Encoding

There is a bit flag (called E flag) in the E Flags field for each adjacent node. The first bit (i.e., the most significant bit) in the E Flags field is for the first adjacent node (e.g., RN1), the second bit is for the second adjacent node (e.g., RN2), and so on. The E flag for an adjacent node RN_i set to one indicates that the links attached to the adjacent node RN_i are included below. The E flag for an adjacent node RN_i set to zero means that no links attached to the adjacent node RN_i are included below.

The links attached to the adjacent node RN_i are represented by the RN_i as a local node and the adjacent nodes of RN_i . The encoding for the adjacent nodes of RN_i is the same as that for the adjacent nodes

of a local node. It consists of an NN field of 3 bits, E Flags field of NN bits, and CNIs for the adjacent nodes of RN_i.

The following is an example of a block encoding for a block (or part) of flooding topology below.



An Example Block of Flooding Topology

It represents 6 links: 3 links between local node LN1 and its 3 adjacent nodes RN1, RN2 and RN3; 1 link between RN1 as a local node and its 1 adjacent node RN11; and 2 links between RN3 as a local node and its 2 adjacent nodes RN31 and RN32.

It starts with the encoding of the links between local node LN1 and 3 adjacent nodes RN1, RN2 and RN3 of the local node LN1. The encoding for the local node LN1 is the same as that for a local node described in Section 6.2.1.1. The encoding for 3 adjacent nodes RN1, RN2 and RN3 of local node LN1 comprises an NN field of 3 bits with value of 3, E Flags field of NN = 3 bits, and the indexes of adjacent nodes RN1, RN2 and RN3.

0 1 2 3 4 5 6 7 8			
+--+--+--+--+--+--+--+			
0 0 0	ENSI (3 bits) [9 bits CNI]	-	
+--+--+--+--+--+--+--+			
LN1 Index Value	CNI (9 bits)	-	Encoding for Local Node LN1
+--+--+--+--+--+--+--+			
0 1 1	NN(3 bits)[3 adjacent nodes]	-	
+--+--+			
1 0 1	E Flags [NN=3 bits]		Encoding for 3 adjacent nodes (RN1, RN2, RN3) of LN1
+--+--+--+--+--+--+--+			
RN1's Index	CNI (9 bits) for RN1		
+--+--+--+--+--+--+--+			
RN2's Index	CNI (9 bits) for RN2		
+--+--+--+--+--+--+--+			
RN3's Index	CNI (9 bits) for RN3	-	
+--+--+--+--+--+--+--+			
0 0 1	NN (3 bits)[1 adjacent node]	-	
+--+--+			
0	E Flags [NN=1 bit]		Encoding for 1 adjacent node (RN11) of RN1
+--+--+--+--+--+--+--+			
RN11's Index	CNI (9 bits) for RN11	-	
+--+--+--+--+--+--+--+			
0 1 0	NN(3 bits)[2 adjacent nodes]	-	
+--+--+			
0 0	E Flags [NN=2 bits]		Encoding for 2 adjacent nodes (RN31, RN32) of RN3 as a local node
+--+--+--+--+--+--+--+			
RN31's Index	CNI (9 bits) for RN31		
+--+--+--+--+--+--+--+			
RN32's Index	CNI (9 bits) for RN32		
+--+--+--+--+--+--+--+			

An Example of Block Encoding

The first E flag in the encoding for adjacent nodes RN1, RN2 and RN3 is set to one, which indicates that the links between the first adjacent node RN1 as a local node and its adjacent nodes are included below. In this example, 1 link between RN1 and its adjacent node RN11 is represented by the encoding for the adjacent node RN11 of RN1 as a local node. The encoding for 1 adjacent node RN11 consists of an NN field of 3 bits with value of 1, E Flags field of NN = 1 bits, and the index of adjacent node RN11. The size of the index of RN11 is the same as that of local node LN1 indicated by the ENSI in the encoding for local node LN1.

The second E flag in the encoding for adjacent nodes RN1, RN2 and RN3 is set to zero, which indicates that no links between the second

adjacent node RN2 as a local node and its adjacent nodes are included below.

The third E flag in the encoding for adjacent nodes RN1, RN2 and RN3 is set to one, which indicates that the links between the third adjacent node RN3 as a local node and its adjacent nodes are included below. In this example, 2 links between RN3 and its 2 adjacent nodes RN31 and RN32 are represented by the encoding for the adjacent nodes RN31 and RN32 of RN3 as a local node. The encoding for 2 adjacent nodes RN31 and RN32 consists of an NN field of 3 bits with value of 2, E Flags field of NN = 2 bits, and the indexes of adjacent nodes RN31 and RN32. The size of the index of RN31 and RN32 is the same as that of local node LN1 indicated by the ENSI in the encoding for local node LN1.

The block encoding may be used in the place of the links encoding in Section 6.2.1.1 for more efficiency. That is that it may be used in a Flooding Topology Links TLV. Alternatively, a new TLV, which is similar to the Flooding Topology Links TLV, may be defined to contain a number of block encodings.

6.2.2. Encodings for Backup Paths

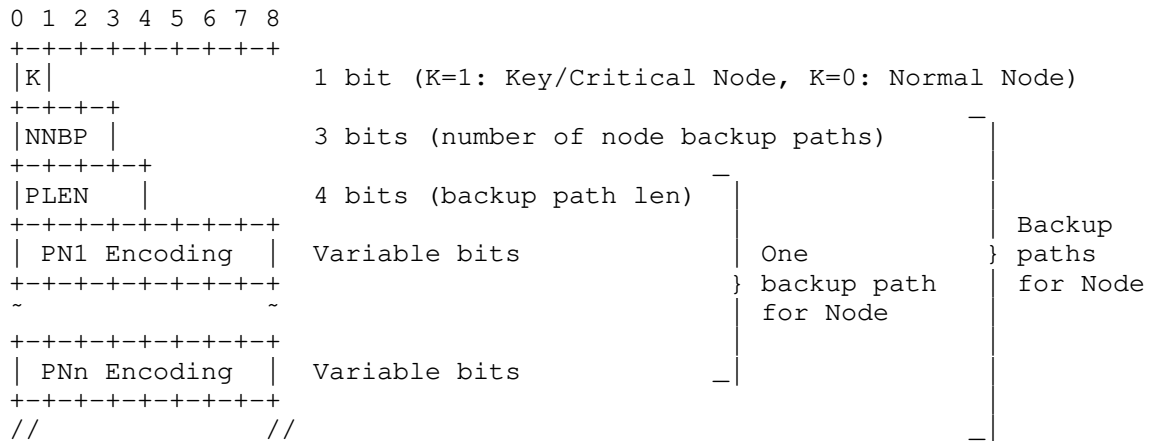
When the leader of an area computes a flooding topology, it may compute a backup path or multiple backup paths for a critical link on the flooding topology. When the critical link fails, a link state can be distributed to every node in the area through one backup path and other links on the flooding topology. In addition, it may compute a backup path or multiple backup paths for a node. When the node fails, a link state can be distributed to the other nodes in the area through the backup paths and the links on the flooding topology.

This section describes two encodings for backup paths: separated encoding and integrated one. In the former, backup paths are encoded in a new message, where the message for the flooding topology described in the previous section is required; In the latter, backup paths are integrated into the flooding topology links encoding, where one message contains the flooding topology and the backup paths.

6.2.2.1. Message for Backup Paths

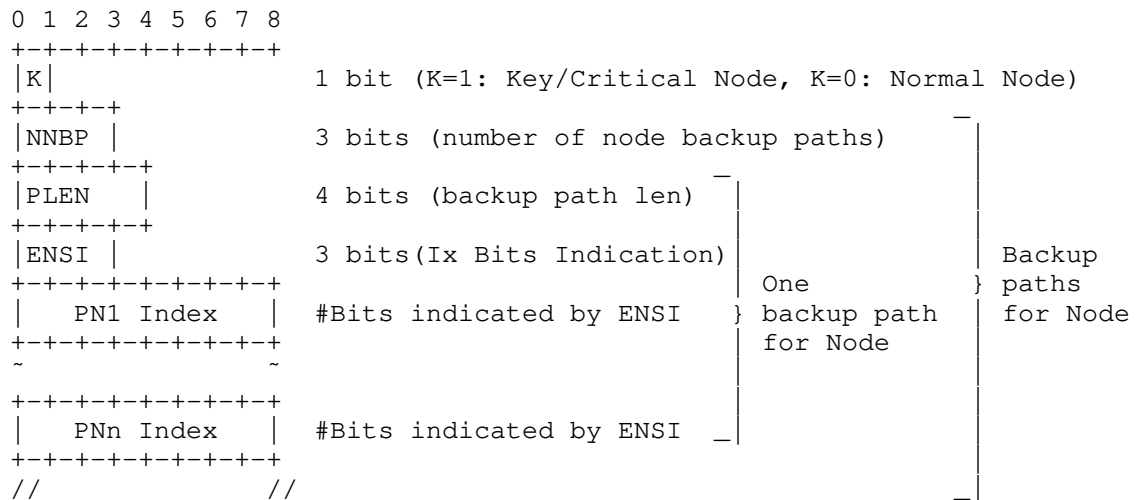
Backup paths for a node (such as Node1) may be represented by the node index encoding and node backup paths encoding. The former is similar to local node index encoding. The latter has the following format. It comprises a K flag (Key/Critical node flag) of 1 bit, a 3 bits NNBP field (number of node backup paths), and each of the backup paths encoding, which consists of the path length PLEN of 4 bits indicating the length of the path (i.e., the number of nodes), and

the encoding of the sequence of nodes along the path such as encodings for nodes PN1, ..., PNn. The encoding of every node may use the encoding of a local node, which comprises encoded node index size indication (ENSI) and compact node index (CNI).



An Example of Node Backup Paths Encoding

Another encoding of the sequence of nodes along the path uses one encoded node index size indication (ENSI) for all the nodes in the path. Thus we have the following Node Backup Paths Encoding.



Another Example of Node Backup Paths Encoding

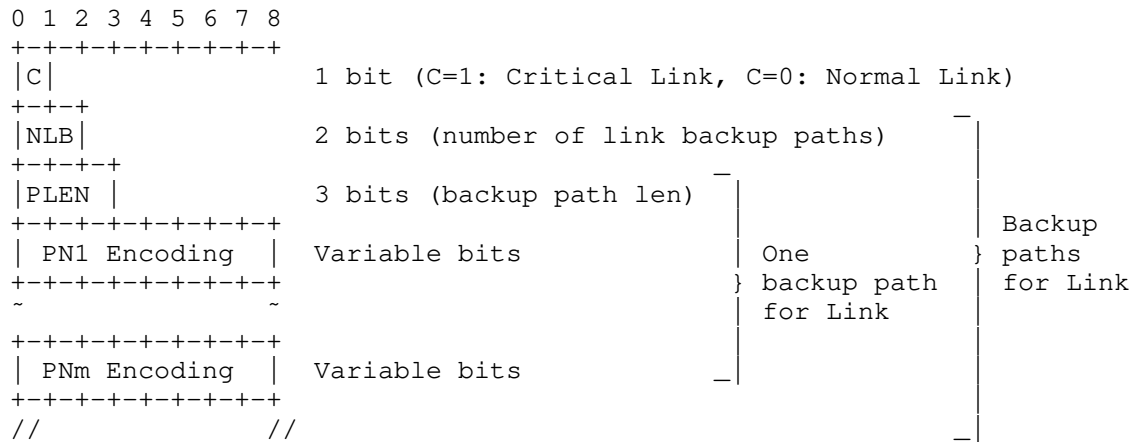
A new TLV called Node Backup Paths TLV is defined below. It may include multiple nodes and their backup paths. Each node is represented by its index encoding, which is followed by its node backup paths encoding.

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|      NBP-TLV-Type (TBD5)      |      TLV-Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|Node1 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:      Node1 backup paths encoding      :
+-----+-----+-----+-----+-----+-----+-----+-----+
|Node2 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:      Node2 backup paths encoding      :
+-----+-----+-----+-----+-----+-----+-----+-----+
//                                         //
                                     Node Backup Paths TLV

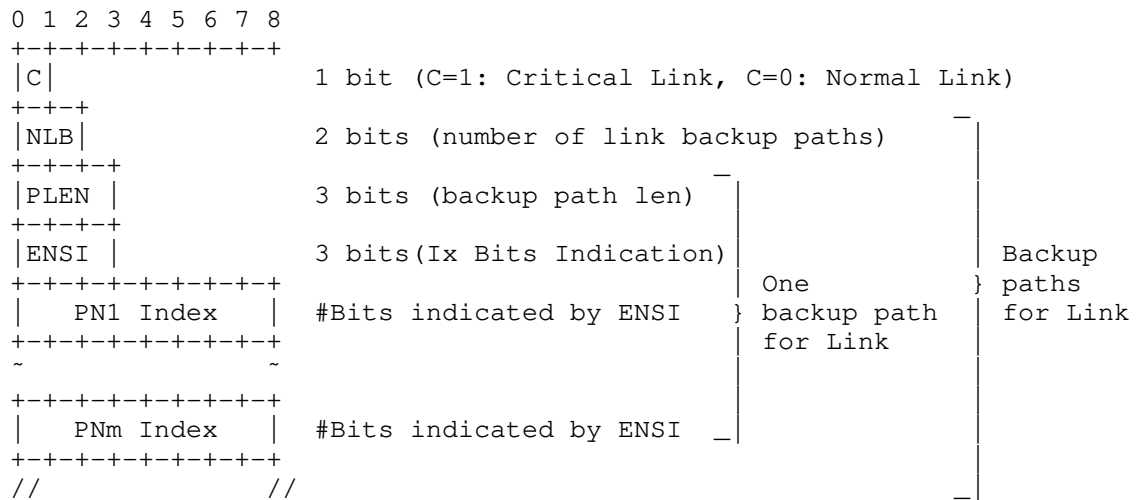
```

The encoding for backup paths for a link (such as Link1) on the flooding topology consists of the link encoding such as Link1 Index Encoding and the link backup paths encoding. The former is similar to local node encoding. It contains encoded link index size indication (ELSI) and compact link index (CLI). The latter has the following format. It comprises a C flag (Critical link flag) of 1 bit, a 2 bits NLB field (number of link backup paths), and each of the backup paths encoding, which consists of the path length PLEN of 3 bits indicating the length of the path (i.e., the number of nodes), and the encoding of the sequence of nodes along the path such as encodings for nodes PN1, ..., PNm. Note that two ends of a link (i.e., the local node and the adjacent/remote node of the link) are not needed in the path. The encoding of every node may use the encoding of a local node, which comprises encoded node index size indication (ENSI) and compact node index (CNI).



An Example of Link Backup Paths Encoding

Another encoding of the sequence of nodes along the path uses one encoded node index size indication (ENSI) for all the nodes in the path. Thus we have the following Link Backup Paths Encoding.



Another Example of Link Backup Paths Encoding

A new TLV called Link Backup Paths TLV is defined below. It may include multiple links and their backup paths. Each link is represented by its index encoding, which is followed by its link backup paths encoding.


```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|      LBP-TLV-Type (TBD6)      |      TLV-Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|Link1 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Link1 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
|Link2 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Link2 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
//                                                                    //
                                Link Backup Paths TLV

```

For OSPFv2, an Opaque LSA of a new opaque type (TBD7), containing node backup paths TLVs and link backup paths TLVs, is used to flood the backup paths from the leader of an area to all the other nodes in the area.

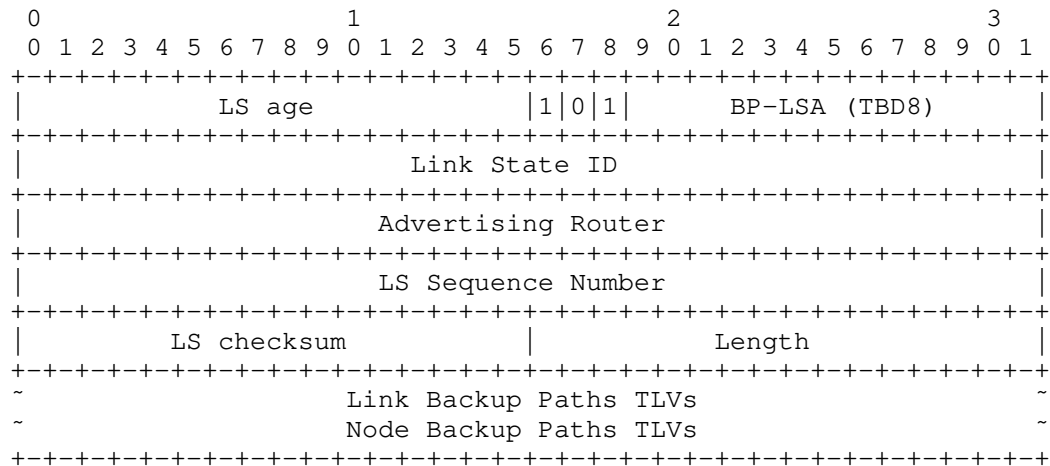
```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|      LS age      |      Options      | LS Type = 10 |
+-----+-----+-----+-----+-----+-----+-----+-----+
| BP-Type(TBD7) |      Instance ID      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|      Advertising Router      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|      LS Sequence Number      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|      LS checksum      |      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               Link Backup Paths TLVs                               ~
~                               Node Backup Paths TLVs                               ~
+-----+-----+-----+-----+-----+-----+-----+-----+

```

OSPFv2: Backup Paths Opaque LSA

For OSPFv3, an area scope LSA of a new LSA function code (TBD8), containing node backup paths TLVs and link backup paths TLVs, is used to flood the backup paths from the leader of an area to all the other nodes in the area.

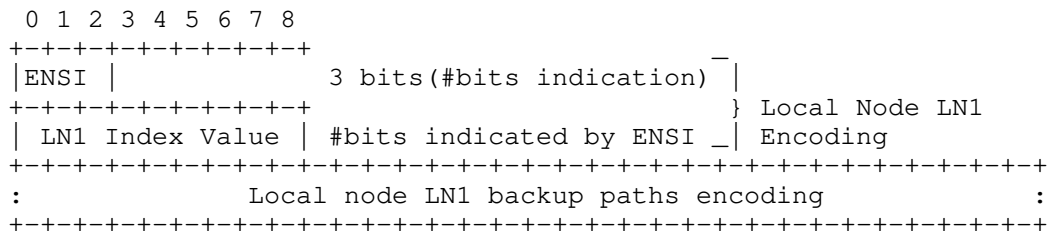


OSPFv3: Backup Paths LSA

The U-bit is set to 1, and the scope is set to 01 for area-scoping.

6.2.2.2. Backup Paths in Links TLV

A local node and its backup paths can be encoded in the following format. It is the local node (such as local node LN1) encoding followed by the local node backup paths encoding, which is the same as the node backup paths encoding described in Section 6.2.2.1.



Local Node with Backup Paths Encoding

A adjacent node and its backup paths can be encoded in the following format. It is the adjacent node (such as adjacent node RN10) index value followed by the adjacent node backup paths encoding, which is the same as the node backup paths encoding described in Section 6.2.2.1.

```

+-----+
|RN10 Index Value |  (#bits indicated by ENSI)
+-----+
:                adjacent node RN10 backup paths encoding                :
+-----+

```

Adjacent Node with Backup Paths Encoding

The links between a local node and a number of its adjacent nodes, the backup paths for each of the nodes, and the backup paths for each of the links can be encoded in the following format. It is called Links from Node with Backup Paths Encoding.

```

      0              1              2              3
      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
+-----+
:                Local Node with backup paths encoding                :
+-----+
| NN |  Number of adjacent Nodes (i.e., Number of links)
+-----+
:                Adjacent Node 1 with backup paths encoding            :
+-----+
:                Link1 backup paths Encoding                            :
+-----+
:                Adjacent Node 2 with backup paths encoding            :
+-----+
:                Link2 backup paths Encoding                            :
+-----+
|

```

Links from Node with Backup Paths Encoding

A new TLV called Links with Backup Paths TLV is defined below. It includes a number of Links from Node with Backup Paths Encodings described above. This TLV contains both the flooding topology and the backup paths for the links and nodes on the flooding topology.

```

      0               1               2               3
    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
+-----+-----+-----+-----+-----+-----+-----+-----+
|  LNSBP-TLV-Type (TBD9)  |          TLV-Length          |
+-----+-----+-----+-----+-----+-----+-----+-----+
:   Links from Node 1 with backup paths encoding   :
+-----+-----+-----+-----+-----+-----+-----+-----+
:   Links from Node 2 with backup paths encoding   :
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                     :
:                                                     :
:   Links with Backup Paths TLV                     :

```

For OSPFv2, an Opaque LSA of a new opaque type (TBDA), called Flooding Topology with Backup Paths (FTBP) Opaque LSA, containing a Links with Backup Paths TLV, is used to flood the flooding topology with backup paths from the leader of an area to all the other nodes in the area.

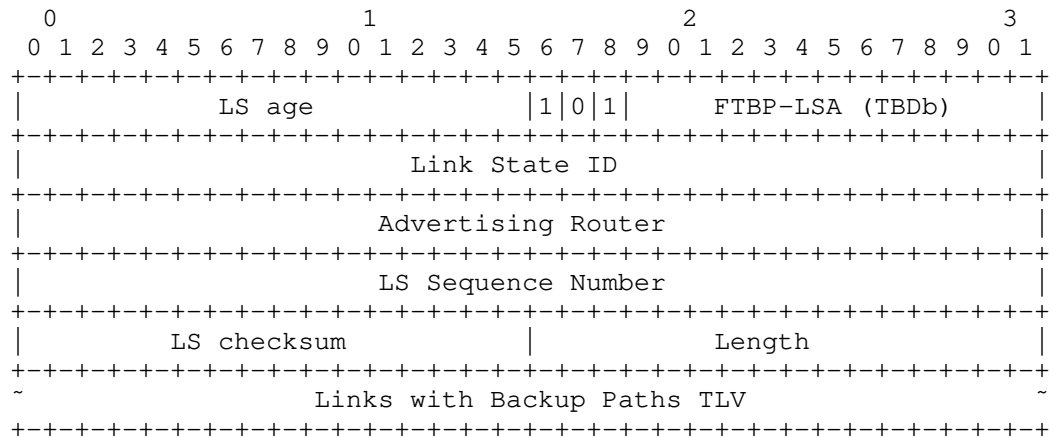
```

      0               1               2               3
    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
+-----+-----+-----+-----+-----+-----+-----+-----+
|          LS age          |      Options      | LS Type = 10 |
+-----+-----+-----+-----+-----+-----+-----+-----+
| FTBP-Type(TBDA) |          Instance ID          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          Advertising Router          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          LS Sequence Number          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          LS checksum          |          Length          |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               ~
~   Links with Backup Paths TLV   ~
~                               ~

```

OSPFv2: Flooding Topology with Backup Paths (FTBP) Opaque LSA

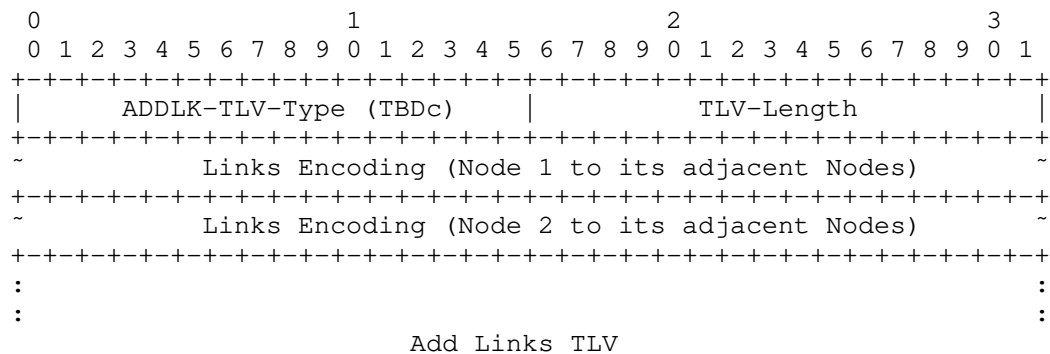
For OSPFv3, an area scope LSA of a new LSA function code (TBDb), containing a Links with Backup Paths TLV, is used to flood the flooding topology with backup paths from the leader of an area to all the other nodes in the area.



OSPFv3: Flooding Topology with Backup Paths (FTBP) LSA

6.2.3. Message for Incremental Changes

For adding some links to the flooding topology, we define a new TLV called Add Links TLVs of the following format. When some new links are added to the flooding topology, the leader may not flood the whole flooding topology with the new links to all the other nodes. It may just flood these new links. After receiving these new links, each of the other nodes adds these new links into the existing flooding topology. When the leader floods the whole flooding topology with the new links to all the other nodes, it removes the LSA for the new links. When removing the LSA for these new links, each of the other nodes does not update the flooding topology (i.e., does not remove these links from the flooding topology).



For deleting some links from the flooding topology, we define a new TLV called Delete Links TLVs of the following format. When some old links are removed from the flooding topology, the leader may not flood the whole flooding topology without the old links to all the other nodes. It may just flood these old links. After receiving these old links, each of the other nodes deletes these old links from the existing flooding topology. When the leader floods the whole flooding topology without the old links to all the other nodes, it removes the LSA for the old links. When removing the LSA for these old links, each of the other nodes does not update the flooding topology (i.e., does not add these links into the flooding topology).

```

      0               1               2               3
    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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|      DELLK-TLV-Type (TBDd)      |      TLV-Length      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
~                               Links Encoding (Node 1 to its adjacent Nodes)                               ~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
~                               Links Encoding (Node 2 to its adjacent Nodes)                               ~
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:                                                                           :
:                                                                           :
                                Delete Links TLV

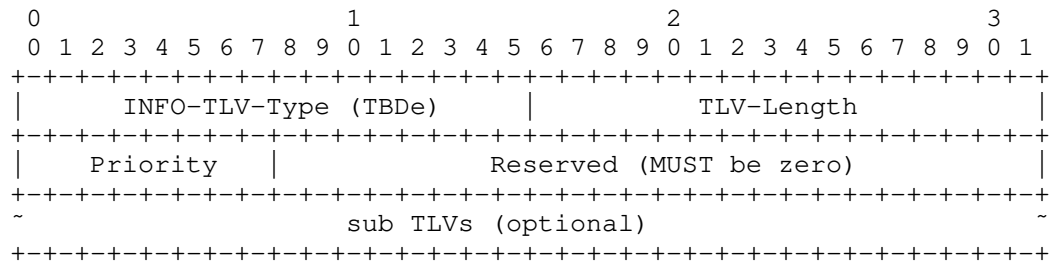
```

The Add Links TLVs and Delete Links TLVs should be in a separate LSA instance. The LSA can be a Flooding Topology LSA defined above. Alternatively, we may define a new LSA for these TLVs.

6.2.4. Leaders Selection

The leader or Designated Router (DR) selection for a broadcast link is about selecting two leaders: a DR and Backup DR. This is generalized to select two or more leaders for an area: the primary/first leader (or leader for short), the secondary leader, the third leader and so on.

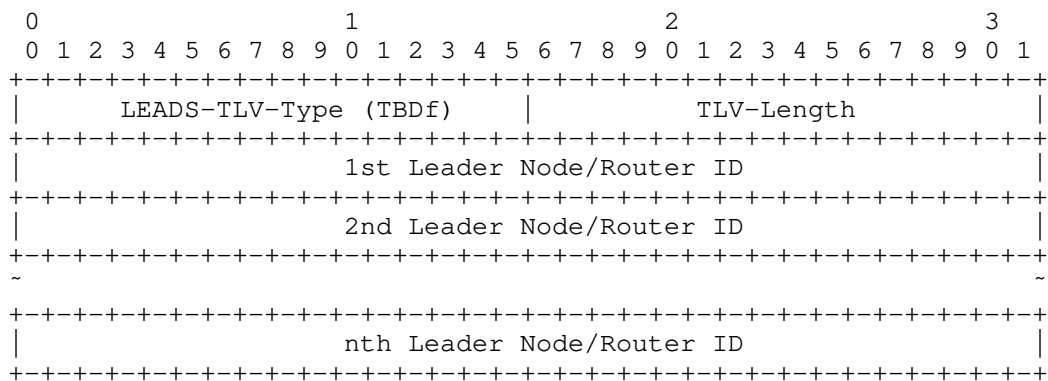
A new TLV is defined to include the information on flooding reduction of a node, which is called Flooding Reduction Information TLV or Information TLV for short. This TLV is generated by every node that supports flooding reduction in general. Every node originates a RI LSA with a Flooding Reduction Information TLV containing its priority to become a leader. The format of the TLV is as follows.



Flooding Reduction Information TLV

A Priority field of eight bits is defined in the TLV to indicate the priority of the node originating the TLV to become the leader node in central mode.

A sub-TLV called leaders sub-TLV is defined. It has the following format.



Leaders sub-TLV

When a node selects itself as a leader, it originates a RI LSA containing the leader in a leaders sub-TLV.

After the first leader node is down, the other leaders will be promoted. The secondary leader becomes the first leader, the third leader becomes the secondary leader, and so on. When a node selects itself as the n-th leader, it originates a RI LSA with a Leaders sub-TLV containing n leaders.

7. Extensions to IS-IS

The extensions to IS-IS is similar to OSPF.

7.1. Extensions for Operations

A new TLV for operations is defined in IS-IS LSP. It has the following format and contains the same contents as the Flooding Reduction Instruction TLV defined in OSPF RI LSA.

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
| INST-Type (TBDi1) | Length |
+-----+-----+-----+-----+-----+-----+-----+-----+
| OP | MOD | Algorithm | Reserved (MUST be zero) | NL |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               sub TLVs (optional)                               ~
+-----+-----+-----+-----+-----+-----+-----+-----+

```

IS-IS Flooding Reduction Instruction TLV

7.2. Extensions for Centralized Mode

7.2.1. TLV for Flooding Topology

A new TLV for the encodings of the links in the flooding topology is defined. It has the following format and contains the same contents as the Flooding Topology Links TLV defined in OSPF Flooding Topology Opaque LSA.

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
| FTL-Type (TBDi2) | Length |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               Links Encoding (Node 1 to its adjacent Nodes)                               ~
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               Links Encoding (Node 2 to its adjacent Nodes)                               ~
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                           :
:                                                                           :

```

IS-IS Flooding Topology Links TLV

7.2.2. Encodings for Backup Paths

7.2.2.1. TLVs for Backup Paths

For flooding backup paths separately, we define two TLVs: IS-IS Node Backup Paths TLV and IS-IS Link Backup Path TLV. The former has the following format and contains the same contents as Node Backup Paths TLV in OSPF.

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|NBP-Type(TBDi3)|   Length   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|Node1 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Node1 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
|Node2 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Node2 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
//                                                                    //
                                IS-IS Node Backup Paths TLV

```

The latter has the following format and contains the same contents as Link Backup Paths TLV in OSPF.

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|LBP-Type(TBDi4)|   Length   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|Link1 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Link1 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
|Link2 Index Enc| Variable bits
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Link2 backup paths encoding                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
//                                                                    //
                                IS-IS Link Backup Paths TLV

```

7.2.2.2. Backup Paths in Links TLV

A new TLV is defined to integrate the backup paths with the links on the flooding topology. It has the following format and contains the same contents as the Links with Backup Paths TLV in OSPF.

```

      0                               1                               2                               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|LSB-Type(TBDi5)|      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
:               Links from Node 1 with backup paths encoding           :
+-----+-----+-----+-----+-----+-----+-----+-----+
:               Links from Node 2 with backup paths encoding           :
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                           :
:                                                                           :
:               IS-IS Links with Backup Paths TLV

```

7.2.3. TLVs for Incremental Changes

Similar to Add Links TLV in OSPF, a new TLV called IS-IS Add Links TLV is defined. It has the following format and contains the same contents as Add Links TLV in OSPF.

```

      0                               1                               2                               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|ADDL-Type(TBDi6)|      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
~               Links Encoding (Node 1 to its adjacent Nodes)           ~
+-----+-----+-----+-----+-----+-----+-----+-----+
~               Links Encoding (Node 2 to its adjacent Nodes)           ~
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                           :
:                                                                           :
:               IS-IS Add Links TLV

```

Similar to Delete Links TLV in OSPF, a new TLV called IS-IS Delete Links TLV is defined. It has the following format and contains the same contents as Delete Links TLV in OSPF.

```

      0               1               2               3
    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
+-----+-----+-----+-----+-----+-----+-----+-----+
|DELL-Type(TBDi7|      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               Links Encoding (Node 1 to its adjacent Nodes) ~
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               Links Encoding (Node 2 to its adjacent Nodes) ~
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                    :
:                                                                    :
                                IS-IS Delete Links TLV

```

7.2.4. Leaders Selection

Similar to Flooding Reduction Information TLV in OSPF, a new TLV called IS-IS Flooding Reduction Information TLV is defined. It has the following format and contains the same contents as Flooding Reduction Information TLV in OSPF.

```

      0               1               2               3
    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
+-----+-----+-----+-----+-----+-----+-----+-----+
|INF-Type(TBDi8)|      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|      Priority      |      Reserved (MUST be zero)      |
+-----+-----+-----+-----+-----+-----+-----+-----+
~                               sub TLVs (optional) ~
+-----+-----+-----+-----+-----+-----+-----+-----+
                                IS-IS Flooding Reduction Information TLV

```

8. Flooding Behavior

This section describes the revised flooding behavior for a node having at least one link on the flooding topology. The revised flooding procedure MUST flood an LS to every node in the network in any case, as the standard flooding procedure does.

8.1. Nodes Perform Flooding Reduction without Failure

8.1.1. Receiving an LS

When a node receives a newer LS that is not originated by itself from one of its interfaces, it floods the LS only to all the other interfaces that are on the flooding topology.

When the LS is received from an interface on the flooding topology, it is flooded only to all the other interfaces that are on the flooding topology. When the LS is received on an interface that is not on the flooding topology, it is also flooded only to all the other interfaces that are on the flooding topology.

In any case, the LS must not be transmitted back to the receiving interface.

Note before forwarding a received LS, the node would do the normal processing as usual.

8.1.2. Originating an LS

When a node originates an LS, it floods the LS to its interfaces on the flooding topology if the LS is a refresh LS (i.e., there is no significant change in the LS comparing to the previous LS); otherwise (i.e., there are significant changes in the LS), it floods the LS to all its interfaces. Choosing flooding the LS with significant changes to all the interfaces instead of limiting to the interfaces on the flooding topology would speed up the distribution of the significant link state changes.

8.1.3. Establishing Adjacencies

Adjacencies being established can be classified into two categories: adjacencies to new nodes and adjacencies to existing nodes.

8.1.3.1. Adjacency to New Node

An adjacency to a new node is an adjacency between a node (say node A) on the flooding topology and the new node (say node Y) which is not on the flooding topology. There is not any adjacency between node Y and a node in the network area.

When new node Y is up and connected to node A, node A assumes that node Y and the link between node Y and node A are on the flooding topology until a new flooding topology is computed and built. Node A may determine whether node Y is a new node through checking if node Y is reachable or on the flooding topology.

The procedure for establishing the adjacency between node A and node Y is the existing normal procedure unchanged. After the status of the adjacency reaches to Exchange or Full, node A sends node Y every new or updated LS that node A receives or originates.

8.1.3.2. Adjacency to Existing Node

An adjacency to an existing node is an adjacency between a node (say node A) on the flooding topology and the existing node (say node X) which exists on the flooding topology. There are some adjacencies between node X and some nodes in the network area.

When existing node X is connected to node A after a link between node X and node A is up, node A assumes that the link connecting node A and node X is not on the flooding topology until a new flooding topology is computed and built. Node A may determine whether node X is an existing node through checking if node X is reachable or on the flooding topology.

The procedure for establishing the adjacency between node A and node X is the existing normal procedure unchanged. Node A does not send node X any new or updated LS that node A receives or originates even after the status of the adjacency reaches to Exchange or Full.

8.2. An Exception Case

During an LS flooding, one or multiple link and node failures may happen. Some failures do not split the flooding topology, thus do not affect the flooding behavior. For example, multiple failures of the links not on the flooding topology do not split the flooding topology and do not affect the flooding behavior. The sections below focus on the failures that may split the flooding topology.

8.2.1. A Critical Failure

For a link failure, if the link is a critical link on the flooding topology, then the LS is flooded through a backup path for the link and the remaining flooding topology until a new flooding topology is computed and built; otherwise, the flooding behavior in Section 8.1 follows.

Similarly, for a node failure, if the node is a critical node on the flooding topology, then the LS is flooded through backup paths for the node and the remaining flooding topology until a new flooding topology is computed and built; otherwise, the flooding behavior in Section 8.1 follows.

8.2.2. Multiple Failures

For multiple link failures, if the number of the failed links on the flooding topology is greater than or equal to two, then the LS is flooded through a backup path for each of the failed links on the flooding topology and the remaining flooding topology until a new

flooding topology is computed and built; otherwise, the flooding behavior in Section 8.1 follows.

If all the backup paths for some of the failed links are broken by some failures, the LS is flooded to all interfaces (except where it is received from) until a new flooding topology is computed and built.

For multiple node failures, the LS is flooded through the backup paths for each of the failed nodes and the remaining flooding topology until a new flooding topology is computed and built; otherwise, the flooding behavior in Section 8.1 follows.

If the backup paths for some of the failed nodes are broken by some failures, the LS is flooded to all interfaces (except where it is received from) until a new flooding topology is computed and built.

Note that if it can be quickly determined that the flooding topology is not split by the failures, the flooding behavior in Section 8.1 may follow.

9. Security Considerations

This document does not introduce any security issue.

10. IANA Considerations

10.1. OSPFv2

Under Registry Name: OSPF Router Information (RI) TLVs [RFC7770], IANA is requested to assign two new TLV values for OSPF flooding reduction as follows:

TLV Value	TLV Name	reference
11	Instruction TLV	This document
12	Information TLV	This document

Under the registry name "Opaque Link-State Advertisements (LSA) Option Types" [RFC5250], IANA is requested to assign new Opaque Type registry values for FT LSA, BP LSA, FTBP LSA as follows:

Registry Value	Opaque Type	reference
10	FT LSA	This document
11	BP LSA	This document
12	FTBP LSA	This document

IANA is requested to create and maintain new registries:

o OSPFv2 FT LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv2 FT LSA TLV Name	Definition
-----	-----	-----
0	Reserved	
1	FT Links TLV	see Section 6.2.1
2-32767	Unassigned	
32768-65535	Reserved	

o OSPFv2 BP LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv2 TBPLSA TLV Name	Definition
-----	-----	-----
0	Reserved	
1	Node Backup Paths TLV	see Section 6.2.2
2	Link Backup Paths TLV	see Section 6.2.2
3-32767	Unassigned	
32768-65535	Reserved	

o OSPFv2 FTBP LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv2 FTBP LSA TLV Name	Definition
-----	-----	-----
0	Reserved	
1	Links with Backup Paths TLV	see Section 6.2.2
2-32767	Unassigned	
32768-65535	Reserved	

10.2. OSPFv3

Under the registry name "OSPFv3 LSA Function Codes", IANA is requested to assign new registry values for FT LSA, BP LSA, FTBP LSA as follows:

Value	LSA Function Code Name	reference
16	FT LSA	This document
17	BP LSA	This document
18	FTBP LSA	This document

IANA is requested to create and maintain new registries:

- o OSPFv3 FT LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv3 FT LSA TLV Name	Definition
0	Reserved	
1	FT Links TLV	see Section 6.2.1
2-32767	Unassigned	
32768-65535	Reserved	

- o OSPFv3 BP LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv3 TBPLSA TLV Name	Definition
0	Reserved	
1	Node Backup Paths TLV	see Section 6.2.2
2	Link Backup Paths TLV	see Section 6.2.2
3-32767	Unassigned	
32768-65535	Reserved	

- o OSPFv3 FTBP LSA TLVs

Initial values for the registry are given below. The future assignments are to be made through IETF Review [RFC5226].

Value	OSPFv3 FTBP LSA TLV Name	Definition
-----	-----	-----
0	Reserved	
1	Links with Backup Paths TLV	see Section 6.2.2
2-32767	Unassigned	
32768-65535	Reserved	

10.3. IS-IS

Under Registry Name: IS-IS TLV Codepoints, IANA is requested to assign new TLV values for IS-IS flooding reduction as follows:

Value	TLV Name	Definition
-----	-----	-----
151	FT Links TLV	see Section 7.2.1
152	Node Backup Paths TLV	see Section 7.2.2
153	Link Backup Paths TLV	see Section 7.2.2
154	Links with Backup Paths TLV	see Section 7.2.2
155	Add Links TLV	see Section 7.2.3
156	Delete Links TLV	see Section 7.2.3
157	Instruction TLV	see Section 7.1
158	Information TLV	see Section 7.2.4

11. Acknowledgements

The authors would like to thank Acee Lindem, Zhibo Hu, Robin Li, Stephane Litkowski and Alvaro Retana for their valuable suggestions and comments on this draft.

12. References

12.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC5250] Berger, L., Bryskin, I., Zinin, A., and R. Coltun, "The OSPF Opaque LSA Option", RFC 5250, DOI 10.17487/RFC5250, July 2008, <<https://www.rfc-editor.org/info/rfc5250>>.

- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC7770] Lindem, A., Ed., Shen, N., Vasseur, JP., Aggarwal, R., and S. Shaffer, "Extensions to OSPF for Advertising Optional Router Capabilities", RFC 7770, DOI 10.17487/RFC7770, February 2016, <<https://www.rfc-editor.org/info/rfc7770>>.

12.2. Informative References

- [I-D.li-dynamic-flooding]
Li, T. and P. Psenak, "Dynamic Flooding on Dense Graphs", draft-li-dynamic-flooding-05 (work in progress), June 2018.
- [I-D.shen-isis-spine-leaf-ext]
Shen, N., Ginsberg, L., and S. Thyamagundalu, "IS-IS Routing for Spine-Leaf Topology", draft-shen-isis-spine-leaf-ext-06 (work in progress), June 2018.
- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", RFC 5226, DOI 10.17487/RFC5226, May 2008, <<https://www.rfc-editor.org/info/rfc5226>>.

Appendix A. Algorithms to Build Flooding Topology

There are many algorithms to build a flooding topology. A simple and efficient one is briefed below.

- o Select a node R according to a rule such as the node with the biggest/smallest node ID;
- o Build a tree using R as root of the tree (details below); and then
- o Connect k ($k \geq 0$) leaves to the tree to have a flooding topology (details follow).

A.1. Algorithms to Build Tree without Considering Others

An algorithm for building a tree from node R as root starts with a candidate queue Cq containing R and an empty flooding topology Ft:

1. Remove the first node A from Cq and add A into Ft
2. If Cq is empty, then return with Ft

3. Suppose that node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in Ft and X_1, X_2, \dots, X_n are in a special order. For example, X_1, X_2, \dots, X_n are ordered by the cost of the link between A and X_i . The cost of the link between A and X_i is less than the cost of the link between A and X_j ($j = i + 1$). If two costs are the same, X_i 's ID is less than X_j 's ID. In another example, X_1, X_2, \dots, X_n are ordered by their IDs. If they are not ordered, then make them in the order.
4. Add X_i ($i = 1, 2, \dots, n$) into the end of Cq, goto step 1.

Another algorithm for building a tree from node R as root starts with a candidate queue Cq containing R and an empty flooding topology Ft:

1. Remove the first node A from Cq and add A into Ft
2. If Cq is empty, then return with Ft
3. Suppose that node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in Ft and X_1, X_2, \dots, X_n are in a special order. For example, X_1, X_2, \dots, X_n are ordered by the cost of the link between A and X_i . The cost of the link between A and X_i is less than the cost of the link between A and X_j ($j = i + 1$). If two costs are the same, X_i 's ID is less than X_j 's ID. In another example, X_1, X_2, \dots, X_n are ordered by their IDs. If they are not ordered, then make them in the order.
4. Add X_i ($i = 1, 2, \dots, n$) into the front of Cq and goto step 1.

A third algorithm for building a tree from node R as root starts with a candidate list Cq containing R associated with cost 0 and an empty flooding topology Ft:

1. Remove the first node A from Cq and add A into Ft
2. If all the nodes are on Ft, then return with Ft
3. Suppose that node A is associated with a cost C_a which is the cost from root R to node A, node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in Ft and the cost of the link between A and X_i is L_{Ci} ($i=1, 2, \dots, n$). Compute $C_i = C_a + L_{Ci}$, check if X_i is in Cq and if C_{xi} (cost from R to X_i) $< C_i$. If X_i is not in Cq, then add X_i with cost C_i into Cq; If X_i is in Cq, then If $C_{xi} > C_i$ then replace X_i with cost C_{xi} by X_i with C_i in Cq; If $C_{xi} == C_i$ then add X_i with cost C_i into Cq.
4. Make sure Cq is in a special order. Suppose that A_i ($i=1, 2, \dots, m$) are the nodes in Cq, C_{ai} is the cost associated with A_i ,

and ID_i is the ID of A_i . One order is that for any $k = 1, 2, \dots, m-1$, $C_{ak} < C_{aj}$ ($j = k+1$) or $C_{ak} = C_{aj}$ and $ID_k < ID_j$. Goto step 1.

A.2. Algorithms to Build Tree Considering Others

An algorithm for building a tree from node R as root with consideration of others's support for flooding reduction starts with a candidate queue Cq containing R associated with previous hop PH=0 and an empty flooding topology Ft:

1. Remove the first node A that supports flooding reduction from the candidate queue Cq if there is such a node A; otherwise (i.e., if there is not such node A in Cq), then remove the first node A from Cq. Add A into the flooding topology Ft.
2. If Cq is empty or all nodes are on Ft, then return with Ft
3. Suppose that node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in the flooding topology Ft and X_1, X_2, \dots, X_n are in a special order considering whether some of them that support flooding reduction (. For example, X_1, X_2, \dots, X_n are ordered by the cost of the link between A and X_i . The cost of the link between A and X_i is less than that of the link between A and X_j ($j = i + 1$). If two costs are the same, X_i 's ID is less than X_j 's ID. The cost of a link is redefined such that 1) the cost of a link between A and X_i both support flooding reduction is much less than the cost of any link between A and X_k where X_k with $F=0$; 2) the real metric of a link between A and X_i and the real metric of a link between A and X_k are used as their costs for determining the order of X_i and X_k if they all (i.e., A, X_i and X_k) support flooding reduction or none of X_i and X_k support flooding reduction.
4. Add X_i ($i = 1, 2, \dots, n$) associated with previous hop PH=A into the end of the candidate queue Cq, and goto step 1.

Another algorithm for building a tree from node R as root with consideration of others' support for flooding reduction starts with a candidate queue Cq containing R associated with previous hop PH=0 and an empty flooding topology Ft:

1. Remove the first node A that supports flooding reduction from the candidate queue Cq if there is such a node A; otherwise (i.e., if there is not such node A in Cq), then remove the first node A from Cq. Add A into the flooding topology Ft.
2. If Cq is empty or all nodes are on Ft, then return with Ft.

3. Suppose that node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in the flooding topology F_t and X_1, X_2, \dots, X_n are in a special order considering whether some of them support flooding reduction. For example, X_1, X_2, \dots, X_n are ordered by the cost of the link between A and X_i . The cost of the link between A and X_i is less than the cost of the link between A and X_j ($j = i + 1$). If two costs are the same, X_i 's ID is less than X_j 's ID. The cost of a link is redefined such that 1) the cost of a link between A and X_i both support flooding reduction is much less than the cost of any link between A and X_k where X_k does not support flooding reduction; 2) the real metric of a link between A and X_i and the real metric of a link between A and X_k are used as their costs for determining the order of X_i and X_k if they all (i.e., A, X_i and X_k) support flooding reduction or none of X_i and X_k supports flooding reduction.
4. Add X_i ($i = 1, 2, \dots, n$) associated with previous hop $PH=A$ into the front of the candidate queue C_q , and goto step 1.

A third algorithm for building a tree from node R as root with consideration of others' support for flooding reduction (using flag $F = 1$ for support, and $F = 0$ for not support in the following) starts with a candidate list C_q containing R associated with low order cost $L_c=0$, high order cost $H_c=0$ and previous hop ID $PH=0$, and an empty flooding topology F_t :

1. Remove the first node A from C_q and add A into F_t .
2. If all the nodes are on F_t , then return with F_t
3. Suppose that node A is associated with a cost C_a which is the cost from root R to node A, node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in F_t and the cost of the link between A and X_i is L_{Ci} ($i=1, 2, \dots, n$). Compute $C_i = C_a + L_{Ci}$, check if X_i is in C_q and if C_{xi} (cost from R to X_i) $< C_i$. If X_i is not in C_q , then add X_i with cost C_i into C_q ; If X_i is in C_q , then If $C_{xi} > C_i$ then replace X_i with cost C_{xi} by X_i with C_i in C_q ; If $C_{xi} == C_i$ then add X_i with cost C_i into C_q .
4. Suppose that node A is associated with a low order cost L_{Ca} which is the low order cost from root R to node A and a high order cost H_{Ca} which is the high order cost from R to A, node X_i ($i = 1, 2, \dots, n$) is connected to node A and not in the flooding topology F_t and the real cost of the link between A and X_i is C_i ($i=1, 2, \dots, n$). Compute L_{Cxi} and H_{Cxi} : $L_{Cxi} = L_{Ca} + C_i$ if both A and X_i have flag F set to one, otherwise $L_{Cxi} = L_{Ca}$ $H_{Cxi} = H_{Ca} + C_i$ if A or X_i does not have flag F set to one, otherwise $H_{Cxi} = H_{Ca}$ If X_i is not in C_q , then add X_i associated with L_{Cxi} , H_{Cxi} and $PH = A$

into Cq; If Xi associated with LCxi' and HCxi' and PHxi' is in Cq, then If HCxi' > HCxi then replace Xi with HCxi', LCxi' and PHxi' by Xi with HCxi, LCxi and PH=A in Cq; otherwise (i.e., HCxi' == HCxi) if LCxi' > LCxi, then replace Xi with HCxi', LCxi' and PHxi' by Xi with HCxi, LCxi and PH=A in Cq; otherwise (i.e., HCxi' == HCxi and LCxi' == LCxi) if PHxi' > PH, then replace Xi with HCxi', LCxi' and PHxi' by Xi with HCxi, LCxi and PH=A in Cq.

5. Make sure Cq is in a special order. Suppose that Ai (i=1, 2, ..., m) are the nodes in Cq, HCai and LCai are low order cost and high order cost associated with Ai, and IDi is the ID of Ai. One order is that for any k = 1, 2, ..., m-1, HCak < HCaj (j = k+1) or HCak = HCaj and LCak < LCaj or HCak = HCaj and LCak = LCaj and IDk < IDj. Goto step 1.

A.3. Connecting Leaves

Suppose that we have a flooding topology Ft built by one of the algorithms described above. Ft is like a tree. We may connect k (k >= 0) leaves to the tree to have a enhanced flooding topology with more connectivity.

Suppose that there are m (0 < m) leaves directly connected to a node X on the flooding topology Ft. Select k (k <= m) leaves through using a deterministic algorithm or rule. One algorithm or rule is to select k leaves that have smaller or larger IDs (i.e., the IDs of these k leaves are smaller/bigger than the IDs of the other leaves directly connected to node X). Since every node has a unique ID, selecting k leaves with smaller or larger IDs is deterministic.

If k = 1, the leaf selected has the smallest/largest node ID among the IDs of all the leaves directly connected to node X.

For a selected leaf L directly connected to a node N in the flooding topology Ft, select a connection/adjacency to another node from node L in Ft through using a deterministic algorithm or rule.

Suppose that leaf node L is directly connected to nodes Ni (i = 1, 2, ..., s) in the flooding topology Ft via adjacencies and node Ni is not node N, IDi is the ID of node Ni, and Hi (i = 1, 2, ..., s) is the number of hops from node L to node Ni in the flooding topology Ft.

One Algorithm or rule is to select the connection to node Nj (1 <= j <= s) such that Hj is the largest among H1, H2, ..., Hs. If there is another node Na (1 <= a <= s) and Hj = Ha, then select the one with smaller (or larger) node ID. That is that if Hj == Ha and IDj < IDa then select the connection to Nj for selecting the one with smaller

node ID (or if $H_j == H_a$ and $ID_j < ID_a$ then select the connection to N_a for selecting the one with larger node ID).

Suppose that the number of connections in total between leaves selected and the nodes in the flooding topology F_t to be added is N_{Lc} . We may have a limit to N_{Lc} .

Authors' Addresses

Huaimo Chen
Huawei Technologies

Email: huaimo.chen@huawei.com

Dean Cheng
Huawei Technologies

Email: dean.cheng@huawei.com

Mehmet Toy
Verizon
USA

Email: mehmet.toy@verizon.com

Yi Yang
IBM
Cary, NC
United States of America

Email: yyietf@gmail.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: May 16, 2022

U. Chunduri
Intel Corporation
R. Li
Futurewei
R. White
Juniper Networks
L. Contreras
Telefonica
J. Tantsura
Microsoft
Y. Qu
Futurewei
November 12, 2021

Preferred Path Routing (PPR) in IS-IS
draft-chunduri-lsr-isis-preferred-path-routing-07

Abstract

This document specifies a Preferred Path Routing (PPR), a routing protocol mechanism to simplify the path description using IS-IS protocol. PPR builds on existing encapsulation to add the path identity to the packet and supports further extensions along the preferred paths. PPR aims to provide path steering, services and support further extensions along the paths. Preferred path routing is achieved through the addition of path descriptions to the IS-IS advertised prefixes, and mapping those to a PPR data-plane identifier.

Requirements Language

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 RFC2119 [RFC2119], RFC8174 [RFC8174] when, and only when they appear in all capitals, as shown here.

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 <https://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 May 16, 2022.

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 (<https://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. Acronyms	3
2. PPR Details	4
2.1. PPR-ID and Data Plane Extensibility	4
2.2. PPR Path Description	4
2.3. ECMP Considerations	5
3. PPR Related TLVs	5
3.1. PPR-Prefix Sub-TLV	7
3.2. PPR-ID Sub-TLV	8
3.3. PPR-PDE Sub-TLV	9
3.4. PPR-Attributes Sub-TLV	12
4. PPR Processing Procedure Example	13
4.1. PPR TLV Processing	14
4.2. Path Fragments	15
5. PPR Data Plane aspects	15
5.1. SR-MPLS with PPR	15
5.2. PPR Native IP Data Planes	16
5.3. SRv6 with PPR	16
6. Acknowledgements	17
7. IANA Considerations	17
7.1. PPR Sub-TLVs	17
7.2. IGP Parameters	18
8. Security Considerations	18
9. Contributing Authors	18

10. References	19
10.1. Normative References	19
10.2. Informative References	19
Appendix A. Appendix	22
A.1. Challenges with Increased SID Depth	22
A.2. Mitigation with MSD	24
Authors' Addresses	25

1. Introduction

PPR involves associating path descriptions to IS-IS advertised prefixes, mapping those to a data-plane identifier and specifying a mechanism to route packets with the abstracted identifier (PPR-ID), as opposed to individual segments on the packet. This is specified in detail in [I-D.chunduri-rtgwg-preferred-path-routing], along with key use cases and deployment scenarios. PPR allows the traffic along an engineered path through the network by replacing the label stack with a path identifier, PPR-ID, in the packet. The PPR-ID can either be a single label or a native destination address. To facilitate the use of a single label to describe an entire path, a new TLV is added to IS-IS, as described below in Section 3.

A PPR could be an SR path, a traffic engineered path computed based on some constraints, an explicitly provisioned Fast Re-Route (FRR) path or a service chained path. A PPR can be signaled by any node, computed by a central controller, or manually configured by an operator. PPR extends the source routing and path steering capabilities to native IP (IPv4 and IPv6) data planes without hardware upgrades; see Section 5.

1.1. Acronyms

EL	- Entropy Label
ELI	- Entropy Label Indicator
LSP	- IS-IS Link State PDU
MPLS	- Multi Protocol Label Switching
MSD	- Maximum SID Depth
MTU	- Maximum Transferrable Unit
NH	- Next-Hop
PPR	- Preferred Path Routing/Route

PPR-ID - Preferred Path Route Identifier, a data plane identifier
SID - Segment Identifier
SPF - Shortest Path First
SR-MPLS - Segment Routing with MPLS data plane
SRH - Segment Routing Header - IPv6 routing Extension header
SRv6 - Segment Routing with IPv6 data plane with SRH
TE - Traffic Engineering

2. PPR Details

2.1. PPR-ID and Data Plane Extensibility

The PPR-ID describes a path through the network. A data plane type and corresponding data plane identifier as specified in Section 3.2 is mapped to PPR-ID to allow data plane extensibility.

For SR-MPLS, PPR-ID is mapped to an MPLS Label/SID and for SRv6, this is mapped to an IPv6-SID. For native IP data planes, this is mapped to either IPv4 or IPv6 address/prefix.

2.2. PPR Path Description

The path identified by the PPR-ID is described as a set of Path Description Elements (PDEs), each of which represents a segment of the path. Each node determines its location in the path as described, and forwards to the next segment/hop or label of the path description (see the Forwarding Procedure Example later in this document).

These PPR-PDEs as defined in Section 3.3, like SR SIDs, can represent topological elements like links/nodes, backup nodes, as well as non-topological elements such as a service, function, or context on a particular node.

A PPR path can be described as a Strict-PPR or a Loose-PPR. In a Strict-PPR all nodes/links on the path are described with SR SIDs for SR data planes or IPv4/IPv6 addresses for native IP data planes. In a Loose-PPR only some of the nodes/links from source to destination are described. More specifics and restrictions around Strict/Loose PPRs are described in respective data planes in Section 5. Each PDE is described as either an MPLS label towards the Next-Hop (NH) in MPLS enabled networks, or as an IP NH, in the case of either

"plain"/"native" IP or SRv6 enabled networks. A PPR path is related to a set of PDEs using the TLVs as specified in Section 3.

2.3. ECMP Considerations

PPR inherently supports Equal Cost Multi Path (ECMP) for both strict and loose paths. If a path is described using nodes, it would have ECMP NHs established for PPR-ID along the path. However, one can avoid ECMP on any segment of the path by pinning the path using a link identifier to the next segment.

3. PPR Related TLVs

This section describes the encoding of PPR TLV. This TLV can be seen as having 4 logical sections viz, encoding of the PPR-Prefix (IS-IS Prefix), encoding of PPR-ID, encoding of path description with an ordered PDE Sub-TLVs and a set of optional PPR attribute Sub-TLVs, which can be used to describe one or more parameters of the path. Multiple instances of this TLV MAY be advertised in IS-IS LSPs with different PPR-ID Type (data plane) and with corresponding PDE Sub-TLVs. The PPR TLV has Type TBD (suggested value xxx), and has the following format:

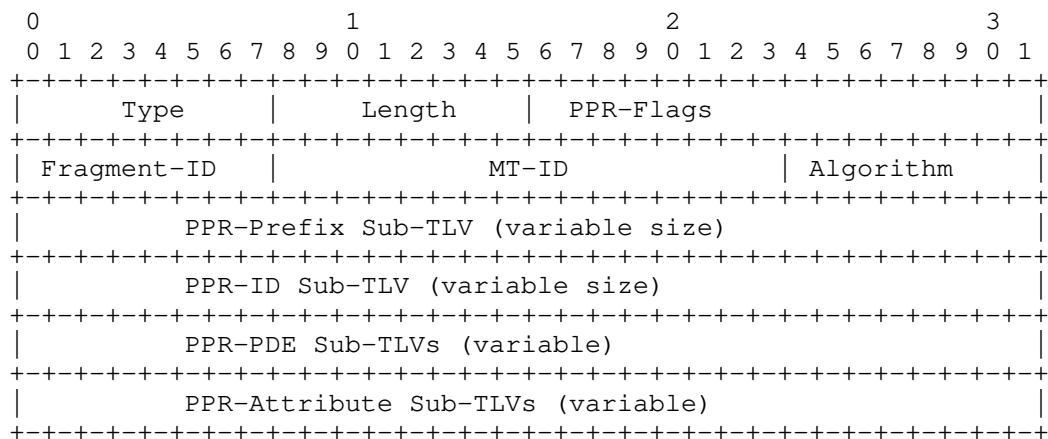


Figure 1: PPR TLV Format

- o Type: 155 (Suggested Value, TBD IANA) from IS-IS top level TLV registry.
- o Length: Total length of the value field in bytes.
- o PPR-Flags: 2 Octet bit-field of flags for this TLV; described below.

- o Fragment-ID: This is an 8-bit Identifier value (0-255) of the TLV fragment. If fragments are not needed to represent the complete path, 'U' bit MUST be set and this value MUST be set to 0.
- o MT-ID: The multi-topology identifier defined in [RFC5120]; the 4 most significant bits MUST be set to 0 on transmit and ignored on receive. The remaining 12-bit field contains the MT-ID.
- o Algorithm: 1 octet value represents the route computation algorithm. Algorithm registry is as defined in [RFC8667]. Computation towards PPR-ID (Section 3.2) happens per MT-ID/ Algorithm pair.
- o PPR-Prefix: A variable size Sub-TLV representing the destination of the path being described. This is defined in Section 3.1.
- o PPR-ID: A variable size Sub-TLV representing the data plane or forwarding identifier of the PPR. Defined in Section 3.2.
- o PPR-PDEs: Variable number of ordered PDE Sub-TLVs which represents the path. This is defined in Section 3.3.
- o PPR-Attributes: Variable number of PPR-Attribute Sub-TLVs which represent the path attributes. These are defined in Section 3.4.

The Flags field has the following flag bits defined:

PPR TLV Flags Format

```

      0 1 2 3 4 5 6 7                               15
      +---+---+---+---+---+---+---+---+---+---+---+---+
      |F|D|A|U|Reserved                               |
      +---+---+---+---+---+---+---+---+---+---+---+---+

```

1. F: Flood bit. If set, the PPR TLV MUST be flooded across the entire routing domain. If the F bit is not set, the PPR TLV MUST NOT be leaked between IS-IS levels. This bit MUST NOT be altered during the TLV leaking
2. D: Down Bit. When the PPR TLV is leaked from IS-IS level-2 to level-1, the D bit MUST be set. Otherwise, this bit MUST be clear. PPR TLVs with the D bit set MUST NOT be leaked from level-1 to level-2. This is to prevent TLV looping across levels.
3. A: Attach bit. The originator of the PPR TLV MUST set the A bit in order to signal that the prefix and PPR-ID advertised in the

PPR TLV are directly connected to the originators. If this bit is not set, this allows any other node in the network to advertise this TLV on behalf of the originating node of the PPR-Prefix. If PPR TLV is leaked to other areas/levels the A-flag MUST be cleared. In case if the originating node of the prefix must be disambiguated for any reason including, if it is a Multi Homed Prefix (MHP) or leaked to a different IS-IS level or because [RFC7794] X-Flag is set, then PPR-Attribute Sub-TLV Source Router ID SHOULD be included.

4. U: Ultimate fragment bit. bit MUST be set if a path has only one fragment or if it is the last Fragment of the path. PPR-ID value for all fragments of the same path MUST be the same.
5. Reserved: For future use; MUST be set to 0 on transmit and ignored on receive.

PPR path description for each IS-IS level is computed and given to one of the nodes for L1 and L2 respectively. Similarly path information when crossing the level boundaries MUST be relevant to the destination level. If there is no path information available for the destination level PPR TLV MUST NOT be leaked regardless of F and D bits as defined above.

The following Sub-TLVs draw from a new registry for Sub-TLV numbers as specified in Section 7.1 and Section 7.2.

3.1. PPR-Prefix Sub-TLV

The structure of PPR-Prefix is:

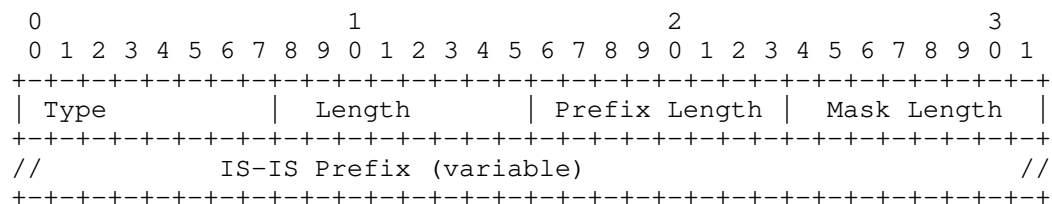


Figure 2: PPR-Prefix Sub-TLV Format

- o Type: 1 (IANA to assign from Sub-TLV registry described above).
- o Length: Total length of the value field in bytes.
- o Prefix Length: The length of the IS-IS Prefix being encoded in bytes. For IPv4 it MUST be 4 and IPv6 it MUST be 16 bytes.

- o Mask Length: The length of the prefix in bits. Only the most significant octets of the Prefix are encoded.
- o IS-IS Prefix: The IS-IS prefix at the tail-end of the advertised PPR. This corresponds to a routable prefix of the originating node and it MAY have one of the [RFC7794] flags set (X-Flag/R-Flag/N-Flag) in the IS-IS reachability TLVs. Length of this field MUST be as per "Prefix Length". Encoding is same as TLV 135 [RFC5305] and TLV 236 [RFC5308] or MT-Capable [RFC5120] IPv4 (TLV 235) and IPv6 Prefixes (TLV 237) respectively.

3.2. PPR-ID Sub-TLV

This is the actual data plane identifier in the packet header and could be of any data plane as defined in the PPR-ID Type field. Both PPR-Prefix and PPR-ID belongs to a same node in the network.

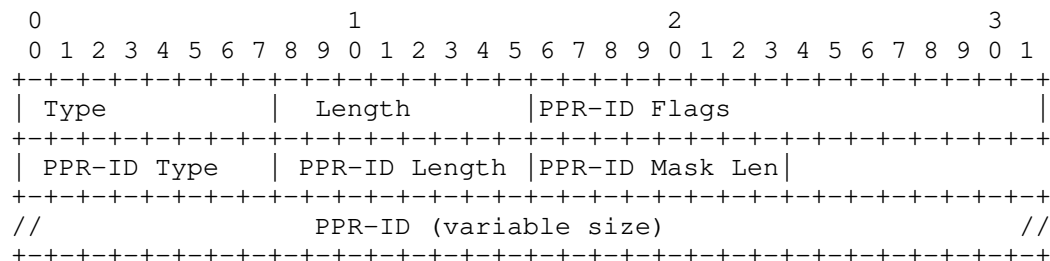
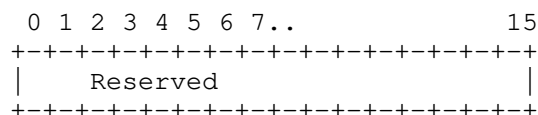


Figure 3: PPR-ID Sub-TLV Format

- o Type: 2 (IANA to assign from Sub-TLV registry described above).
- o Length: Total length of the value field in bytes.
- o PPR-ID Flags: 2 Octet field for PPR-ID flags:

PPR-ID Flags Format



Reserved: For future use; MUST be set to 0 on transmit and ignored on receive.

- o PPR-ID Type: Data plane type of PPR-ID. This is a new registry (TBD IANA - Suggested values as below) for this Sub-TLV and the defined types are as follows:
 - Type: 1 SR-MPLS SID/Label
 - Type: 2 Native IPv4 Address/Prefix
 - Type: 3 Native IPv6 Address/Prefix
 - Type: 4 IPv6 SID in SRv6 with SRH
- o PPR-ID Length: Length of the PPR-ID field in octets and this depends on the PPR-ID type.
- o PPR-ID Mask Len: It is applicable only for PPR-ID Type 2, 3 and 4. For Type 1 this value MUST be set to zero. It contains the length of the PPR-ID Prefix in bits. Only the most significant octets of the Prefix are encoded. This is needed, if PPR-ID followed by an IPv4/IPv6 Prefix instead of 4/16 octet Address respectively.
- o PPR-ID: This is the Preferred Path forwarding identifier that would be on the data packet. The value of this field is variable and it depends on the PPR-ID Type - for Type 1, this is encoded as SR-MPLS SID/Label. For Type 2 this is a 4 byte IPv4 address. For Type 3 this is a 16 byte IPv6 address. For Type 2 and Type 3 encoding is similar to "IS-IS Prefix" as specified in Section 3.1. For Type 4, this is encoded as 16 byte SRv6 SID.

For PPR-ID Type 2, 3 or 4, PPR-ID MUST NOT be advertised as a routable prefix in TLV 135, TLV 235, TLV 236 and TLV 237. PPR-ID MUST belong to the node, from where the PPR-Prefix (Section 3.1) is advertised.

3.3. PPR-PDE Sub-TLV

This Sub-TLV represents the PPR Path Description Element (PDE). PPR-PDEs are used to describe the path in the form of a set of contiguous and ordered Sub-TLVs, where first Sub-TLV represents (the top of the stack in MPLS data plane or) first node/segment of the path. These sets of ordered Sub-TLVs can have both topological elements and non-topological elements (e.g., service segments).

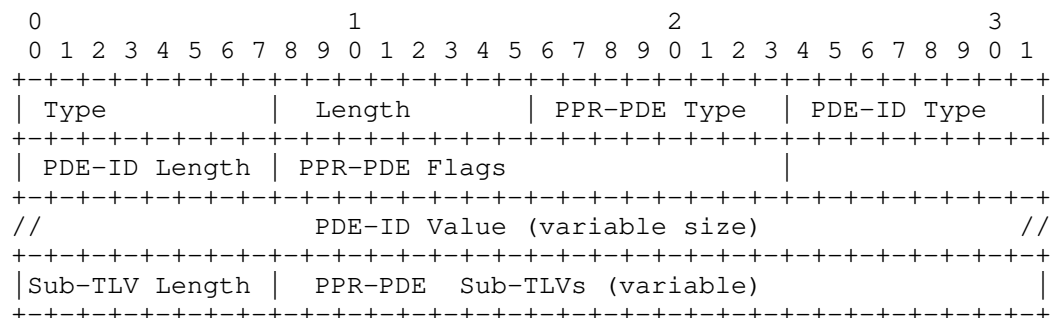


Figure 4: PPR-PDE Sub-TLV Format

- o Type: 3 (See IANA for suggested value) from IS-IS PPR TLV Section 3 Sub-TLV registry.
- o Length: Total length of the value field in bytes.
- o PPR-PDE Type: A new registry (TBD IANA) for this Sub-TLV and the defined types are as follows:

Type: 1 Topological

Type: 2 Non-Topological

- o PDE-ID Type: 1 Octet PDE-forwarding IDentifier Type. A new registry (Suggested Values as listed, IANA TBD) for this Sub-TLV and the defined types and corresponding PDE-ID Length, PDE-ID Value are as follows:

Type 0: This value MUST be set only when PPR-PDE Type is Non-Topological. PDE-ID Length indicates the length of the PDE-ID Value field in bytes. For this type, PDE-ID value represents a service/function. This information is provisioned on the immediate topological PDE preceding this PDE based on the 'E' bit.

Type 1: SID/Label type as defined in [RFC8667]. PDE-ID Length and PDE-ID Value fields are per Section 2.3 of the referenced document.

Type 2: SR-MPLS Prefix SID. PDE-ID Length and PDE-ID Value are same as Type 1.

Type 3: SR-MPLS Adjacency SID. PDE-ID Length and PDE-ID Value are same as Type 1.

Type 4: IPv4 Node Loopback Address. PDE-ID Length 4 bytes and PDE-ID Value is full 4 bytes IPv4 address encoded as specified in "4-octet IPv4 address" of Sub-TLV 6/TLV 22 in [RFC5305].

Type 5: IPv4 Interface Address. PDE-ID Length is 4 bytes and PDE-ID Value is full 4 bytes IPv4 address encoded as specified in "4-octet IPv4 address" of Sub-TLV 6/TLV 22 in [RFC5305]. This PDE-ID in the path description represents the egress interface of the path segment and corresponding adjacency is set as nexthop for the PPR-ID.

Type 6: IPv6 Node Loopback Address. PDE-ID Length and PDE-ID Value are encoded as specified in "Prefix Len" and "prefix" portion of TLV 236 in [RFC5308] respectively.

Type 7: IPv6 Interface Address. PDE-ID Length is 16 bytes and PDE-ID Value is full 16 bytes IPv6 address encoded as specified in "Interface Address 1" portion of TLV 232 in [RFC5308]. This PDE-ID in the path description represents the egress interface of the path segment and corresponding adjacency is set as nexthop for the PPR-ID.

Type 8: SRv6 Node SID as defined in [I-D.ietf-lsr-isis-srv6-extensions]. PDE-ID Length and PDE-ID Value are as defined in SRv6 SID from the referenced draft.

Type 9: SRv6 Adjacency-SID. PDE-ID Length and PDE-ID Values are similar to SRv6 Node SID above.

- o PPR-PDE Flags: 2 Octet bit-field of flags; described below:

PPR-PDE Flags Format

```

0 1 2 3 4 5 6 7 .. 15
+---+---+---+---+---+---+---+---+---+---+---+---+
|L|N|E|           Reserved           |
+---+---+---+---+---+---+---+---+---+---+---+---+

```

L: Loose Bit. Indicates the type of next "Topological PDE-ID" in the path description. This bit MUST be set for only Node/Prefix PDE type. If this flag is unset, the next Topological PDE is Strict Type.

N: Node Bit. By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is 1 i.e., Topological and this PDE represents the node, where PPR-Prefix (Section 3.1) belongs to

(if there is no further PDE specific Sub-TLVs to override PPR-Prefix and PPR-ID values).

E: Egress Bit. By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is 2 i.e., Non-Topological and the service needs to be applied on the egress side of the topological PDE preceding this PDE.

Reserved: Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

- o Sub-TLV Length: 1 byte length of all Sub-TLVs followed. It MUST be set to 0 if no further Sub-TLVs are present.
- o PPR-PDE Sub-TLVs: These have 1 octet type, 1 octet length and value field is defined per the type field. Types are as defined in PPR-TLV Sub-TLVs (Section 7), encoded further as sub-sub-TLVs of PPR-PDE and the length field represents the total length of the value field in bytes.

IS-IS System-ID Sub-TLV: Type 4 (IANA TBD), Length Total length of value field in bytes, Value: IS-IS System-ID of length "ID Length" as defined in [ISO.10589.1992]. This Sub-TLV MUST NOT be present, if the PPR-PDE Type is not Topological. Though the type for this comes from the PPR Sub-TLV registry, here this is a sub-sub-TLV and is part of PPR-ID/PPR-PDE Sub-TLV.

3.4. PPR-Attributes Sub-TLV

PPR-Attribute Sub-TLVs describe the attributes of the path. This document defines the following optional PPR-Attribute Sub-TLVs:

- o Type 5 (Suggested Value - IANA TBD): PPR-Prefix originating node's IPv4 Router ID Sub-TLV. Length and Value fields are as specified in [RFC7794].
- o Type 6 (Suggested Value - IANA TBD): PPR-Prefix originating node's IPv6 Router ID Sub-TLV. Length and Value fields are as specified in [RFC7794].
- o Type 7 (Suggested Value - IANA TBD): PPR-Metric Sub-TLV. Length 4 bytes, and Value is the metric of this path represented through the PPR-ID. Different nodes can advertise the same PPR-ID for the same Prefix with a different set of PPR-PDE Sub-TLVs and the receiving node MUST consider the lowest metric value.

4. PPR Processing Procedure Example

As specified in [I-D.chunduri-rtgwg-preferred-path-routing], a PPR can be a TE path, locally provisioned by the operator or by a controller. Consider the following IS-IS network to describe the operation of PPR TLV as defined in Section 3:

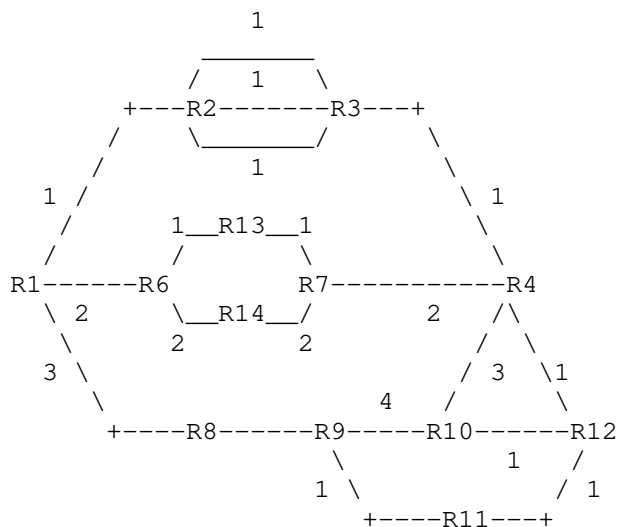


Figure 5: IS-IS Network

In the (Figure 5), consider node R1 as an ingress node, or a head-end node, and the node R4 may be an egress node or another head-end node. The numbers shown on links between nodes indicate the bi-directional IS-IS metric as provisioned. R1 may be configured to receive TE source routed path information from a central entity (PCE [RFC5440], Netconf [RFC6241] or a Controller) that comprises PPR information which relates to sources that are attached to R1. It is also possible to have a PPR provisioned locally by the operator for non-TE needs (e.g. FRR or for chaining certain services).

The PPR TLV (as specified in Section 3) is encoded as an ordered list of PPR-PDEs from source to a destination node in the network and is represented with a PPR-ID (Section 3.2). The PPR TLV includes PPR-PDE Sub-TLVs Section 3.3, which represent both topological and non-topological elements and specifies the actual path towards a PPR-Prefix at R4.

- o The shortest path towards R4 from R1 are through the following sequence of nodes: R1-R2-R3-R4 based on the provisioned metrics.

- o The central entity can define a few PPRs from R1 to R4 that deviate from the shortest path based on other network characteristic requirements as requested by an application or service. For example, the network characteristics or performance requirements may include bandwidth, jitter, latency, throughput, error rate, etc.
- o A first PPR may be identified by PPR-ID = 1 (value) and may include the path of R1-R6-R7-R4 for a Prefix advertised by R4. This is an example for a Loose-PPR and 'L' bit MUST be set appropriately at Section 3.3.
- o A second PPR may be identified by PPR-ID = 2 (value) and may include the path of R1-R8-R9-R10-R4. This is an example for a Strict-PPR and 'L' bit MUST be unset appropriately at Section 3.3. Though this example shows PPR with all nodal SIDs, it is possible to have a PPR with combination of node and adjacency SIDs (local or global) or with PPR-PDE Type set to Non-Topological as defined in Section 3.3 elements along with these.

4.1. PPR TLV Processing

The first topological sub-object or PDE (Section 3.3) relative to the beginning of PPR Path contains the information about the first node (e.g. in SR-MPLS it's the topmost label). The last topological sub-object or PDE contains information about the last node (e.g. in SR-MPLS it's the bottommost label).

Each receiving node determines whether an advertised PPR includes information regarding the receiving node. Before processing any further, validation MUST be done to see if any PPR topological PDE is seen more than once (possible loop), if yes, this PPR TLV MUST be ignored. Processing of PPR TLVs may be done, during the end of the SPF computation (for MTID that is advertised in this TLV) and for each prefix described through PPR TLV. For example, node R9 receives the PPR information, and ignores the PPR-ID=1 (Section 4) because this PPR TLV does not include node R9 in the path description/ordered PPR-PDE list.

However, node R9 may determine that the second PPR identified by PPR-ID = 2 does include the node R9 in its PDE list. Therefore, node R9 updates the local forwarding database to include an entry for the destination address that R4 indicates, so that when a data packet comprising a PPR-ID of 2 is received, forward the data packet to node R10 instead of R11. This is done, even though from R9 the shortest path to reach R4 via R11 (Cost 3: R9-R11-R12-R4) it chooses the NH to R10 to reach R4 as specified in the PPR path description. Same

process happens to all nodes or all topological PDEs as described in the PPR TLV.

In summary, the receiving node checks first, if this node is on the path by checking the node's topological elements (with PPR-PDE Type set to Topological) in the path list. If yes, it adds/adjusts the PPR-ID's shortest path NH towards the next topological PDE in the PPR's Path.

4.2. Path Fragments

A complete PPR path may not fit into the maximum allowable size of the IS-IS TLV. To overcome this a 7 bit Fragment-ID field is defined in Section 3 . With this, a single PPR path is represented via one or more fragmented PPR path TLVs, with all having the same PPR-ID. Each fragment carries the PPR-ID as well as a numeric Fragment-ID from 0 to (N-1), when N fragments are needed to describe the PPR Graph (where N>1). In this case Fragment (N-1) MUST set the 'U' bit (PPR-Flags) to indicate it is the last fragment. If Fragment-ID is non-zero in the TLV, then it MUST not carry PPR-Prefix Sub-TLV. The optional PPR Attribute Sub-TLVs which describe the path overall MUST be included in the last fragment only (i.e., when the 'U' bit is set).

5. PPR Data Plane aspects

Data plane for PPR-ID is selected by the entity (e.g., a controller, locally provisioned by operator), which selects a particular PPR in the network. Section 3.2 defines various data plane identifier types and a corresponding data plane identifier is selected by the entity which selects the PPR.

5.1. SR-MPLS with PPR

If PPR-ID Type is 1, then the PPR belongs to SR-MPLS data plane and the complete PPR stack is represented with a unique SR SID/Label and this gets programmed on the data plane of each node, with the appropriate NH computed as specified in Section 4. PPR-ID here is a label/index from the SRGB (like another node SID or global ADJ-SID). PPR path description here is a set of ordered SIDs represented with PPR-PDE (Section 3.2) Sub-TLVs. Non-Topological segments are also programmed in the forwarding to enable specific function/service, when the data packet hits with corresponding PPR-ID.

Based on the 'L' flag in PPR-ID Flags (Section 3.2), for SR-MPLS data plane either 1 label or 2 labels need to be provisioned on individual nodes on the path description. For the example network in Section 4, for PPR-ID=1, which is a loose path, node R6 programs the bottom

label as PPR-ID and the top label as the next topological PPR-PDE in the path, which is a node SID of R7. The NH computed at R6 would be the shortest path towards R7 i.e., the interface towards R13. If 'L' flag is unset, only PPR-ID is programmed on the data plane with NH set to the shortest path towards the next topological PPR-PDE.

5.2. PPR Native IP Data Planes

If PPR-ID Type is 2 then source routing and packet steering can be done in IPv4 data plane (PPR-IPv4), along the path as described in PPR Path description. This is achieved by setting the destination IP address as PPR-ID, which is an IPv4 address in the data packet (tunneled/encapsulated). There is no data plane change or upgrade needed to support this.

Similarly for PPR-ID Type is 3, then source routing and packet steering can be done in the IPv6 data plane (PPR-IPv6), along the path as described in PPR Path description. Whatever specified above for IPv4 applies here too, except that the destination IP address of the data packet is an IPv6 Address (PPR-ID). This doesn't require any IPv6 extension headers (EH), if there is no metadata/TLVs need to be carried in the data packet.

Based on 'L' flag in PPR-ID Flags (Section 3.2), for PPR-ID Type 2 or 3 (Native IPv4 or IPv6 data planes respectively) the packet has to be encapsulated using the capabilities (either dynamically signaled through [I-D.ietf-isis-encapsulation-cap] or statically provisioned on the nodes) of the next loose PDE in the path description.

For the example network in Section 4, for PPR-ID=1, which is a loose path, node R6 programs to encapsulate the data packet towards the next loose topological PPR-PDE in the path, which is R7. The NH computed at R6 would be the shortest path towards R7 i.e., the interface towards R13. If the 'L' flag is unset, only PPR-ID is programmed on the data plane with NH set to the shortest path towards the next topological PPR-PDE, with no further encapsulation of the data packet.

5.3. SRv6 with PPR

If PPR-ID Type is 4, the PPR belongs to SRv6 with SRH data plane and the complete PPR stack is represented with IPv6 SIDs and this gets programmed on the data plane with the appropriate NH computed as specified in Section 4. PPR-ID here is a SRv6 SID. PPR path description here is a set of ordered SID TLVs similar to as specified in Section 5.1. One way PPR-ID would be used in this case is by setting it as the destination IPv6 address and SL field in SRH would

be set to 0; however SRH [RFC8754] can contain any other TLVs and non-topological SIDs as needed.

6. Acknowledgements

Thanks to Alex Clemm, Lin Han, Toerless Eckert, Asit Chakraborti, Stewart Bryant and Kiran Makhijani for initial discussions on this topic. Thanks to Kevin Smith and Stephen Johnson for various deployment scenarios applicability from ETSI WGs perspective. Authors also acknowledge Alexander Vainshtein for detailed discussions and few suggestions on this topic.

Earlier versions of [RFC8667] have a mechanism to advertise EROs through Binding SID.

7. IANA Considerations

This document requests the following new TLV in IANA IS-IS TLV code-point registry.

TLV #	Name
-----	-----
155	PPR TLV (Suggested Value, IANA TBD)

7.1. PPR Sub-TLVs

This document requests IANA to create a new Sub-TLV registry for PPR TLV Section 3 with the following initial entries (suggested values). Though these are defined as Sub-TLVs of PPR TLV, these can be part of another Sub-TLV as a nested sub-sub-TLV (e.g. IS-IS System-ID).

Sub-TLV #	Sub-TLV Name
-----	-----
1	PPR-Prefix (Section 3.1)
2	PPR-ID (Section 3.2)
3	PPR-PDE (Section 3.3)
4	IS-IS System-ID (Section 3.3)
5	PPR-Prefix Source IPv4 Router ID (Section 3.4)
6	PPR-Prefix Source IPv6 Router ID (Section 3.4)
7	PPR-Metric (Section 3.4)

7.2. IGP Parameters

This document requests additional IANA registries in an IANA managed registry "Interior Gateway Protocol (IGP) Parameters" for various PPR TLV parameters. The registration procedure is based on the "Expert Review" as defined in [RFC8126]. The suggested registry names are:

- o "PPR-Type" - Types are unsigned 8 bit numbers. Values are as defined in Section 3 of this document.
- o "PPR-Flags" - 1 Octet. Bits as described in Section 3 of this document.
- o "PPR-ID Type" - Types are unsigned 8 bit numbers. Values are as defined in Section 3.2 of this document.
- o "PPR-ID Flags" - 1 Octet. Bits as described in Section 3.2 of this document.
- o "PPR-PDE Type" - Types are unsigned 8 bit numbers. Values are as defined in Section 3.3 of this document.
- o "PPR-PDE Flags" - 1 Octet. Bits as described in Section 3.3 of this document.
- o "PDE-ID Type" - Types are unsigned 8 bit numbers. Values are as defined in Section 3.3 of this document.

8. Security Considerations

Security concerns for IS-IS are addressed in [RFC5304] and [RFC5310]. Further security analysis for the IS-IS protocol is done in [RFC7645] with detailed analysis of various security threats and why [RFC5304] should not be used in the deployments. Advertisement of the additional information defined in this document introduces no new security concerns in IS-IS protocol. However, for extensions related to SR-MPLS and SRH data planes, those particular data plane security considerations do apply here.

9. Contributing Authors

The following people contributed substantially to the content of this document and should be considered co-authors.

Yingzhen Qu
Futurewei
2330 Central Expressway
Santa Clara
CA 95050
USA
Email: yingzhen.qu@futurewei.com

10. References

10.1. Normative References

- [I-D.chunduri-rtgwg-preferred-path-routing]
Bryant, S., Chunduri, U., and A. Clemm, "Preferred Path Routing Framework", draft-chunduri-rtgwg-preferred-path-routing-01 (work in progress), October 2021.
- [ISO.10589.1992]
International Organization for Standardization,
"Intermediate system to intermediate system intra-domain-routing routine information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)",
ISO Standard 10589, 1992.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119,
DOI 10.17487/RFC2119, March 1997,
<<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174,
May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

10.2. Informative References

- [I-D.bryant-rtgwg-plfa]
Bryant, S., Chunduri, U., and T. Eckert, "Preferred Path Loop-Free Alternate (pLFA)", draft-bryant-rtgwg-plfa-02 (work in progress), June 2021.
- [I-D.chunduri-dmm-5g-mobility-with-ppr]
Chunduri, U., Contreras, L. M., Bhaskaran, S., Tantsura, J., and P. Muley, "Transport aware 5G mobility with PPR", draft-chunduri-dmm-5g-mobility-with-ppr-00 (work in progress), November 2020.

- [I-D.ietf-dmm-tn-aware-mobility]
Chunduri, U., Kaippallimalil, J., Bhaskaran, S., Tantsura, J., and P. Muley, "Mobility aware Transport Network Slicing for 5G", draft-ietf-dmm-tn-aware-mobility-02 (work in progress), October 2021.
- [I-D.ietf-isis-encapsulation-cap]
Xu, X., Decraene, B., Raszuk, R., Chunduri, U., Contreras, L. M., and L. Jalil, "Advertising Tunnelling Capability in IS-IS", draft-ietf-isis-encapsulation-cap-01 (work in progress), April 2017.
- [I-D.ietf-isis-mpls-elc]
Xu, X., Kini, S., Psenak, P., Filsfils, C., Litkowski, S., and M. Bocci, "Signaling Entropy Label Capability and Entropy Readable Label Depth Using IS-IS", draft-ietf-isis-mpls-elc-13 (work in progress), May 2020.
- [I-D.ietf-lsr-isis-srv6-extensions]
Psenak, P., Filsfils, C., Bashandy, A., Decraene, B., and Z. Hu, "IS-IS Extensions to Support Segment Routing over IPv6 Dataplane", draft-ietf-lsr-isis-srv6-extensions-18 (work in progress), October 2021.
- [I-D.ietf-mpls-sfc]
Farrel, A., Bryant, S., and J. Drake, "An MPLS-Based Forwarding Plane for Service Function Chaining", draft-ietf-mpls-sfc-07 (work in progress), March 2019.
- [I-D.ietf-teas-enhanced-vpn]
Dong, J., Bryant, S., Li, Z., Miyasaka, T., and Y. Lee, "A Framework for Enhanced Virtual Private Network (VPN+) Services", draft-ietf-teas-enhanced-vpn-09 (work in progress), October 2021.
- [I-D.xuclad-spring-sr-service-chaining]
Clad, F., Xu, X., Filsfils, C., Bernier, D., Li, C., Decraene, B., Ma, S., Yadlapalli, C., Henderickx, W., and S. Salsano, "Segment Routing for Service Chaining", draft-xuclad-spring-sr-service-chaining-01 (work in progress), March 2018.
- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.

- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5308] Hopps, C., "Routing IPv6 with IS-IS", RFC 5308, DOI 10.17487/RFC5308, October 2008, <<https://www.rfc-editor.org/info/rfc5308>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC5440] Vasseur, JP., Ed. and JL. Le Roux, Ed., "Path Computation Element (PCE) Communication Protocol (PCEP)", RFC 5440, DOI 10.17487/RFC5440, March 2009, <<https://www.rfc-editor.org/info/rfc5440>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6790] Kompella, K., Drake, J., Amante, S., Henderickx, W., and L. Yong, "The Use of Entropy Labels in MPLS Forwarding", RFC 6790, DOI 10.17487/RFC6790, November 2012, <<https://www.rfc-editor.org/info/rfc6790>>.
- [RFC7645] Chunduri, U., Tian, A., and W. Lu, "The Keying and Authentication for Routing Protocol (KARP) IS-IS Security Analysis", RFC 7645, DOI 10.17487/RFC7645, September 2015, <<https://www.rfc-editor.org/info/rfc7645>>.
- [RFC7794] Ginsberg, L., Ed., Decraene, B., Previdi, S., Xu, X., and U. Chunduri, "IS-IS Prefix Attributes for Extended IPv4 and IPv6 Reachability", RFC 7794, DOI 10.17487/RFC7794, March 2016, <<https://www.rfc-editor.org/info/rfc7794>>.
- [RFC8126] Cotton, M., Leiba, B., and T. Narten, "Guidelines for Writing an IANA Considerations Section in RFCs", BCP 26, RFC 8126, DOI 10.17487/RFC8126, June 2017, <<https://www.rfc-editor.org/info/rfc8126>>.

- [RFC8402] Filsfils, C., Ed., Previdi, S., Ed., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", RFC 8402, DOI 10.17487/RFC8402, July 2018, <<https://www.rfc-editor.org/info/rfc8402>>.
- [RFC8491] Tantsura, J., Chunduri, U., Aldrin, S., and L. Ginsberg, "Signaling Maximum SID Depth (MSD) Using IS-IS", RFC 8491, DOI 10.17487/RFC8491, November 2018, <<https://www.rfc-editor.org/info/rfc8491>>.
- [RFC8667] Previdi, S., Ed., Ginsberg, L., Ed., Filsfils, C., Bashandy, A., Gredler, H., and B. Decraene, "IS-IS Extensions for Segment Routing", RFC 8667, DOI 10.17487/RFC8667, December 2019, <<https://www.rfc-editor.org/info/rfc8667>>.
- [RFC8754] Filsfils, C., Ed., Dukes, D., Ed., Previdi, S., Leddy, J., Matsushima, S., and D. Voyer, "IPv6 Segment Routing Header (SRH)", RFC 8754, DOI 10.17487/RFC8754, March 2020, <<https://www.rfc-editor.org/info/rfc8754>>.

Appendix A. Appendix

A.1. Challenges with Increased SID Depth

SR label stacks carried in the packet header create challenges in the design and deployment of networks and networking equipment. Following examples illustrates the need for increased SID depth in various use cases:

(a). Consider the following network where SR-MPLS data plane is in use and with same SRGB (5000-6000) on all nodes i.e., A1 to A11 and B1 to B7 for illustration:

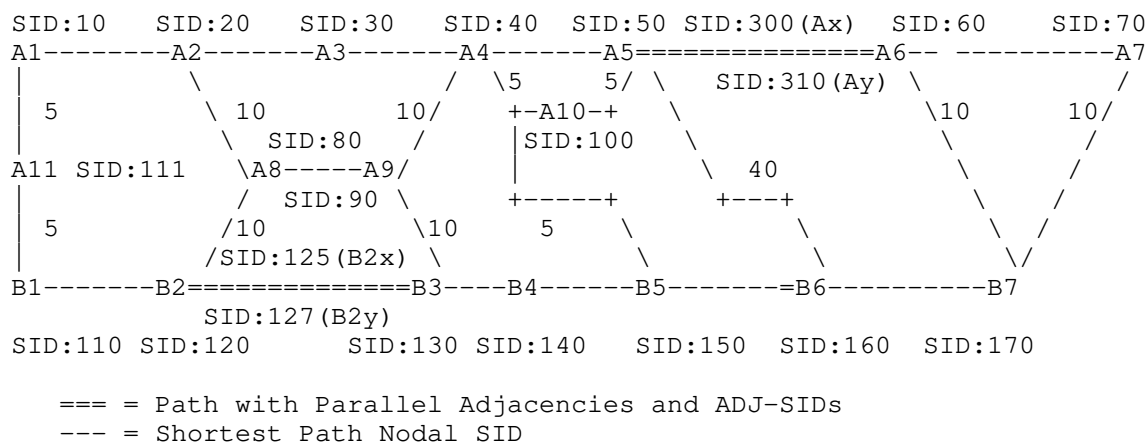


Figure 6: SR-MPLS Network

Global ADJ-SIDs are provisioned between A5-A6 and B2-B3 (with parallel adjacencies). All other SIDs shown are nodal SID indices.

All metrics of the links are set to 1, unless marked otherwise.

Shortest Path from A1 to A7: A2-A3-A4-A5-A6-A7

Path-x: From A1 to A7 - A2-A8-B2-B2x-A9-A10-Ax-A7; Pushed Label Stack @A1: 5020:5080:5120:5125:5090:5100:5300:5070 (where B2x is a local ADJ-SID and Ax is a global ADJ-SID).

In this example, the traffic engineered path is represented with a combination of Adjacency and Node SIDs with a stack of 8 labels. However, this value can be larger, if the use of entropy label [RFC6790] is desired and based on the Readable Label Depth (Appendix A.2) capabilities of each node and additional labels required to insert ELI/EL at appropriate places.

Though above network is shown with SR-MPLS data plane, if the network were to use SRv6 data plane, path size would be increased even more because of the size of the IPv6 SID (16 bytes) in SRH.

(b). Apart from the TE case above, when deploying [I-D.ietf-mpls-sfc] or [I-D.xuclad-spring-sr-service-chaining], with the inclusion of services, or non-topological segments on the label stack, can also make the size of the stack much larger.

Overall the additional path overhead in various SR deployments may cause the following issues:

- a. **HW Capabilities:** Not all nodes in the path can support the ability to push or read label stack (with additional non-topological and special labels) needed to satisfy user/operator requirements. Alternate paths, which meet these user/operator requirements may not be available.
- b. **Line Rate:** Potential performance issues in deployments, which use data plane with extension header as both size of the SIDs in the extension header and the fixed extension header size itself needs to be factored by the hardware.
- c. **MTU:** Larger SID stacks on the data packet can cause potential MTU/fragmentation issues (SRH).
- d. **Header Tax:** Some deployments, such as 5G, require minimal packet overhead in order to conserve network resources. Carrying 40 or 50 octets of data in a packet with hundreds of octet of header would be an unacceptable use of available bandwidth.

With the solution proposed in this document, for Path-x in Figure 6 above, SID stack would be reduced from 8 SIDs to a single SID without any additional overhead.

A.2. Mitigation with MSD

The number of SIDs in the stack a node can impose is referred as Maximum SID Depth (MSD) capability [RFC8491], which must be taken into consideration when computing a path to transport a data packet in a network implementing segment routing. [I-D.ietf-isis-mpls-elc] defines another MSD type, Readable Label Depth (RLD) that is used by a head-end to insert Entropy Label pair (ELI/EL) at appropriate depth, so it could be read by transit nodes. There are situations where the source routed path can be excessive as path represented by SR SIDs need to describe all the nodes and ELI/EL based on the readability of the nodes in that path. Registries set forth in [RFC8491] applicable for MPLS data plane and for IPv6 data plane with SRH.

MSDs (and RLD type) capabilities advertisement help mitigate the problem for a central entity to create the right source routed path per application/operator requirements. However the availability of actual paths meeting these requirements are still limited by the underlying hardware and their MSD capabilities in the data path.

Authors' Addresses

Uma Chunduri
Intel Corporation

Email: umac.ietf@gmail.com

Richard Li
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
USA

Email: richard.li@futurewei.com

Russ White
Juniper Networks
Oak Island, NC 28465
USA

Email: russ@riw.us

Luis M. Contreras
Telefonica
Sur-3 building, 3rd floor
Madrid 28050
Spain

Email: luismiguel.contrerasmurillo@telefonica.com

Jeff Tantsura
Microsoft

Email: jefftanti.ietf@gmail.com

Yingzhen Qu
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
USA

Email: yingzhen.qu@futurewei.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: September 9, 2020

U. Chunduri
Y. Qu
Futurewei
R. White
Juniper Networks
J. Tantsura
Apstra Inc.
L. Contreras
Telefonica
March 8, 2020

Preferred Path Routing (PPR) in OSPF
draft-chunduri-lsr-ospf-preferred-path-routing-04

Abstract

This document specifies a Preferred Path Routing (PPR), a routing protocol mechanism to simplify the path description of data plane traffic in Segment Routing (SR) deployments with OSPFv2 and OSPFv3 protocols. PPR aims to mitigate the MTU and data plane processing issues that may result from SR packet overheads; and also supports further extensions along the paths. Preferred path routing is achieved through the addition of path descriptions to the OSPF advertised prefixes, and mapping those to a PPR data-plane identifier.

Requirements Language

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 RFC2119 [RFC2119], RFC8174 [RFC8174] when, and only when they appear in all capitals, as shown here".

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 <https://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 September 9, 2020.

Copyright Notice

Copyright (c) 2020 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 (<https://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. Acronyms	3
2. OSPFv2 PPR TLV	4
2.1. PPR-Flags	6
2.2. PPR-Prefix Sub-TLV	6
2.3. PPR-ID Sub-TLV	7
2.4. PPR-PDE Sub-TLV	9
2.5. PPR-Attributes Sub-TLV	11
3. OSPFv3 PPR TLV	12
3.1. OSPFv3 PPR-Prefix Sub-TLV	13
3.2. OSPFv3 PPR-ID Sub-TLVs	14
3.3. OSPFv3 PPR-PDE Sub-TLV	16
3.4. OSPFv3 PPR-Attributes Sub-TLV	19
4. Other Considerations	19
5. Acknowledgements	19
6. IANA Considerations	19
7. Security Considerations	20
8. References	20
8.1. Normative References	20
8.2. Informative References	20
Authors' Addresses	22

1. Introduction

In a network implementing Segment Routing (SR), packets are steered through the network using Segment Identifiers (SIDs) carried in the packet header. Each SID uniquely identifies a segment as defined in [I-D.ietf-spring-segment-routing]. SR capabilities are defined for MPLS and IPv6 data planes called SR-MPLS and SRv6 respectively.

In SR-MPLS, a segment is encoded as a label and an ordered list of segments is encoded as a stack of labels on the data packet. In SRv6, a segment is encoded as an IPv6 address, with in a new type of IPv6 hop-by-hop routing header/extension header (EH) called SRH [I-D.ietf-6man-segment-routing-header], where an ordered list of IPv6 addresses/segments is encoded in SRH.

Preferred path routing can be described as a) enabling route computation based on the specific path described along with the prefix as opposed to shortest path towards the prefix and b) forwarding based on the abstracted path identifier as opposed to the individual segments on the packet. This is also further described in Section 2 of [I-D.chunduri-lsr-isis-preferred-path-routing].

Any prefix advertised with a path description from any node in the network is called PPR. A PPR could be an SR path, an explicitly provisioned Fast Re-Route (FRR) path or a service chained path. A PPR can be signaled by any node, which receives the SR path computed by a central controller, or by statically configuring the same on a node in the network.

The issues caused by the large SID depth, and existing methods for mitigation are introduced in [I-D.chunduri-lsr-isis-preferred-path-routing] in Appendix A.1 and A.2. To mitigate these issues and also to facilitate forwarding plane extensibility, this draft proposes a new OSPFv2 PPR TLV (Section 2), OSPFv3 PPR TLV (Section 3) to use the path with a corresponding data plane identifier.

1.1. Acronyms

EL	- Entropy Label
ELI	- Entropy Label Indicator
MPLS	- Multi Protocol Label Switching
MSD	- Maximum SID Depth
MTU	- Maximum Transferrable Unit

PPR	- Preferred Path Route
SID	- Segment Identifier
SPF	- Shortest Path First
SR	- Segment Routing
SRH	- Segment Routing Header
SR-MPLS	- Segment Routing with MPLS data plane
SRv6	- Segment Routing with Ipv6 data plane with SRH
SRH	- IPv6 Segment Routing Header
TE	- Traffic Engineering

2. OSPFv2 PPR TLV

Extended Prefix Opaque LSAs defined in [RFC7684] are used for advertisements of PPRs. This section describes the encoding of PPR TLV. This TLV can be seen as having 4 logical sections viz., encoding of the OSPFv2 Prefix, encoding of PPR-ID, encoding of path description with an ordered PDE Sub-TLVs and a set of optional PPR attribute Sub-TLVs, which can be used to describe one or more parameters of the path. Multiple OSPF PPR TLVs MAY be advertised in each OSPF Extended Prefix Opaque LSA, but all TLVs included in a single OSPF Extended Prefix Opaque LSA MUST have the same flooding scope.

The PPR TLV has Type TBD (suggested value xxx), and has the following format:

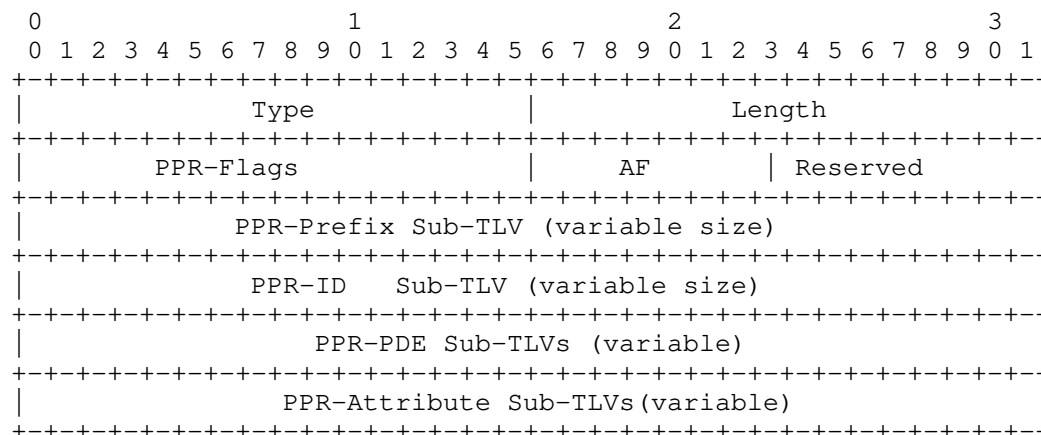
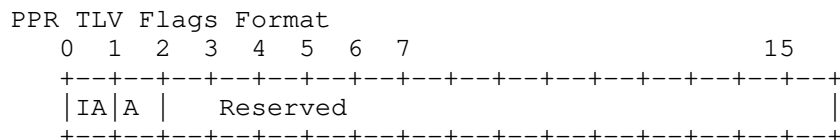


Figure 1: OSPFV2 PPR TLV Format

- o Type - TBD (IANA) from OSPF Extended Prefix Opaque LSA registry.
- o Length - Total length of the value field in bytes (variable).
- o PPR-Flags - 2 Octet flags for this TLV are described below.
- o AF - Address family for the prefix. Currently, the only supported value is 0 for IPv4 unicast. The inclusion of address family in this TLV allows for future extension.
- o Reserved - 1 Octet reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.
- o PPR-Prefix - This is a variable size Sub-TLV, which represents the prefix for which path description is being attached to. This is defined in Section 2.2.
- o PPR-ID - This is a variable size Sub-TLV, which represents the data plane or forwarding identifier of the PPR. This is defined in Section 2.3.
- o PPR-PDEs - Variable number of ordered PDE Sub-TLVs which represents the path. This is defined in Section 2.4.
- o PPR-Attributes - Variable number of PPR-Attribute Sub-TLVs which represent the path attributes. These are defined in Section 2.5.

2.1. PPR-Flags

Flags: 2 octet field of PPR TLV has following flags defined:



w=Where:

IA-Flag: Inter-Area flag. If set, advertisement is of inter-area type. An Area Boarder Router (ABR) that is advertising the OSPF PPR TLV between areas MUST set this bit.

A: The originator of the PPR TLV MUST set the A bit in order to signal that the prefixes and PPR-IDs advertised in the PPR TLV are directly connected to the originators. If this bit is not set, this allows any other node in the network advertise this TLV on behalf of the originating node of the "OSPF Prefix". If PPR TLV is propagated to other areas the A-flag MUST be cleared. In case if the originating node of the prefix has to be disambiguated for any reason including, if it is a Multi Homed Prefix (MHP) or propagated to a different OSPF area, then PPR-Attribute Sub-TLV Source Router ID SHOULD be included.

Reserved: Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

PPR path description for each OSPF area is computed and given to one of the nodes in that area for dissemination. Similarly path information when crossing the area boundaries MUST be relevant to the destination area. If there is no path information available for the destination area, PPR TLV MUST NOT be leaked regardless of the IA bit status.

2.2. PPR-Prefix Sub-TLV

The structure of PPR-Prefix, for which path description is attached to is as follows:

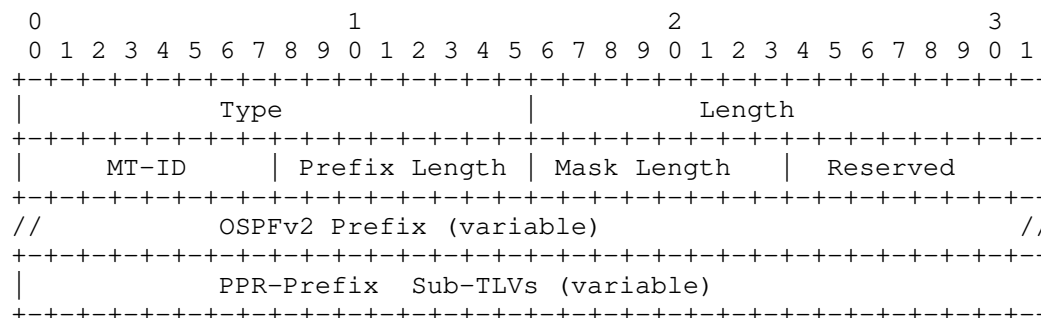


Figure 2: PPR-Prefix Sub-TLV Format

- o Type - 1 (suggested value, IANA TBD) from OSPFv2 PPR TLV Section 2 Sub-TLV registry.
- o Length - Total length of the value field in bytes (variable).
- o MT-ID - Multi-Topology ID (as defined in [RFC4915]).
- o Prefix Len - contains the length of the OSPF prefix being encoded in bytes.
- o Mask Length - The length of the prefix in bits. Only the most significant octets of the Prefix are encoded.
- o OSPFv2 Prefix - represents the OSPFv2 prefix at the tail-end of the advertised PPR. For the address family IPv4 unicast, the prefix itself is encoded as a 32-bit value. The default route is represented by a prefix of length 0.
- o PPR-Prefix Sub-TLVs have 2 octet type, 2 octet length and value field is defined per type.

2.3. PPR-ID Sub-TLV

This represents the actual data plane identifier in the packet and could be of any data plane as defined in PPR-ID-type field. Both OSPF Prefix and PPR-ID MUST belong to a same node in the network.

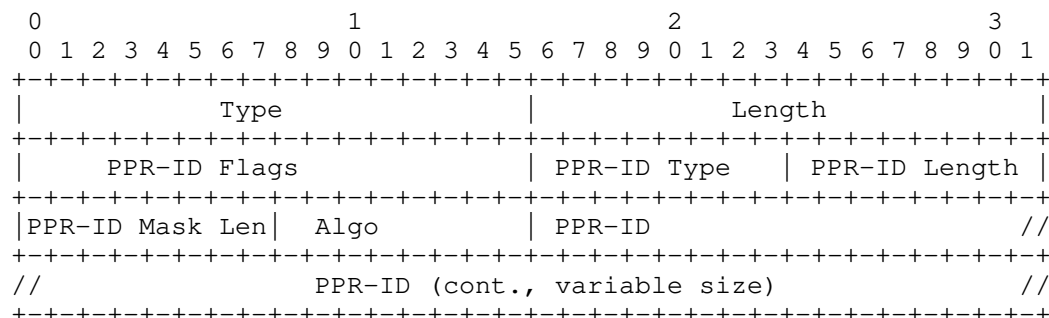


Figure 3: PPR-ID Sub-TLV Format

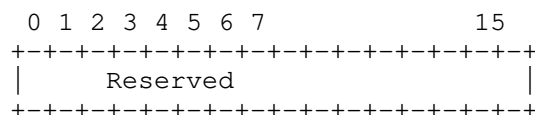
- o Type - 2 (suggested value, IANA TBD) from OSPFv2 PPR TLV Section 2 Sub-TLV registry.
- o Length - Total length of the value field in bytes (variable).
- o PPR-ID Type - Data plane type of PPR-ID. This is a new registry (TBD IANA) for this Sub-TLV and the defined types are as follows:

Type: 1 SR-MPLS SID/Label

Type: 2 Native IPv4 Address/Prefix

- o PPR-ID Flags - 2 Octet field for PPR-ID flags:

PPR-ID Flags Format



Reserved - Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

- o PPR-ID Type - Data plane type of PPR-ID. Values are defined in [I-D.chunduri-lsr-isis-preferred-path-routing]. Only Type 1 and Type 2 are applicable here.
- o PPR-ID Length - Length of the PPR-ID field in octets and this depends on the PPR-ID type. See PPR-ID below for the length of this field and other considerations.

- o PPR-ID Mask Len - It is applicable for only for PPR-ID Type 2. For Type 1 this value MUST be set to zero. It contains the length of the PPR-ID Prefix in bits. Only the most significant octets of the Prefix are encoded. This is needed, if PPR-ID followed is an IPv4 Prefix instead of 4 octet Address respectively.
- o Algo - 1 octet value represents the SPF algorithm. Algorithm registry is as defined in [I-D.ietf-ospf-segment-routing-extensions].
- o PPR-ID - This is the Preferred Path forwarding identifier that would be on the data packet. The value of this field is variable and it depends on the PPR-ID Type - for Type 1, this is encoded as SR-MPLS SID/Label. For Type 2 this is a 4 byte IPv4 address encoded similar to PPR-Prefix.

2.4. PPR-PDE Sub-TLV

This is a new Sub-TLV type in PPR TLV Section 2 and is called as PPR Path Description Element (PDE). PPR-PDEs are used to describe the path in the form of set of contiguous and ordered Sub-TLVs, where first Sub-TLV represents (the top of the stack in MPLS data plane or) first node/segment of the path. These set of ordered Sub-TLVs can have both topological SIDs and non-topological SIDs (e.g., service segments).

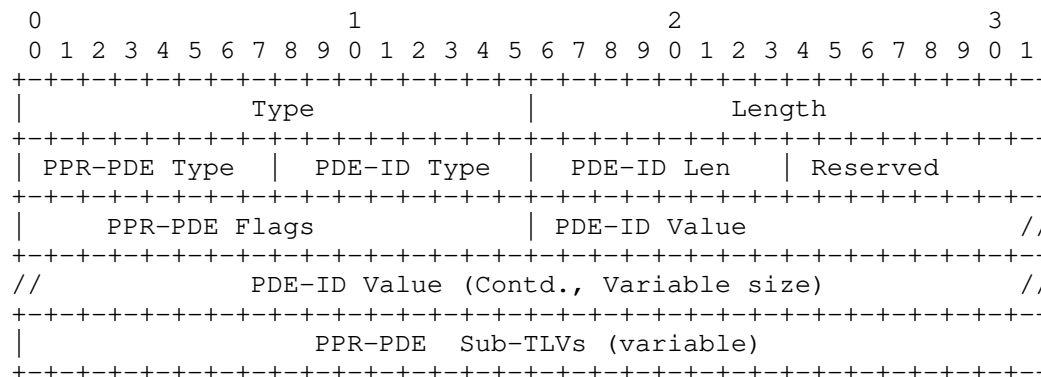


Figure 4: PPR-PDE Sub-TLV Format

- o Type - 3 (TBD IANA, suggested value) from OSPFv2 PPR TLV Section 2 Sub-TLV registry.
- o Length - Total length of the value field in bytes (variable).

- o PPR-PDE Type - This is a new registry (TBD IANA) for this Sub-TLV and the defined types are as follows:
 - Type: 1 Topological
 - Type: 2 Non-Topological
- o PDE-ID Type - 1 Octet PDE-forwarding IDentifier Type. This is a new registry (TBD IANA) for this Sub-TLV and the defined types and corresponding PDE-ID Len, PDE-ID Value are as follows:
 - Type 0: This value MUST be set only when PPR-PDE Type is Non-Topological. PDE-ID Len specified in bytes and encoded in NBO in PDE-ID Value field which can represent a service/function. This information is provisioned on the immediate topological PDE preceding to this PDE based on the 'E' bit.
 - Type 1: SID/Label Sub-TLV as defined in [I-D.ietf-ospf-segment-routing-extensions]. PDE-ID Len and PDE-ID Value fields are per Section 2.1 of the referenced document.
 - Type 2: SR-MPLS Prefix SID. PDE-ID Len and PDE-ID Value are same as Type 1.
 - Type 3: SR-MPLS Adjacency SID. PDE-ID Len and PDE-ID Value are same as Type 1.
 - Type 4: IPv4 Node Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2.
 - Type 5: IPv4 P2P interface Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2.
 - Type 6: IPv4 LAN interface Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2. This type MUST have OSPF Neighbor ID sub-TLV in the PDE.
- o PDE-ID Len - 1 Octet. Length of PDE-ID field.
- o Reserved - 1 Octet reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.
- o PPR-PDE Flags - 2 Octet flags for this TLV are described below:

PPR-PDE Flags Format

```

      0 1 2 3 4 5 6 7...          15
+---+---+---+---+---+---+---+---+
|L|D|E|   Reserved               |
+---+---+---+---+---+---+---+---+

```

L: Loose Bit: This bit indicates the type of next "Topological PDE-ID" in the path description. If set, the next PDE is Loose. If this flag is unset, the next Topological PDE is Strict Type.

D: Destination Bit: By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is Topological and this PDE represents the PDE-ID corresponding to the PPR-Prefix Section 2.2.

E: Egress Bit. By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is 2 i.e., Non-Topological and the service needs to be applied on the egress side of the topological PDE preceding this PDE.

Reserved: Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

- o PPR-PDE Sub-TLVs have 2 octet type, 2 octet length and value field is defined per type.
- o PPR-PDE Sub-TLV: Type 4 (IANA TBD), Length Total length of value field in bytes, Value: The Router ID of the neighbor for which the LAN interface is advertised. This Sub-TLV MUST NOT be present, if the PPR-PDE Type is not equal to 1 i.e., Topological PDE and PDE-ID Type 6.

2.5. PPR-Attributes Sub-TLV

PPR-Attribute Sub-TLVs describe the attributes of the path. The following Sub-TLVs draw from a new registry for Sub-TLV numbers; this registry is to be created by IANA, and administered using the first come first serve process:

- o Type 1 (Suggested Value, IANA TBD): PPR-Metric Sub-TLV. Length 4 bytes, and Value is metric of this path represented through the PPR-ID. Different nodes can advertise the same PPR-ID for the same Prefix with a different set of PPR-PDE Sub-TLVs and the receiving node MUST consider the lowest metric value.

3. OSPFv3 PPR TLV

The OSPFv3 PPR TLV is a top level TLV of the following LSAs defined in [I-D.ietf-ospf-ospfv3-lsa-extend].

E-Intra-Area-Prefix-LSA

E-Inter-Area-Prefix-LSA

E-AS-External-LSA

E-Type-7-LSA

Multiple OSPFv3 PPR TLVs MAY be advertised in each LSA mentioned above. The OSPFv3 PPR TLV has the following format:

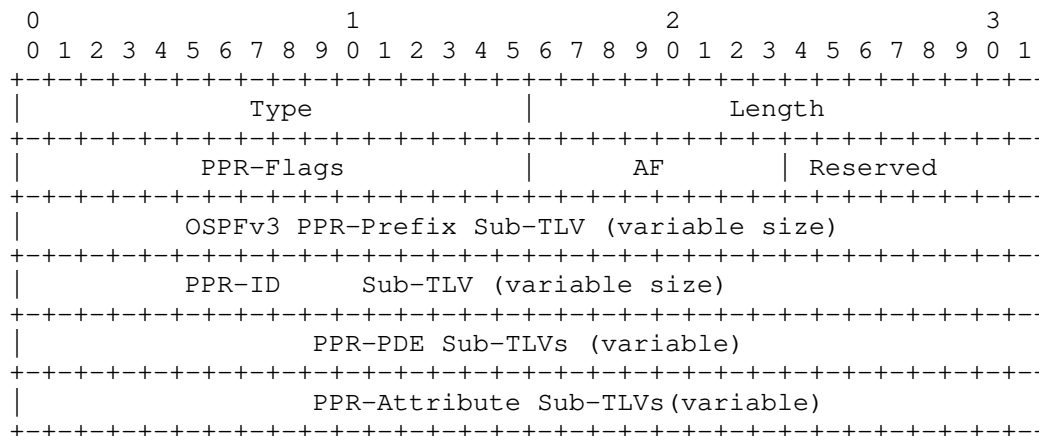


Figure 5: OSPFv3 PPR TLV Format

- o Type - TBD (IANA) from OSPF Extended Prefix Opaque LSA registry.
- o Length - Total length of the value field in bytes (variable).
- o PPR-Flags - 2 Octet flags for this TLV are described below.
- o AF: Address family for the prefix.

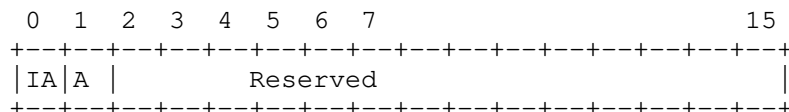
AF: 0 - IPv4 unicast

AF: 1 - IPv6 unicast

- o Reserved - 1 Octet reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

Flags: 2 octet field. The following flags are defined:

OSPFv3 PPR TLV Flags Format



IA-Flag: Inter-Area flag. If set, advertisement is of inter-area type. An ABR that is advertising the OSPF PPR TLV between areas MUST set this bit.

[I-D.ietf-ospf-ospfv3-segment-routing-extensions]

A: The originator of the PPR TLV MUST set the A bit in order to signal that the prefixes and PPR-IDs advertised in the PPR TLV are directly connected to the originators. If this bit is not set, this allows any other node in the network advertise this TLV on behalf of the originating node of the "OSPF Prefix". If PPR TLV is propagated to other areas the A-flag MUST be cleared. In case if the originating node of the prefix has to be disambiguated for any reason including, if it is a Multi Homed Prefix (MHP) or propagated to a different OSPF area, then PPR-Attribute Sub-TLV Source Router ID SHOULD be included.

Reserved - reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

PPR path description for each OSPF area is computed and given to one of the nodes in that area for dissemination. Similarly path information when crossing the area boundaries MUST be relevant to the destination area. If there is no path information available for the destination area, PPR TLV MUST NOT be leaked regardless of the IA bit status.

3.1. OSPFv3 PPR-Prefix Sub-TLV

The structure of OSPFv3 PPR-Prefix, for which path description is attached to is as follows:

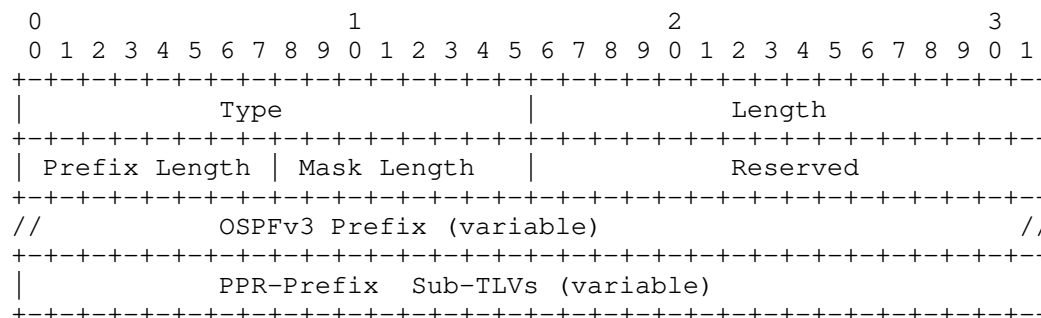


Figure 6: OSPFv3 PPR-Prefix Sub-TLV Format

- o Type - 1 (suggested value, IANA TBD) from OSPFv3 PPR TLV Section 3 Sub-TLV registry.
- o Length - Total length of the value field in bytes (variable).
- o Prefix Len - contains the length of the prefix in bits. Only the most significant octets of the Prefix are encoded.
- o Mask Length - The length of the prefix in bits. Only the most significant octets of the Prefix are encoded.
- o OSPFv3 Prefix - represents the OSPFv3 prefix at the tail-end of the advertised PPR. For the address family IPv4 unicast, the prefix itself is encoded as a 32-bit value. The default route is represented by a prefix of length 0. For the address family (AF in OSPFv3 PPR TLV) in IPv6 unicast, the prefix, encoded as an even multiple of 32-bit words, padded with zeroed bits as necessary. This encoding consumes $((\text{PrefixLength} + 31) / 32)$ 32-bit words.
- o PPR-Prefix Sub-TLVs have 2 octet type, 2 octet length and value field is defined per type.

3.2. OSPFv3 PPR-ID Sub-TLVs

This represents the actual data plane identifier in the packet and could be of any data plane as defined in PPR-ID-type field. Both OSPF Prefix and PPR-ID MUST belong to a same node in the network.

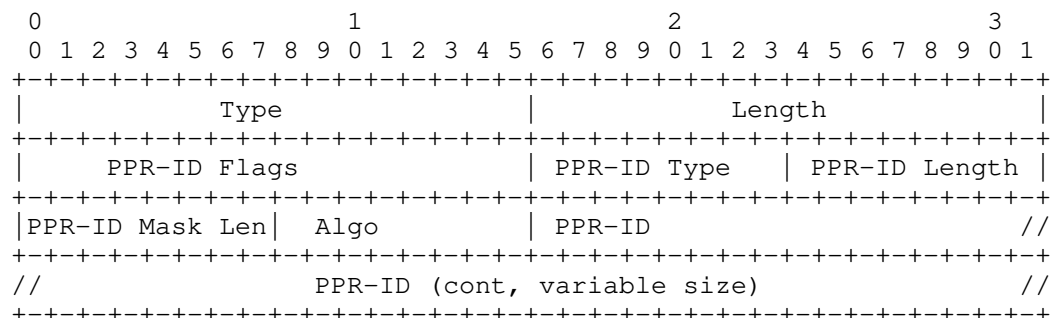


Figure 7: OSPFv3 PPR-ID Sub-TLV Format

- o Type - 2 (suggested value, IANA TBD) from OSPFv3 PPR TLV Section 3 Sub-TLV registry.
- o Length - Total length of the value field in bytes (variable).
- o PPR-ID Type - Data plane type of PPR-ID. This is a new registry (TBD IANA) for this Sub-TLV and the defined types are as follows:

Type: 1 SR-MPLS SID/Label

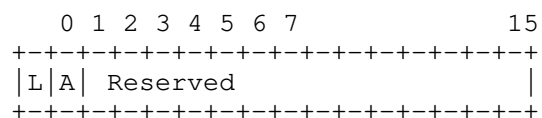
Type: 2 Native IPv4 Address/Prefix

Type: 3 Native IPv6 Address/Prefix

Type: 4 IPv6 SID in SRv6 with SRH

- o PPR-ID Flags - 2 Octet field for PPR-ID flags:

PPR-ID Flags Format



Reserved - Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

- o PPR-ID Length - Length of the PPR-ID field in octets and this depends on the PPR-ID type. See PPR-ID below for the length of this field and other considerations.

- o PPR-ID Mask Len - It is applicable for only for PPR-ID Type 2, 3 and 4. For Type 1 this value MUST be set to zero. It contains the length of the PPR-ID Prefix in bits. Only the most significant octets of the Prefix are encoded. This is needed, if PPR-ID followed is an IPv4/IPv6 Prefix instead of 4/16 octet Address respectively.
- o Algo - 1 octet value represents the SPF algorithm. Algorithm registry is as defined in [I-D.ietf-ospf-ospfv3-segment-routing-extensions].
- o PPR-ID - This is the Preferred Path forwarding identifier that would be on the data packet. The value of this field is variable and it depends on the PPR-ID Type - for Type 1, this is encoded as SR-MPLS SID/Label. For Type 2 this is encoded as 4 byte IPv4 address. For Type 3 this is encoded as 16 byte IPv6 address. For Type 2 and Type 3 encoding is similar to OSPF Prefix as specified in Section 2.2. For Type 4, this is encoded as 16 byte IPv6 SID.

3.3. OSPFv3 PPR-PDE Sub-TLV

This is a new Sub-TLV type in PPR TLV Section 3 and is called as PPR Path Description Element (PDE). PPR-PDEs are used to describe the path in the form of set of contiguous and ordered Sub-TLVs, where first Sub-TLV represents (the top of the stack in MPLS data plane or) first node/segment of the path. These set of ordered Sub-TLVs can have both topological SIDs and non-topological SIDs (e.g., service segments).

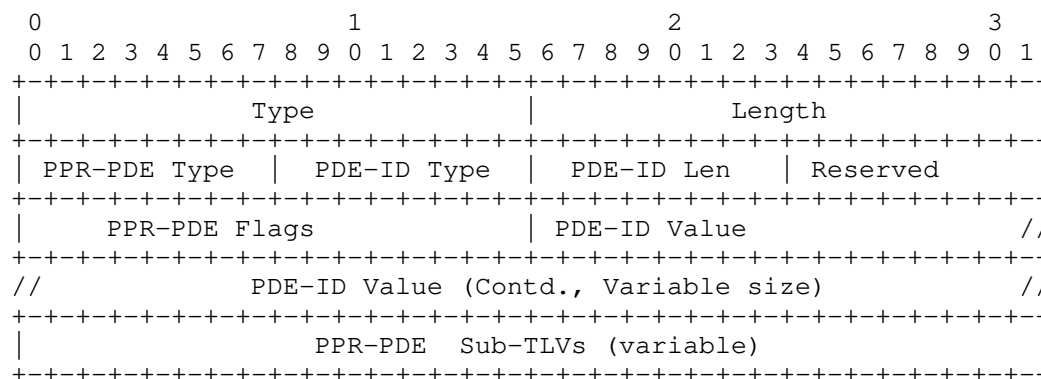


Figure 8: OSPFv3 PPR-PDE Sub-TLV Format

- o Type - 3 (suggested value, IANA TBD) from OSPFv3 PPR TLV Section 3 Sub-TLV registry.

- o Length - Total length of the value field in bytes (variable).
- o PPR-PDE Type - This is a new registry (TBD IANA) for this Sub-TLV and the defined types are as follows:
 - Type: 1 Topological
 - Type: 2 Non-Topological
- o PDE-ID Type - 1 Octet PDE-forwarding IDentifier Type. This is a new registry (TBD IANA) for this Sub-TLV and the defined types and corresponding PDE-ID Len, PDE-ID Value are as follows:

Type 0: This value MUST be set only when PPR-PDE Type is Non-Topological. PDE-ID Len specified in bytes and encoded in NBO in PDE-ID Value field which can represent a service/function. This information is provisioned on the immediate topological PDE preceding to this PDE based on the 'E' bit.

Type 1: SID/Label Sub-TLV as defined in [I-D.ietf-ospf-segment-routing-extensions]. PDE-ID Len and PDE-ID Value fields are per Section 2.1 of the referenced document.

Type 2: SR-MPLS Prefix SID. PDE-ID Len and PDE-ID Value are same as Type 1.

Type 3: SR-MPLS Adjacency SID. PDE-ID Len and PDE-ID Value are same as Type 1.

Type 4: IPv4 Node Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2.

Type 5: IPv4 P2P interface Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2.

Type 6: IPv4 LAN interface Address. PDE-ID Len is 4 bytes and PDE-ID Value is 4 bytes IPv4 address encoded similar to IPv4 Prefix described in Section 2.2. This type MUST have OSPF Neighbor ID Sub-TLV in the PDE.

Type 7: IPv6 Node Address. PDE-ID Len is 16 bytes and PDE-ID Value is 16 bytes IPv6 address encoded similar to IPv6 Prefix described in Section 2.2.

Type 8: IPv6 P2P interface Address. PDE-ID Len is 16 bytes and PDE-ID Value is 16 bytes IPv6 address encoded similar to IPv6 Prefix described in Section 2.2.

Type 9: IPv6 LAN interface Address. PDE-ID Len is 16 bytes and PDE-ID Value is 16 bytes IPv6 address encoded similar to IPv6 Prefix described in Section 2.2. This type MUST have OSPF Neighbor ID Sub-TLV in the PDE.

Type 10: SRv6 Node SID as defined in [I-D.li-ospf-ospfv3-srv6-extensions]. PDE-ID Len and PDE-ID Value are as defined in SRv6 SID.

Type 11: SRv6 Adjacency-SID. PDE-ID Len and PDE-ID Value are as defined in Type 6.

- o PDE-ID Len - 1 Octet. Length of PDE-ID field.
- o Reserved - 1 Octet reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.
- o PPR-PDE Flags - 2 Octet flags for this TLV are described below:

PPR-PDE Flags Format

```

    0 1 2 3 4 5 6 7...          15
+---+---+---+---+---+---+---+
|L|D|E|   Reserved               |
+---+---+---+---+---+---+---+
```

L: Loose Bit. This bit indicates the type of next "Topological PDE-ID" in the path description and overrides the L bit in Section 3.2. If set, the next PDE is Loose. If this flag is unset, the next Topological PDE is Strict Type.

D: Destination Bit. By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is Topological and this PDE represents the PDE-ID corresponding to the PPR-Prefix Section 3.1.

E: Egress Bit. By default this bit MUST be unset. This bit MUST be set only for PPR-PDE Type is 2 i.e., Non-Topological and the service needs to be applied on the egress side of the topological PDE preceding this PDE.

Reserved - Reserved bits for future use. Reserved bits MUST be reset on transmission and ignored on receive.

- o PPR-PDE Sub-TLVs have 2 octet type, 2 octet length and value field is defined per type.
- o PPR-PDE Sub-TLV: Type 4 (IANA TBD), Length Total length of value field in bytes, Value: The Router ID of the neighbor for which the LAN interface is advertised. This Sub-TLV MUST NOT be present, if the PPR-PDE Type is not equal to 1 i.e., Topological PDE and PDE-ID Type 6/9.

3.4. OSPFv3 PPR-Attributes Sub-TLV

PPR-Attribute Sub-TLVs describe the attributes of the path. The following Sub-TLVs draw from a new registry for Sub-TLV numbers; this registry is to be created by IANA, and administered using the first come first serve process:

- o Type 1 (suggested value, IANA TBD): PPR-Metric Sub-TLV. Length 4 bytes, and Value is metric of this path represented through the PPR-ID. Different nodes can advertise the same PPR-ID for the same Prefix with a different set of PPR-PDE Sub-TLVs and the receiving node MUST consider the lowest metric value.

4. Other Considerations

Please refer to [I-D.chunduri-isis-preferred-path-routing] section 4, 5, 6 and 7.

5. Acknowledgements

Thanks to Richard Li, Alex Clemm, Padma Pillay-Esnault, Toerless Eckert, Kiran Makhijani and Lin Han for initial discussions on this topic. Thanks to Kevin Smith and Stephen Johnson for various deployment scenarios applicability from ETSI WGs perspective. Authors also acknowledge Alexander Vainshtein for detailed discussions and suggestions on this topic.

Earlier versions of draft-ietf-ospf-segment-routing-extensions have a mechanism to advertise EROs through Binding SID.

6. IANA Considerations

This document requests the following new TLV in IANA OSPFv2 and OSPFv3 TLV code-point registry as specified in Section 2 Section 3 respectively .

TLV #	Name
-----	-----
TBD	PPR TLV

This document also requests IANA to create new registries for PPR TLV Flags field, PPR Flags, and PPR Sub-TLVs in PPR TLV as described in Section 2 and Section 3.

7. Security Considerations

Existing security extensions as described in [RFC2328] and [RFC7684] apply to the extensions specified in this document. While OSPF is under a single administrative domain, there can be deployments where potential attackers have access to one or more networks in the OSPF routing domain. In these deployments, stronger authentication mechanisms such as those specified in [RFC7474] SHOULD be used.

Advertisement of the additional information defined in this document introduces no new security concerns in OSPF protocol. However as this extension is related to SR-MPLS and SRH data planes as defined in [I-D.ietf-spring-segment-routing], those particular data plane security considerations does apply here.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

8.2. Informative References

- [I-D.chunduri-lsr-isis-preferred-path-routing] Chunduri, U., Li, R., White, R., Tantsura, J., Contreras, L., and Y. Qu, "Preferred Path Routing (PPR) in IS-IS", draft-chunduri-lsr-isis-preferred-path-routing-04 (work in progress), July 2019.

- [I-D.ietf-6man-segment-routing-header]
Filsfils, C., Dukes, D., Previdi, S., Leddy, J., Matsushima, S., and D. Voyer, "IPv6 Segment Routing Header (SRH)", draft-ietf-6man-segment-routing-header-26 (work in progress), October 2019.
- [I-D.ietf-ospf-ospfv3-lsa-extend]
Lindem, A., Roy, A., Goethals, D., Vallem, V., and F. Baker, "OSPFv3 LSA Extendibility", draft-ietf-ospf-ospfv3-lsa-extend-23 (work in progress), January 2018.
- [I-D.ietf-ospf-ospfv3-segment-routing-extensions]
Psenak, P. and S. Previdi, "OSPFv3 Extensions for Segment Routing", draft-ietf-ospf-ospfv3-segment-routing-extensions-23 (work in progress), January 2019.
- [I-D.ietf-ospf-segment-routing-extensions]
Psenak, P., Previdi, S., Filsfils, C., Gredler, H., Shakir, R., Henderickx, W., and J. Tantsura, "OSPF Extensions for Segment Routing", draft-ietf-ospf-segment-routing-extensions-27 (work in progress), December 2018.
- [I-D.ietf-spring-segment-routing]
Filsfils, C., Previdi, S., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", draft-ietf-spring-segment-routing-15 (work in progress), January 2018.
- [I-D.li-ospf-ospfv3-srv6-extensions]
Li, Z., Hu, Z., Cheng, D., Talaulikar, K., and P. Psenak, "OSPFv3 Extensions for SRv6", draft-li-ospf-ospfv3-srv6-extensions-07 (work in progress), November 2019.
- [RFC4915] Psenak, P., Mirtorabi, S., Roy, A., Nguyen, L., and P. Pillay-Esnault, "Multi-Topology (MT) Routing in OSPF", RFC 4915, DOI 10.17487/RFC4915, June 2007, <<https://www.rfc-editor.org/info/rfc4915>>.
- [RFC7474] Bhatia, M., Hartman, S., Zhang, D., and A. Lindem, Ed., "Security Extension for OSPFv2 When Using Manual Key Management", RFC 7474, DOI 10.17487/RFC7474, April 2015, <<https://www.rfc-editor.org/info/rfc7474>>.
- [RFC7684] Psenak, P., Gredler, H., Shakir, R., Henderickx, W., Tantsura, J., and A. Lindem, "OSPFv2 Prefix/Link Attribute Advertisement", RFC 7684, DOI 10.17487/RFC7684, November 2015, <<https://www.rfc-editor.org/info/rfc7684>>.

Authors' Addresses

Uma Chunduri
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
USA

Email: umac.ietf@gmail.com

Yingzhen Qu
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
USA

Email: yingzhen.qu@futurewei.com

Russ White
Juniper Networks
Oak Island, NC 28465
USA

Email: russ@riw.us

Jeff Tantsura
Apstra Inc.
333 Middlefield Road
Menlo Park, CA 94025
USA

Email: jefftant.ietf@gmail.com

Luis M. Contreras
Telefonica
Sur-3 building, 3rd floor
Madrid 28050
Spain

Email: luismiguel.contrerasmurillo@telefonica.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: 3 August 2022

J. Dong
Z. Hu
Z. Li
Huawei Technologies
X. Tang
R. Pang
China Unicom
L. JooHeon
LG U+
S. Bryant
Futurewei Technologies
30 January 2022

IGP Extensions for Scalable Segment Routing based Enhanced VPN
draft-dong-lsr-sr-enhanced-vpn-07

Abstract

Enhanced VPN (VPN+) aims to provide enhanced VPN services to support some application's needs of enhanced isolation and stringent performance requirements. VPN+ requires integration between the overlay VPN connectivity and the characteristics provided by the underlay network. A Virtual Transport Network (VTN) is a virtual underlay network which has a customized network topology and a set of network resources allocated from the physical network. A VTN could be used to support one or a group of VPN+ services.

This document specifies the IGP mechanisms with necessary extensions to advertise the associated topology and resource attributes for scalable Segment Routing (SR) based VTNs. Each VTN can have a customized topology and a set of network resources allocated from the physical network. Multiple VTNs may shared the same topology, and multiple VTNs may share the same set of network resources on some network segments. A group of resource-aware SIDs are allocated for each VTN. This allows flexible combination of the network topology and network resource attributes to build a relatively large number of VTNs with a small number of logical topologies. The proposed mechanism is applicable to both Segment Routing with MPLS data plane (SR-MPLS) and segment routing with IPv6 data plane (SRv6). This document also describes the mechanisms of using dedicated VTN-ID in the data plane instead of the per-VTN resource-aware SIDs to further reduce the control plane overhead.

Requirements Language

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 [RFC2119].

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 <https://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 3 August 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	3
2. VTN Definition Advertisement	4
3. Advertisement of VTN Topology Attribute	6
3.1. MTR based Topology Advertisement	6
3.2. Flex-Algo based Topology Advertisement	7
4. Advertisement of VTN Resource Attribute	8
4.1. Option 1: L2 Bundle based Approach	8
4.2. Option 2: Per-VTN Link TE Attributes	10
5. Advertisement of VTN specific Data Plane Identifiers	12
5.1. Advertisement of VTN-specific SR-MPLS SIDs	12
5.2. Advertisement of VTN-specific SRv6 Locators and SIDs	14

5.2.1. VTN-specific SRv6 Locators and End SIDs	14
5.2.2. VTN-specific SRv6 End.X SIDs	17
5.3. Advertisement of Dedicated Data Plane VTN IDs	17
6. Security Considerations	18
7. IANA Considerations	18
8. Contributors	19
9. Acknowledgments	19
10. References	19
10.1. Normative References	19
10.2. Informative References	21
Authors' Addresses	21

1. Introduction

Enhanced VPN (VPN+) is an enhancement to VPN services to support the needs of new applications, particularly the applications that are associated with 5G services. These applications require enhanced isolation and have more stringent performance requirements than that can be provided with traditional overlay VPNs. These properties require integration between the underlay and the overlay networks. [I-D.ietf-teas-enhanced-vpn] specifies the framework of enhanced VPN and describes the candidate component technologies in different network planes and layers. An enhanced VPN can be used for 5G network slicing, and will also be of use in more generic scenarios.

To meet the requirement of different enhanced VPN services, a number of virtual underlay networks need to be created, each with a customized network topology and a set of network resources allocated from the physical network to meet the requirement of one or a group of VPN+ services. Such a virtual underlay network is called Virtual Transport Network (VTN) in [I-D.ietf-teas-enhanced-vpn].

[I-D.ietf-spring-resource-aware-segments] introduces resource-aware segments by associating existing type of SIDs with network resource attributes (e.g. bandwidth, processing or storage resources). These resource-aware SIDs retain their original functionality, with the additional semantics of identifying the set of network resources available for the packet processing action. [I-D.ietf-spring-sr-for-enhanced-vpn] describes the use of resource-aware segments to build SR based VTNs. To allow the network controller and network nodes to perform VTN-specific explicit path computation and/or shortest path computation, the group of resource-aware SIDs allocated by network nodes to each VTN and the associated topology and resource attributes need to be distributed using the control plane.

[I-D.dong-teas-nrp-scalability] analyzes the scalability requirements and the control plane and data plane scalability considerations of enhanced VPN, more specifically, the scalability of the VTNs. In order to support a relatively large number of VTNs in the network, one proposed approach is to separate the topology and resource attributes of the VTN in control plane, so that the advertisement and processing of each type of attribute could be decoupled. Multiple VTNs may share the same topology, and multiple VTNs may share the same set of network resources on some network segments, while the difference in either the topology or resource attributes makes them different VTNs. This allows flexible combination of network topology and network resource attributes to build a large number of VTNs with a relatively small number of logical topologies.

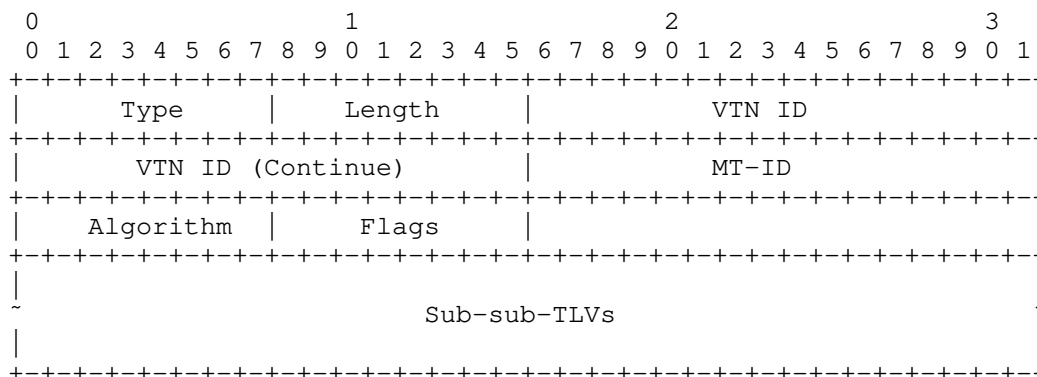
This document specifies the IGP control plane mechanisms with necessary extensions for scalable SR based VTNs. The proposed mechanism is applicable to both segment routing with MPLS data plane (SR-MPLS) and segment routing with IPv6 data plane (SRv6). This document also describes the mechanisms of using dedicated VTN-ID in the data plane instead of the per-VTN resource-aware SIDs to further reduce the control plane overhead.

In general this approach applies to both IS-IS and OSPF, while the specific protocol extensions and encodings are different. In the current version of this document, the required IS-IS extensions are described. The required OSPF extensions will be described in a future version or in a separate document.

2. VTN Definition Advertisement

According to [I-D.ietf-teas-enhanced-vpn], a VTN is associated with a customized network topology and a set of dedicated or shared network resources. Thus a VTN can be defined as the combination of a set of network attributes, which include the topology attribute and other attributes, such as the network resources. IS-IS Virtual Transport Network Definition (VTND) sub-TLV is used to advertise the definition of a VTN. It is a sub-TLV of the IS-IS Router-Capability TLV 242 as defined in [RFC7981].

The format of IS-IS VTND sub-TLV is as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the included sub-TLVs.
- * VTN ID: A global significant 32-bit identifier which is used to identify a VTN.
- * MT-ID: 16-bit field which indicates the multi-topology identifier as defined in [RFC5120]. The first 4-bit are set to zero.
- * Algorithm: 8-bit identifier which indicates the algorithm which applies to this VTN. It can be either a normal algorithm [RFC8402] or a Flexible Algorithm [I-D.ietf-lsr-flex-algo].
- * Flags: 8-bit flags. Currently all the flags are reserved for future use. They SHOULD be set to zero on transmission and MUST be ignored on receipt.
- * Sub-sub-TLVs: optional sub-sub-TLVs to specify the additional attributes of a VTN. Currently no sub-sub-TLV is defined in this document.

The VTND Sub-TLV MAY be advertised in an LSP of any number. A node MUST NOT advertise more than one VTND Sub-TLV for a given VTN ID.

3. Advertisement of VTN Topology Attribute

This section describes the mechanisms used to advertise the topology attribute associated with SR based VTNs. Basically the topology of a VTN can be determined by the MT-ID and/or the algorithm ID included in the VTN definition. In practice, it could be described using two optional approaches.

The first approach is to use Multi-Topology Routing (MTR) [RFC4915] [RFC5120] with the segment routing extensions to advertise the topology associated with the SR based VTNs. Different algorithms MAY be used to further specify the computation algorithm or the metric type used for path computation within the topology. Multiple VTNs can be associated with the same <topology, algorithm>, and the IGP computation with the <topology, algorithm> tuple can be shared by these VTNs.

The second approach is to use Flex- Algo [I-D.ietf-lsr-flex-algo] to describe the topological constraints of SR based VTNs on a shared network topology (e.g. the default topology). Multiple VTNs can be associated with the same Flex- Algo, and the IGP computation with this Flex- Algo can be shared by these VTNs.

3.1. MTR based Topology Advertisement

Multi-Topology Routing (MTR) has been defined in [RFC4915] and [RFC5120] to create different network topologies in one network. It also has the capability of specifying customized attributes for each topology. The traditional use cases of multi-topology are to maintain separate topologies for unicast and multicast services, or to create different topologies for IPv4 and IPv6 in a network. There are some limitations when MTR is used with native IP forwarding, the considerations about MT based IP forwarding are described in [RFC5120].

MTR can be used with SR-MPLS data plane. [RFC8667] specifies the IS-IS extensions to support SR-MPLS data plane, in which the Prefix-SID sub-TLVs can be carried in IS-IS TLV 235 (MT IP Reachability) and TLV 237 (MT IPv6 IP Reachability), and the Adj-SID sub-TLVs can be carried in IS-IS TLV 222 (MT-ISN) and TLV 223 (MT IS Neighbor Attribute).

MTR can also be used with SRv6 data plane.

[I-D.ietf-lsr-isis-srv6-extensions] specifies the IS-IS extensions to support SRv6 data plane, in which the MT-ID is carried in the SRv6 Locator TLV. The SRv6 End SIDs are carried as sub-TLVs in the SRv6 Locator TLV, and inherit the topology/algorithm from the parent locator. The SRv6 End.X SIDs are carried as sub-TLVs in the IS-IS TLV 222 (MT-ISN) and TLV 223 (MT IS Neighbor Attribute), and inherit the topology/algorithm from the parent locator.

These IGP extensions for SR-MPLS and SRv6 can be used to advertise and build the topology for a group of SR based VTNs.

An algorithm ID MAY be used to further specify the computation algorithm or the metric type used for path computation within the topology.

3.2. Flex-Algo based Topology Advertisement

[I-D.ietf-lsr-flex-algo] specifies the mechanisms to provide distributed computation of constraint-based paths, and how the SR-MPLS prefix-SIDs and SRv6 locators can be used to steer packets along the constraint-based paths.

The Flex-Algo Definition (FAD) can be used to describe the topological constraints for path computation on a network topology. According to the network nodes' participation of a Flex-Algo, and the rules of including or excluding specific Administrative Groups (colors) and the Shared Risk Link Groups (SRLGs), the topology of a VTN can be determined using the associated Flex-Algo on a particular topology (e.g. the default topology).

With the mechanisms defined in[RFC8667] [I-D.ietf-lsr-flex-algo], prefix-SID advertisement can be associated with a <topology, algorithm> tuple, in which the algorithm can be a Flex-Algo. This allows network nodes to use the prefix-SID to steer traffic along distributed computed paths according to the identified Flex-Algo in the topology.

[I-D.ietf-lsr-isis-srv6-extensions] specifies the IS-IS extensions to support SRv6 data plane, in which the SRv6 locators advertisement can be associated with a specific topology and a specific algorithm, which can be a Flex-Algo. With the mechanism defined in [I-D.ietf-lsr-flex-algo], The SRv6 locator can be used to steer traffic along distributed computed paths according to the identified Flex-Algo in the topology. In addition, topology/algorithm specific SRv6 End SID and End.X SID can be used to enforce traffic over the LFA computed backup path.

Multiple Flex-Algos MAY be defined to describe the topological constraints on a shared network topology (e.g. the default topology).

4. Advertisement of VTN Resource Attribute

This section specifies the mechanisms to advertise the network resource attributes associated with the VTNs. The mechanism of advertising the link resources and attributes associated with VTNs is described. The mechanism of advertising node resources and attributes associated with VTNs are for further study. Two optional approaches are described in the following sub-sections: the first option is the L2 Bundle [RFC8668] based approach, the second option is to extend IGP to advertise per-VTN link TE attributes.

4.1. Option 1: L2 Bundle based Approach

On a Layer-3 interface, each VTN can be allocated with a subset of link resources (e.g. bandwidth). A subset of link resources may be dedicated to a VTN, or may be shared by a group of VTNs. Each subset of link resource can be represented as a virtual layer-2 member link under the Layer-3 interface, and the Layer-3 interface is considered as a virtual Layer-2 bundle. The Layer-3 interface may also be a physical Layer 2 link bundle, in this case a subset of link resources allocated to a VTN may be provided by one of the physical Layer-2 member links.

[RFC8668] describes the IS-IS extensions to advertise the link attributes of the Layer 2 member links which comprise a Layer 3 interface. Such mechanism can be extended to advertise the attributes of each physical or virtual member links, and its associated VTNs.

A new flag "E" (Exclusive) is defined in the flag field of the Parent L3 Neighbor Descriptor in the L2 Bundle Member Attributes TLV (25).

```

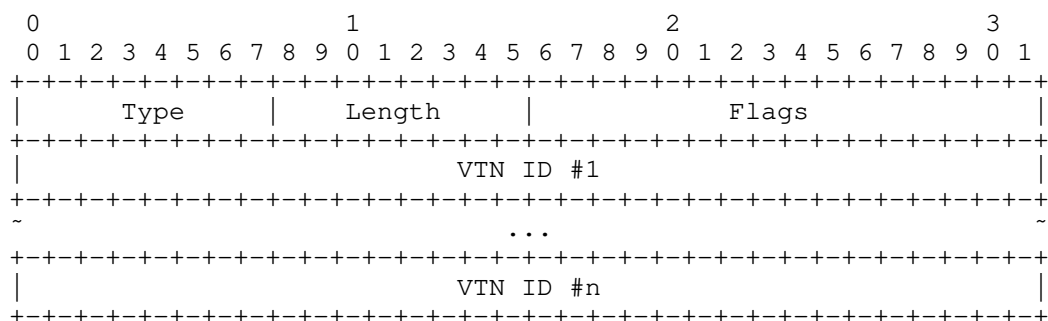
0 1 2 3 4 5 6 7
+--+--+--+--+--+--+
|P|E|          |
+--+--+--+--+--+--+

```

E flag: When the E flag is set, it indicates each member link under the Parent L3 link are used exclusively for one or a specific group of VTNs, and load sharing among the member links is not allowed. When the E flag is clear, it indicates load balancing and sharing among the member links are allowed.

A new VTN-IDs sub-TLV is carried under the L2 Bundle Attribute Descriptors to describe the mapping relationship between the VTNs and the virtual or physical member links. As one or more VTNs may use the same set of link resource on a specific network segment, these VTN IDs will be advertised under the same virtual or physical member link.

The format of the VTN-IDs Sub-TLV is as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the number of VTN IDs included.
- * Flags: 16 bit flags. All the bits are reserved for future use, which SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * VTN IDs: One or more 32-bit identifier to identify the VTNs this member link belongs to.

Each physical or virtual member link MAY be associated with a different group of VTNs. Thus each L2 Bundle Attribute Descriptor may carry the link local identifier and attributes of only one Layer 2 member link. Multiple L2 Bundle Attribute Descriptors will be used to carry the attributes and the associated VTN-IDs of all the Layer 2 member links.

The TE attributes of each virtual or physical member link, such as the bandwidth attributes and the SR SIDs, can be advertised using the mechanism as defined in [RFC8668].

4.2. Option 2: Per-VTN Link TE Attributes

A Layer-3 interface can participate in multiple VTNs, each of which is allocated with a subset of the forwarding resources of the interface. For each VTN, the associated resources can be described using per-VTN TE attributes. A new VTN-specific TE attribute sub-TLV is defined to advertise the link attributes associated with a VTN. This sub-TLV MAY be advertised as a sub-TLV of the following TLVs:

TLV-22 (Extended IS reachability) [RFC5305]

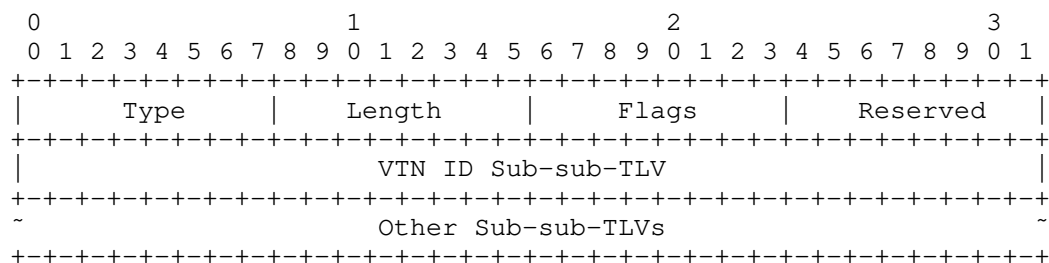
TLV-23 (IS Neighbor Attribute) [RFC5311]

TLV-141 (Inter-AS Reachability Information) [RFC5316]

TLV-222 (MT ISN) [RFC5120]

TLV-223 (MT IS Neighbor Attribute) [RFC5311]

The format of the sub-TLV is shown as below:

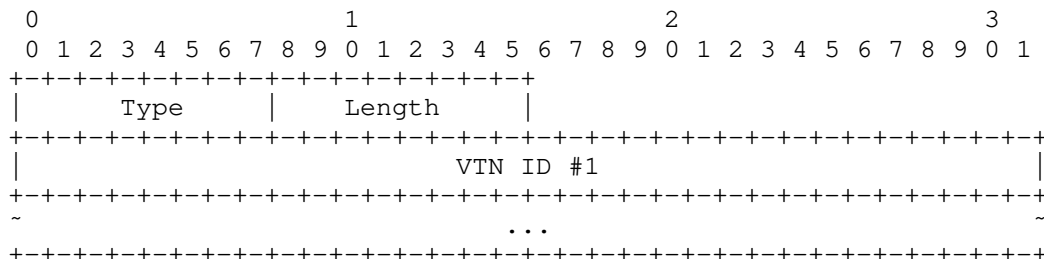


Where:

- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the length of the Sub-sub-TLVs field.
- * Flags: 8-bit flags. All the 8 bits are reserved for future use, which SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * Reserved: 8-bit field reserved for future use, SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * VTN ID Sub-sub-TLV: contains one or more VTN IDs which is associated with the same group of TE attributes.

- * Other Sub-sub-TLVs: the TLVs which carry the TE attributes associated with the VTNs.

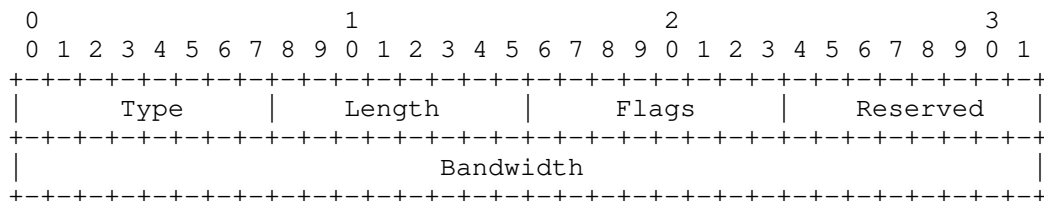
The format of the VTN ID sub-sub-TLV is shown as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-sub-TLV. It is the number of the VTN IDs in the TLV multiplied by 4.
- * VTN ID: A global significant 32-bit identifier which is used to identify a VTN.

One sub-sub-TLV "VTN bandwidth sub-sub-TLV" is defined in this document. Its format is shown as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-sub-TLV. It is set to 6.
- * Flags: 8-bit flags. All the 8 bits are reserved for future use, which SHOULD be set to 0 on transmission and MUST be ignored on receipt.

- * **Reserved:** 8-bit field reserved for future use, SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * **Bandwidth:** The bandwidth allocated to the VTN, encoded in 32 bits in IEEE floating point format.

The VTN-specific Bandwidth sub-sub-TLV is optional. This sub-sub-TLV SHOULD appear once at most in each VTN-specific TE attribute sub-TLV.

5. Advertisement of VTN specific Data Plane Identifiers

In order to steer packets to the VTN-specific paths which are computed taking the topology and network resources of the VTN as the constraints, some fields in the data packet needs to be used to infer or identify the VTN the packet belongs to. As multiple VTNs may share the same topology or Flex-Algo, the topology/Flex-Algo specific SR SIDs or Locators cannot be used to distinguish the packets which belong to different VTNs. Some additional data plane identifiers would be needed to identify the VTN a packet belongs to.

This section describes the mechanisms to advertise the VTN identifiers in different data plane encapsulations.

5.1. Advertisement of VTN-specific SR-MPLS SIDs

With SR-MPLS data plane, the VTN identification information can be implicitly carried in the VTN-specific SIDs. Each node SHOULD allocate a unique Prefix-SID for each VTN it participates in. On a Layer-3 interface, if each Layer 2 member link is associated with only one VTN, the adj-SIDs of the L2 member links could also identify the VTNs. If a member link is associated with multiple VTNs, VTN-specific adj-SIDs MAY need to be allocated to help the VTN-specific local protection.

A new VTN-specific prefix-SID sub-TLV is defined to advertise the prefix-SID and its associated VTN. This sub-TLV MAY be advertised as a sub-TLV of the following TLVs:

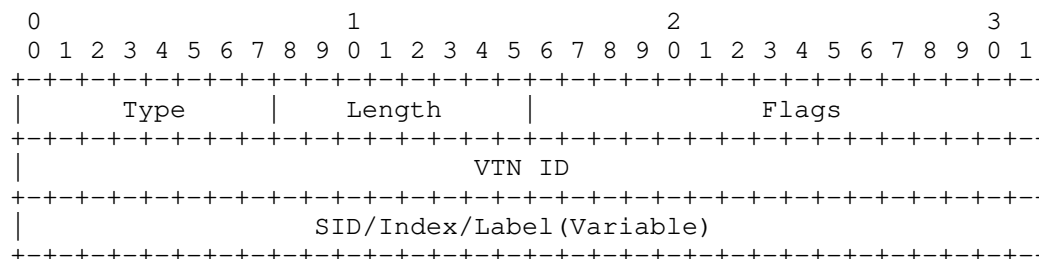
TLV-135 (Extended IPv4 Reachability) defined in [RFC5305].

TLV-235 (MT IP Reachability) defined in [RFC5120].

TLV-236 (IPv6 IP Reachability) defined in [RFC5308].

TLV-237 (MT IPv6 IP Reachability) defined in [RFC5120].

The format of the sub-TLV is shown as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the length of the SID/Index/Label field.
- * Flags: 16-bit flags. The high-order 8 bits are the same as in the Prefix-SID sub-TLV defined in [RFC8667]. The lower-order 8 bits are reserved for future use, which SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * VTN ID: A 32-bit identifier to identify the VTN this prefix-SID associates with.
- * SID/Index/Label: The same as defined in [RFC8667].

One or more of VTN-specific Prefix-SID sub-TLVs MAY be carried in the Multi-topology IP Reachability TLVs (TLV 235 or TLV 237), the MT-ID of the TLV SHOULD be the same as the MT-ID in the definition of these VTNs.

A new VTN-specific Adj-SID sub-TLV is defined to advertise the adj-SID and its associated VTN. This sub-TLV may be advertised as a sub-TLV of the following TLVs:

TLV-22 (Extended IS reachability) [RFC5305]

TLV-23 (IS Neighbor Attribute) [RFC5311]

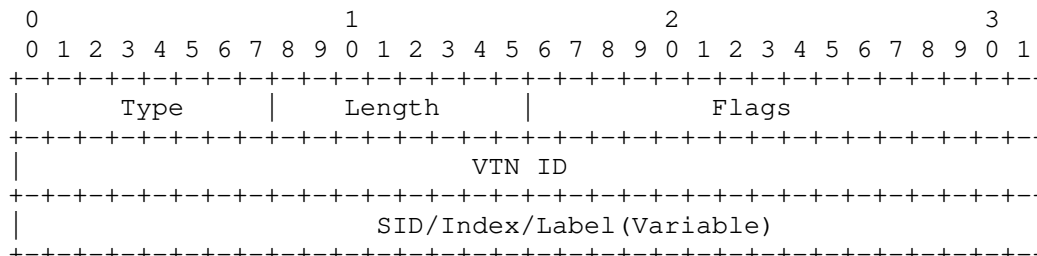
TLV-25 (L2 Bundle Member Attributes) [RFC8668]

TLV-141 (Inter-AS Reachability Information) [RFC5316]

TLV-222 (MT ISN) [RFC5120]

TLV-223 (MT IS Neighbor Attribute) [RFC5311]

The format of the sub-TLV is shown as below:



Where:

- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the length of the SID/Index/Label field.
- * Flags: 16-bit flags. The high-order 8 bits are the same as in the Adj-SID sub-TLV defined in [RFC8667]. The lower-order 8 bits are reserved for future use, which SHOULD be set to 0 on transmission and MUST be ignored on receipt.
- * VTN ID: A 32-bit global identifier to identify the VTN this Adj-SID associates with.
- * SID/Index/Label: The same as defined in [RFC8667].

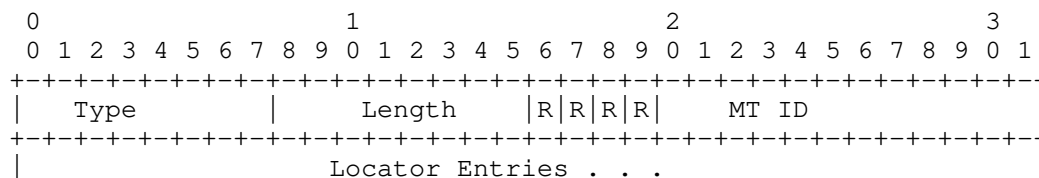
One or more VTN-specific Adj-SID sub-TLV MAY be carried in the Multi-topology ISN or Multi-topology IS Attribute TLVs (TLV 222 or TLV 223), the MT-ID of the TLV SHOULD be the same as the MT-ID in the definition of these VTNs.

5.2. Advertisement of VTN-specific SRv6 Locators and SIDs

5.2.1. VTN-specific SRv6 Locators and End SIDs

With SRv6 data plane, the VTN identification information can be implicitly or explicitly carried in the SRv6 Locator of the corresponding VTN, this is to ensure that all network nodes (including both the end nodes and the transit nodes) can identify the VTN to which a packet belongs to. Network nodes SHOULD allocate VTN-specific Locators for each VTN it participates in. The VTN-specific Locators are used as the covering prefix of VTN-specific SRv6 End SIDs, End.X SIDs and other types of SIDs.

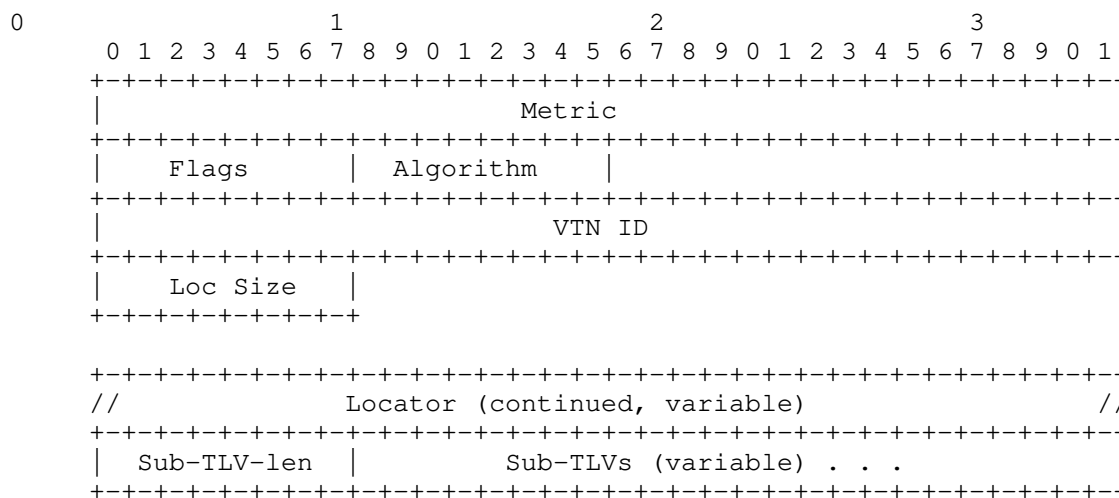
In one possible approach, each VTN-specific Locator is advertised in a separate TLV called "VTN specific SRv6 Locator TLV". Its format is shown as below:



Where:

- * Type: TBD
- * The semantics of the Length field, the R bits and the MT ID field are the same as those defined in [I-D.ietf-lsr-isis-srv6-extensions].

Followed by one or more locator entries of the form:



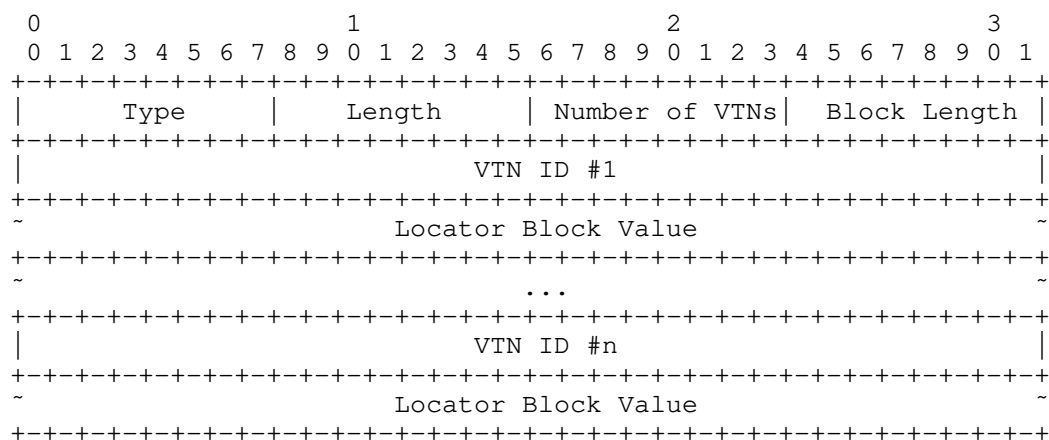
Where:

- * VTN ID: A 32-bit global identifier to identify the VTN this Locator associates with.
- * All the other fields are the same as those defined in [I-D.ietf-lsr-isis-srv6-extensions].

The VTN-specific SRv6 End SIDs are carried in the VTN-specific SRv6 Locator TLV, and inherits the topology, algorithm and VTN from the parent VTN-specific Locator.

In another possible approach, when a group of VTNs share the same topology/algorithm, the topology/algorithm specific Locator is the covering prefix of a group of VTN-specific Locators. Then the advertisement of VTN-specific locators can be optimized to reduce the amount of Locator TLVs advertised in the control plane.

A new VTN locator-block sub-TLV under the SRv6 Locator TLV is defined to advertise a set of sub-blocks which follows the topology/algorithm specific Locator. Each VTN locator-block value is assigned to one of the VTNs which share the same topology/algorithm.



Where:

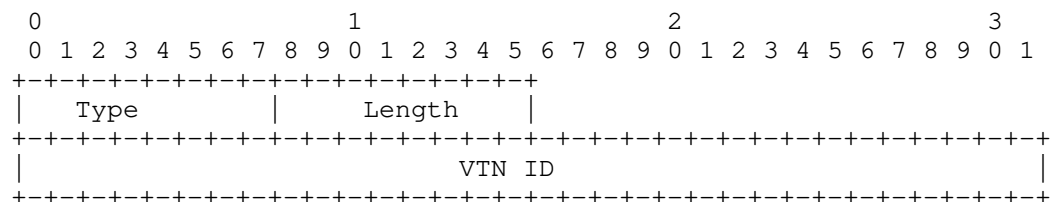
- * Type: TBD
- * Length: The length of the value field of the sub-TLV. It is variable dependent on the number of VTNs and the Block Length.
- * Number of VTNs: The number of VTNs which share the same topology/algorithm specific Locator as the covering prefix.
- * Block Length: The length of the VTN locator-block which follows the length of the topology/algorithm specific Locator.
- * VTN ID: A 32-bit global identifier to identify the VTN the locator-block is associates with.
- * Block Value: The value of the VTN locator-block for each VTN.

With the VTN locator-block sub-TLV, the VTN-specific Locator can be obtained by concatenating the topology/algorithm specific locator and the locator-block value advertised for the VTN.

The VTN-specific SRv6 End SIDs inherit the topology, algorithm and the VTN from the parent VTN-specific Locator.

5.2.2. VTN-specific SRv6 End.X SIDs

The SRv6 End.X SIDs are advertised as sub-TLVs of TLV 22, 23, 25, 141, 222, and 223. In order to distinguish the End.X SIDs which belong to different VTNs, a new "VTN ID sub-sub-TLV" is introduced under the SRv6 End.X SID sub-TLV and SRv6 LAN End.X SID sub-TLV defined in [I-D.ietf-lsr-isis-srv6-extensions]. Its format is shown as below:



Where:

- * Type: TBD.
- * Length: the length of the Value field of the TLV. It is set to 4.
- * VTN ID: A 32-bit global identifier to identify the VTN this End.X SID associates with.

5.3. Advertisement of Dedicated Data Plane VTN IDs

As the number of VTNs increases, with the mechanism described in [I-D.ietf-spring-sr-for-enhanced-vpn], the number of SR SIDs and SRv6 Locators allocated for different VTNs would also increase. In network scenarios where the number of SIDs or Locators becomes a concern, some data plane optimization may be needed to reduce the amount of SR SIDs and Locators allocated. As described in [I-D.dong-teas-nrp-scalability], one approach is to decouple the data plane identifiers used for topology based forwarding and the identifiers used for the VTN-specific processing. Thus a dedicated data plane VTN-ID could be encapsulated in the packet. One possible encapsulation of VTN-ID in IPv6 data plane is proposed in [I-D.dong-6man-enhanced-vpn-vtn-id]. One possible encapsulation of VTN-ID in MPLS data plane is proposed in

[I-D.li-mpls-enhanced-vpn-vtn-id].

In that case, the VTN-ID encapsulated in data plane can have the same value as the VTN-ID in control plane, so that the overhead of advertising the mapping between the control plane VTN-IDs and the corresponding data plane identifiers could be saved.

6. Security Considerations

This document introduces no additional security vulnerabilities to IS-IS.

The mechanism proposed in this document is subject to the same vulnerabilities as any other protocol that relies on IGPs.

7. IANA Considerations

IANA is requested to assign a new code point in the "sub-TLVs for TLV 242 registry".

Type: TBD1

Description: Virtual Transport Network Definition

IANA is requested to assign three new code points in the "sub-TLVs for TLVs 22, 23, 25, 141, 222, and 223 registry".

Type: TBD2

Description: Virtual Transport Network Identifiers

Type: TBD3

Description: VTN-specific TE attribute sub-TLV

Type: TBD4

Description: VTN-specific Adj-SID

IANA is requested to assign two new code points in the "Sub-TLVs for TLVs 27, 135, 235, 236 and 237 registry".

Type: TBD5

Description: VTN-specific Prefix-SID

Type: TBD6

Description: VTN locator-block

IANA is requested to assign a new code point in the "IS-IS TLV Codepoints registry".

Type: TBD7

Description: VTN-specific SRv6 Locator TLV

IANA is requested to assign a new code point in the "sub-sub-TLVs for SRv6 End SID and SRv6 End.X SID registry".

Type: TBD8

Description: VTN ID Sub-sub-TLV

8. Contributors

Hongjie Yang

Email: hongjie.yang@huawei.com

9. Acknowledgments

The authors would like to thank Mach Chen, Dean Cheng and Guoqi Xu for their review and discussion of this document.

10. References

10.1. Normative References

[I-D.ietf-lsr-flex-algo]

Psenak, P., Hegde, S., Filsfils, C., Talaulikar, K., and A. Gulko, "IGP Flexible Algorithm", Work in Progress, Internet-Draft, draft-ietf-lsr-flex-algo-18, 25 October 2021, <<https://www.ietf.org/archive/id/draft-ietf-lsr-flex-algo-18.txt>>.

[I-D.ietf-lsr-isis-srv6-extensions]

Psenak, P., Filsfils, C., Bashandy, A., Decraene, B., and Z. Hu, "IS-IS Extensions to Support Segment Routing over IPv6 Dataplane", Work in Progress, Internet-Draft, draft-ietf-lsr-isis-srv6-extensions-18, 20 October 2021, <<https://www.ietf.org/archive/id/draft-ietf-lsr-isis-srv6-extensions-18.txt>>.

[I-D.ietf-spring-resource-aware-segments]

Dong, J., Bryant, S., Miyasaka, T., Zhu, Y., Qin, F., Li, Z., and F. Clad, "Introducing Resource Awareness to SR Segments", Work in Progress, Internet-Draft, draft-ietf-spring-resource-aware-segments-03, 12 July 2021, <<https://www.ietf.org/archive/id/draft-ietf-spring-resource-aware-segments-03.txt>>.

- [I-D.ietf-spring-sr-for-enhanced-vpn]
Dong, J., Bryant, S., Miyasaka, T., Zhu, Y., Qin, F., Li, Z., and F. Clad, "Segment Routing based Virtual Transport Network (VTN) for Enhanced VPN", Work in Progress, Internet-Draft, draft-ietf-spring-sr-for-enhanced-vpn-01, 12 July 2021, <<https://www.ietf.org/archive/id/draft-ietf-spring-sr-for-enhanced-vpn-01.txt>>.
- [I-D.ietf-teas-enhanced-vpn]
Dong, J., Bryant, S., Li, Z., Miyasaka, T., and Y. Lee, "A Framework for Enhanced Virtual Private Network (VPN+) Services", Work in Progress, Internet-Draft, draft-ietf-teas-enhanced-vpn-09, 25 October 2021, <<https://www.ietf.org/archive/id/draft-ietf-teas-enhanced-vpn-09.txt>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC4915] Psenak, P., Mirtorabi, S., Roy, A., Nguyen, L., and P. Pillay-Esnault, "Multi-Topology (MT) Routing in OSPF", RFC 4915, DOI 10.17487/RFC4915, June 2007, <<https://www.rfc-editor.org/info/rfc4915>>.
- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.
- [RFC8402] Filsfils, C., Ed., Previdi, S., Ed., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", RFC 8402, DOI 10.17487/RFC8402, July 2018, <<https://www.rfc-editor.org/info/rfc8402>>.
- [RFC8667] Previdi, S., Ed., Ginsberg, L., Ed., Filsfils, C., Bashandy, A., Gredler, H., and B. Decraene, "IS-IS Extensions for Segment Routing", RFC 8667, DOI 10.17487/RFC8667, December 2019, <<https://www.rfc-editor.org/info/rfc8667>>.

[RFC8668] Ginsberg, L., Ed., Bashandy, A., Filsfils, C., Nanduri, M., and E. Aries, "Advertising Layer 2 Bundle Member Link Attributes in IS-IS", RFC 8668, DOI 10.17487/RFC8668, December 2019, <<https://www.rfc-editor.org/info/rfc8668>>.

10.2. Informative References

[I-D.dong-6man-enhanced-vpn-vtn-id]
Dong, J., Li, Z., Xie, C., Ma, C., and G. Mishra,
"Carrying Virtual Transport Network (VTN) Identifier in
IPv6 Extension Header", Work in Progress, Internet-Draft,
draft-dong-6man-enhanced-vpn-vtn-id-06, 24 October 2021,
<<https://www.ietf.org/archive/id/draft-dong-6man-enhanced-vpn-vtn-id-06.txt>>.

[I-D.dong-teas-nrp-scalability]
Dong, J., Li, Z., Gong, L., Yang, G., Guichard, J. N.,
Mishra, G., and F. Qin, "Scalability Considerations for
Network Resource Partition", Work in Progress, Internet-
Draft, draft-dong-teas-nrp-scalability-00, 17 December
2021, <<https://www.ietf.org/archive/id/draft-dong-teas-nrp-scalability-00.txt>>.

[I-D.li-mpls-enhanced-vpn-vtn-id]
Li, Z. and J. Dong, "Carrying Virtual Transport Network
Identifier in MPLS Packet", Work in Progress, Internet-
Draft, draft-li-mpls-enhanced-vpn-vtn-id-01, 14 April
2021, <<https://www.ietf.org/archive/id/draft-li-mpls-enhanced-vpn-vtn-id-01.txt>>.

Authors' Addresses

Jie Dong
Huawei Technologies

Email: jie.dong@huawei.com

Zhibo Hu
Huawei Technologies

Email: huzhibo@huawei.com

Zhenbin Li
Huawei Technologies

Email: lizhenbin@huawei.com

Xiongyan Tang
China Unicom

Email: tangxy@chinaunicom.cn

Ran Pang
China Unicom

Email: pangran@chinaunicom.cn

Lee JooHeon
LG U+

Email: playgame@lguplus.co.kr

Stewart Bryant
Futurewei Technologies

Email: stewart.bryant@gmail.com

IS-IS for IP Internets
Internet-Draft
Obsoletes: 5306 (if approved)
Intended status: Standards Track
Expires: December 30, 2018

L. Ginsberg
P. Wells
Cisco Systems, Inc.
June 28, 2018

Restart Signaling for IS-IS
draft-ginsberg-isis-rfc5306bis-01

Abstract

This document describes a mechanism for a restarting router to signal to its neighbors that it is restarting, allowing them to reestablish their adjacencies without cycling through the down state, while still correctly initiating database synchronization.

This document additionally describes a mechanism for a router to signal its neighbors that it is preparing to initiate a restart while maintaining forwarding plane state. This allows the neighbors to maintain their adjacencies until the router has restarted, but also allows the neighbors to bring the adjacencies down in the event of other topology changes.

This document additionally describes a mechanism for a restarting router to determine when it has achieved Link State Protocol Data Unit (LSP) database synchronization with its neighbors and a mechanism to optimize LSP database synchronization, while minimizing transient routing disruption when a router starts.

This document obsoletes RFC 5306.

Requirements Language

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.

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 <https://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 December 30, 2018.

Copyright Notice

Copyright (c) 2018 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 (<https://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. Overview	3
2. Approach	4
2.1. Timers	4
2.2. Restart TLV	5
2.2.1. Use of RR and RA Bits	6
2.2.2. Use of the SA Bit	7
2.2.3. Use of PR and PA Bits	8
2.3. Adjacency (Re)Acquisition	10
2.3.1. Adjacency Reacquisition during Restart	10
2.3.2. Adjacency Acquisition during Start	12
2.3.3. Multiple Levels	14
2.4. Database Synchronization	14
2.4.1. LSP Generation and Flooding and SPF Computation	15
3. State Tables	17
3.1. Running Router	18
3.2. Restarting Router	18
3.3. Starting Router	19
4. IANA Considerations	20
5. Security Considerations	21
6. Manageability Considerations	21
7. Acknowledgements	21

8. Normative References	22
Appendix A. Summary of Changes from RFC 5306	23
Authors' Addresses	23

1. Overview

The Intermediate System to Intermediate System (IS-IS) routing protocol [RFC1195] [ISO10589] is a link state intra-domain routing protocol. Normally, when an IS-IS router is restarted, temporary disruption of routing occurs due to events in both the restarting router and the neighbors of the restarting router.

The router that has been restarted computes its own routes before achieving database synchronization with its neighbors. The results of this computation are likely to be non-convergent with the routes computed by other routers in the area/domain.

Neighbors of the restarting router detect the restart event and cycle their adjacencies with the restarting router through the down state. The cycling of the adjacency state causes the neighbors to regenerate their LSPs describing the adjacency concerned. This in turn causes a temporary disruption of routes passing through the restarting router.

In certain scenarios, the temporary disruption of the routes is highly undesirable. This document describes mechanisms to avoid or minimize the disruption due to both of these causes.

When an adjacency is reinitialized as a result of a neighbor restarting, a router does three things:

1. It causes its own LSP(s) to be regenerated, thus triggering SPF runs throughout the area (or in the case of Level 2, throughout the domain).
2. It sets SRMflags on its own LSP database on the adjacency concerned.
3. In the case of a Point-to-Point link, it transmits a complete set of Complete Sequence Number PDUs (CSNPs), over the adjacency.

In the case of a restarting router process, the first of these is highly undesirable, but the second is essential in order to ensure synchronization of the LSP database.

The third action above minimizes the number of LSPs that must be exchanged and, if made reliable, provides a means of determining when the LSP databases of the neighboring routers have been synchronized. This is desirable whether or not the router is being restarted (so

that the overload bit can be cleared in the router's own LSP, for example).

This document describes a mechanism for a restarting router to signal that it is restarting to its neighbors, and allow them to reestablish their adjacencies without cycling through the down state, while still correctly initiating database synchronization.

This document additionally describes a mechanism for a restarting router to determine when it has achieved LSP database synchronization with its neighbors and a mechanism to optimize LSP database synchronization and minimize transient routing disruption when a router starts.

It is assumed that the three-way handshake [RFC5303] is being used on Point-to-Point circuits.

2. Approach

2.1. Timers

Three additional timers, T1, T2, and T3, are required to support the functionality defined in this document.

An instance of the timer T1 is maintained per interface, and indicates the time after which an unacknowledged (re)start attempt will be repeated. A typical value might be 3 seconds.

An instance of the timer T2 is maintained for each LSP database (LSPDB) present in the system, i.e., for a Level 1/2 system, there will be an instance of the timer T2 for Level 1 and an instance for Level 2. This is the maximum time that the system will wait for LSPDB synchronization. A typical value might be 60 seconds.

A single instance of the timer T3 is maintained for the entire system. It indicates the time after which the router will declare that it has failed to achieve database synchronization (by setting the overload bit in its own LSP). This is initialized to 65535 seconds, but is set to the minimum of the remaining times of received IS-IS Hellos (IIHs) containing a restart TLV with the Restart Acknowledgement (RA) set and an indication that the neighbor has an adjacency in the "UP" state to the restarting router.

NOTE: The timer T3 is only used by a restarting router.

2.2. Restart TLV

A new TLV is defined to be included in IIH PDUs. The presence of this TLV indicates that the sender supports the functionality defined in this document and it carries flags that are used to convey information during a (re)start. All IIHs transmitted by a router that supports this capability MUST include this TLV.

Type 211

Length: Number of octets in the Value field (1 to (3 + ID Length))
Value

	No. of octets
+-----+ Flags +-----+	1
+-----+ Remaining Time +-----+	2
+-----+ Restarting Neighbor ID +-----+	ID Length

Flags (1 octet)

0	1	2	3	4	5	6	7
+---	+---	+---	+---	+---	+---	+---	+---
Reserved	PA	PR	SA	RA	RR		
+---	+---	+---	+---	+---	+---	+---	+---

RR - Restart Request
RA - Restart Acknowledgement
SA - Suppress adjacency advertisement
PR - Restart is planned
PA - Planned restart acknowledgement

(Note: Remaining fields are required when the RA bit is set.)
Remaining Time (2 octets)

Remaining holding time (in seconds)

Restarting Neighbor System ID (ID Length octets)

The System ID of the neighbor to which an RA refers. Note: Implementations based on earlier versions of this document may not include this field in the TLV when the RA is set. In this case, a router that is expecting an RA on a LAN circuit SHOULD assume that the acknowledgement is directed at the local system.

2.2.1. Use of RR and RA Bits

The RR bit is used by a (re)starting router to signal to its neighbors that a (re)start is in progress, that an existing adjacency SHOULD be maintained even under circumstances when the normal operation of the adjacency state machine would require the adjacency to be reinitialized, to request a set of CSNPs, and to request setting of the SRMflags.

The RA bit is sent by the neighbor of a (re)starting router to acknowledge the receipt of a restart TLV with the RR bit set.

When the neighbor of a (re)starting router receives an IIH with the restart TLV having the RR bit set, if there exists on this interface an adjacency in state "UP" with the same System ID, and in the case of a LAN circuit, with the same source LAN address, then, irrespective of the other contents of the "Intermediate System Neighbors" option (LAN circuits) or the "Point-to-Point Three-Way Adjacency" option (Point-to-Point circuits):

- a. the state of the adjacency is not changed. If this is the first IIH with the RR bit set that this system has received associated with this adjacency, then the adjacency is marked as being in "Restart mode" and the adjacency holding time is refreshed -- otherwise, the holding time is not refreshed. The "remaining time" transmitted according to (b) below MUST reflect the actual time after which the adjacency will now expire. Receipt of a normal IIH with the RR bit reset will clear the "Restart mode" state. This procedure allows the restarting router to cause the neighbor to maintain the adjacency long enough for restart to successfully complete, while also preventing repetitive restarts from maintaining an adjacency indefinitely. Whether or not an adjacency is marked as being in "Restart mode" has no effect on adjacency state transitions.
- b. immediately (i.e., without waiting for any currently running timer interval to expire, but with a small random delay of a few tens of milliseconds on LANs to avoid "storms") transmit over the corresponding interface an IIH including the restart TLV with the RR bit clear and the RA bit set, in the case of Point-to-Point adjacencies having updated the "Point-to-Point Three-Way Adjacency" option to reflect any new values received from the (re)starting router. (This allows a restarting router to quickly acquire the correct information to place in its hellos.) The "Remaining Time" MUST be set to the current time (in seconds) before the holding timer on this adjacency is due to expire. If the corresponding interface is a LAN interface, then the Restarting Neighbor System ID SHOULD be set to the System ID of

the router from which the IIH with the RR bit set was received. This is required to correctly associate the acknowledgement and holding time in the case where multiple systems on a LAN restart at approximately the same time. This IIH SHOULD be transmitted before any LSPs or SNPs are transmitted as a result of the receipt of the original IIH.

- c. if the corresponding interface is a Point-to-Point interface, or if the receiving router has the highest LnRouterPriority (with the highest source MAC (Media Access Control) address breaking ties) among those routers to which the receiving router has an adjacency in state "UP" on this interface whose IIHs contain the restart TLV, excluding adjacencies to all routers which are considered in "Restart mode" (note the actual DIS is NOT changed by this process), initiate the transmission over the corresponding interface of a complete set of CSNPs, and set SRMflags on the corresponding interface for all LSPs in the local LSP database.

Otherwise (i.e., if there was no adjacency in the "UP" state to the System ID in question), process the IIH as normal by reinitializing the adjacency and setting the RA bit in the returned IIH.

2.2.2. Use of the SA Bit

The SA bit is used by a starting router to request that its neighbor suppress advertisement of the adjacency to the starting router in the neighbor's LSPs.

A router that is starting has no maintained forwarding function state. This may or may not be the first time the router has started. If this is not the first time the router has started, copies of LSPs generated by this router in its previous incarnation may exist in the LSP databases of other routers in the network. These copies are likely to appear "newer" than LSPs initially generated by the starting router due to the reinitialization of LSP fragment sequence numbers by the starting router. This may cause temporary blackholes to occur until the normal operation of the update process causes the starting router to regenerate and flood copies of its own LSPs with higher sequence numbers. The temporary blackholes can be avoided if the starting router's neighbors suppress advertising an adjacency to the starting router until the starting router has been able to propagate newer versions of LSPs generated by previous incarnations.

When a router receives an IIH with the restart TLV having the SA bit set, if there exists on this interface an adjacency in state "UP" with the same System ID, and in the case of a LAN circuit, with the same source LAN address, then the router MUST suppress advertisement

of the adjacency to the neighbor in its own LSPs. Until an IIH with the SA bit clear has been received, the neighbor advertisement MUST continue to be suppressed. If the adjacency transitions to the "UP" state, the new adjacency MUST NOT be advertised until an IIH with the SA bit clear has been received.

Note that a router that suppresses advertisement of an adjacency MUST NOT use this adjacency when performing its SPF calculation. In particular, if an implementation follows the example guidelines presented in [ISO10589], Annex C.2.5, Step 0:b) "pre-load TENT with the local adjacency database", the suppressed adjacency MUST NOT be loaded into TENT.

2.2.3. Use of PR and PA Bits

The PR bit is used by a router which is planning to initiate a restart to signal to its neighbors that it will be restarting.

The PA bit is sent by the neighbor of a router planning to restart to acknowledge receipt of a restart TLV with the PR bit set.

When the neighbor of a router planning a restart receives an IIH with the restart TLV having the PR bit set, if there exists on this interface an adjacency in state "UP" with the same System ID, and in the case of a LAN circuit, with the same source LAN address, then:

- a. if this is the first IIH with the PR bit set that this system has received associated with this adjacency, then the adjacency is marked as being in "Planned Restart state" and the adjacency holding time is refreshed -- otherwise, the holding time is not refreshed. The "remaining time" transmitted according to (b) below MUST reflect the actual time after which the adjacency will now expire. Receipt of a normal IIH with the PR bit reset will clear the "Planned Restart mode" state. This procedure allows the router planning a restart to cause the neighbor to maintain the adjacency long enough for restart to successfully complete. Whether or not an adjacency is marked as being in "Planned Restart mode" has no effect on adjacency state transitions.
- b. immediately (i.e., without waiting for any currently running timer interval to expire, but with a small random delay of a few tens of milliseconds on LANs to avoid "storms") transmit over the corresponding interface an IIH including the restart TLV with the PR bit clear and the PA bit set. The "Remaining Time" MUST be set to the current time (in seconds) before the holding timer on this adjacency is due to expire. If the corresponding interface is a LAN interface, then the Restarting Neighbor System ID SHOULD be set to the System ID of the router from which the IIH with the

PR bit set was received. This is required to correctly associate the acknowledgement and holding time in the case where multiple systems on a LAN are planning a restart at approximately the same time.

While a control plane restart is in progress it is expected that the restarting router will be unable to respond to topology changes. It is therefore useful to signal a planned restart (if the forwarding plane on the restarting router is maintained) so that the neighbors of the restarting router can determine whether it is safe to maintain the adjacency if other topology changes occur prior to the completion of the restart. Signalling a planned restart in the absence of maintained forwarding plane state is likely to lead to significant traffic loss and MUST NOT be done.

Neighbors of the router which has signaled planned restart SHOULD maintain the adjacency in a planned restart state until it receives an IIH with the RR bit set, receives an IIH with both PR and RR bits clear, or the adjacency holding time expires - whichever occurs first.

While the adjacency is in planned restart state the following actions MAY be taken:

- a. If additional topology changes occur, the adjacency which is in planned restart state MAY be brought down even though the hold time has not yet expired. Given that the neighbor which has signaled a planned restart is not expected to update its forwarding plane in response to signaling of the topology changes (since it is restarting) traffic which transits that node is at risk of being improperly forwarded. On a LAN circuit, if the router in planned restart state is the DIS at any supported level, the adjacency(ies) SHOULD be brought down whenever any LSP update is either generated or received so as to trigger a new DIS election. Failure to do so will compromise the reliability of the Update Process on that circuit. What other criteria are used to determine what topology changes will trigger bringing the adjacency down is a local implementation decision.
- b. If a BFD session to the neighbor which signals a planned restart is in the UP state and subsequently goes DOWN, the event MAY be ignored since it is possible this is an expected side effect of the restart. Use of the Control Plane Independent state as signalled in BFD control packets [RFC5880] SHOULD be considered in the decision to ignore a BFD Session DOWN event

- c. On a Point-to-Point circuit, transmission of LSPs, CSNPs, and PSNPs MAY be suppressed. It is expected that the PDUs will not be received.

2.3. Adjacency (Re)Acquisition

Adjacency (re)acquisition is the first step in (re)initialization. Restarting and starting routers will make use of the RR bit in the restart TLV, though each will use it at different stages of the (re)start procedure.

2.3.1. Adjacency Reacquisition during Restart

The restarting router explicitly notifies its neighbor that the adjacency is being reacquired, and hence that it SHOULD NOT reinitialize the adjacency. This is achieved by setting the RR bit in the restart TLV. When the neighbor of a restarting router receives an IIH with the restart TLV having the RR bit set, if there exists on this interface an adjacency in state "UP" with the same System ID, and in the case of a LAN circuit, with the same source LAN address, then the procedures described in Section 3.2.1 are followed.

A router that does not support the restart capability will ignore the restart TLV and reinitialize the adjacency as normal, returning an IIH without the restart TLV.

On restarting, a router initializes the timer T3, starts the timer T2 for each LSPDB, and for each interface (and in the case of a LAN circuit, for each level) starts the timer T1 and transmits an IIH containing the restart TLV with the RR bit set.

On a Point-to-Point circuit, the restarting router SHOULD set the "Adjacency Three-Way State" to "Init", because the receipt of the acknowledging IIH (with RA set) MUST cause the adjacency to enter the "UP" state immediately.

On a LAN circuit, the LAN-ID assigned to the circuit SHOULD be the same as that used prior to the restart. In particular, for any circuits for which the restarting router was previously DIS, the use of a different LAN-ID would necessitate the generation of a new set of pseudonode LSPs, and corresponding changes in all the LSPs referencing them from other routers on the LAN. By preserving the LAN-ID across the restart, this churn can be prevented. To enable a restarting router to learn the LAN-ID used prior to restart, the LAN-ID specified in an IIH with RR set MUST be ignored.

Transmission of "normal" IIHs is inhibited until the conditions described below are met (in order to avoid causing an unnecessary

adjacency initialization). Upon expiry of the timer T1, it is restarted and the IIH is retransmitted as above.

When a restarting router receives an IIH a local adjacency is established as usual, and if the IIH contains a restart TLV with the RA bit set (and on LAN circuits with a Restart Neighbor System ID that matches that of the local system), the receipt of the acknowledgement over that interface is noted. When the RA bit is set and the state of the remote adjacency is "UP", then the timer T3 is set to the minimum of its current value and the value of the "Remaining Time" field in the received IIH.

On a Point-to-Point link, receipt of an IIH not containing the restart TLV is also treated as an acknowledgement, since it indicates that the neighbor is not restart capable. However, since no CSNP is guaranteed to be received over this interface, the timer T1 is cancelled immediately without waiting for a complete set of CSNPs. Synchronization may therefore be deemed complete even though there are some LSPs which are held (only) by this neighbor (see Section 3.4). In this case, we also want to be certain that the neighbor will reinitialize the adjacency in order to guarantee that the SRMflags have been set on its database, thus ensuring eventual LSPDB synchronization. This is guaranteed to happen except in the case where the Adjacency Three-Way State in the received IIH is "UP" and the Neighbor Extended Local Circuit ID matches the extended local circuit ID assigned by the restarting router. In this case, the restarting router MUST force the adjacency to reinitialize by setting the local Adjacency Three-Way State to "DOWN" and sending a normal IIH.

In the case of a LAN interface, receipt of an IIH not containing the restart TLV is unremarkable since synchronization can still occur so long as at least one of the non-restarting neighboring routers on the LAN supports restart. Therefore, T1 continues to run in this case. If none of the neighbors on the LAN are restart capable, T1 will eventually expire after the locally defined number of retries.

In the case of a Point-to-Point circuit, the "LocalCircuitID" and "Extended Local Circuit ID" information contained in the IIH can be used immediately to generate an IIH containing the correct three-way handshake information. The presence of "Neighbor Extended Local Circuit ID" information that does not match the value currently in use by the local system is ignored (since the IIH may have been transmitted before the neighbor had received the new value from the restarting router), but the adjacency remains in the initializing state until the correct information is received.

In the case of a LAN circuit, the source neighbor information (e.g., SNPAAddress) is recorded and used for adjacency establishment and maintenance as normal.

When BOTH a complete set of CSNPs (for each active level, in the case of a Point-to-Point circuit) and an acknowledgement have been received over the interface, the timer T1 is cancelled.

Once the timer T1 has been cancelled, subsequent IIHs are transmitted according to the normal algorithms, but including the restart TLV with both RR and RA clear.

If a LAN contains a mixture of systems, only some of which support the new algorithm, database synchronization is still guaranteed, but the "old" systems will have reinitialized their adjacencies.

If an interface is active, but does not have any neighboring router reachable over that interface, the timer T1 would never be cancelled, and according to Section 3.4.1.1, the SPF would never be run. Therefore, timer T1 is cancelled after some predetermined number of expirations (which MAY be 1).

2.3.2. Adjacency Acquisition during Start

The starting router wants to ensure that in the event that a neighboring router has an adjacency to the starting router in the "UP" state (from a previous incarnation of the starting router), this adjacency is reinitialized. The starting router also wants neighboring routers to suppress advertisement of an adjacency to the starting router until LSP database synchronization is achieved. This is achieved by sending IIHs with the RR bit clear and the SA bit set in the restart TLV. The RR bit remains clear and the SA bit remains set in subsequent transmissions of IIHs until the adjacency has reached the "UP" state and the initial T1 timer interval (see below) has expired.

Receipt of an IIH with the RR bit clear will result in the neighboring router utilizing normal operation of the adjacency state machine. This will ensure that any old adjacency on the neighboring router will be reinitialized.

Upon receipt of an IIH with the SA bit set, the behavior described in Section 3.2.2 is followed.

Upon starting, a router starts timer T2 for each LSPDB.

For each interface (and in the case of a LAN circuit, for each level), when an adjacency reaches the "UP" state, the starting router

starts a timer T1 and transmits an IIH containing the restart TLV with the RR bit clear and SA bit set. Upon expiry of the timer T1, it is restarted and the IIH is retransmitted with both RR and SA bits set (only the RR bit has changed state from earlier IIHs).

Upon receipt of an IIH with the RR bit set (regardless of whether or not the SA bit is set), the behavior described in Section 2.2.1 is followed.

When an IIH is received by the starting router and the IIH contains a restart TLV with the RA bit set (and on LAN circuits with a Restart Neighbor System ID that matches that of the local system), the receipt of the acknowledgement over that interface is noted.

On a Point-to-Point link, receipt of an IIH not containing the restart TLV is also treated as an acknowledgement, since it indicates that the neighbor is not restart capable. Since the neighbor will have reinitialized the adjacency, this guarantees that SRMflags have been set on its database, thus ensuring eventual LSPDB synchronization. However, since no CSNP is guaranteed to be received over this interface, the timer T1 is cancelled immediately without waiting for a complete set of CSNPs. Synchronization may therefore be deemed complete even though there are some LSPs that are held (only) by this neighbor (see Section 2.4).

In the case of a LAN interface, receipt of an IIH not containing the restart TLV is unremarkable since synchronization can still occur so long as at least one of the non-restarting neighboring routers on the LAN supports restart. Therefore, T1 continues to run in this case. If none of the neighbors on the LAN are restart capable, T1 will eventually expire after the locally defined number of retries. The usual operation of the update process will ensure that synchronization is eventually achieved.

When BOTH a complete set of CSNPs (for each active level, in the case of a Point-to-Point circuit) and an acknowledgement have been received over the interface, the timer T1 is cancelled. Subsequent IIHs sent by the starting router have the RR and RA bits clear and the SA bit set in the restart TLV.

Timer T1 is cancelled after some predetermined number of expirations (which MAY be 1).

When the T2 timer(s) are cancelled or expire, transmission of "normal" IIHs (with RR, RA, and SA bits clear) will begin.

2.3.3. Multiple Levels

A router that is operating as both a Level 1 and a Level 2 router on a particular interface MUST perform the above operations for each level.

On a LAN interface, it MUST send and receive both Level 1 and Level 2 IIHs and perform the CSNP synchronizations independently for each level.

On a Point-to-Point interface, only a single IIH (indicating support for both levels) is required, but it MUST perform the CSNP synchronizations independently for each level.

2.4. Database Synchronization

When a router is started or restarted, it can expect to receive a complete set of CSNPs over each interface. The arrival of the CSNP(s) is now guaranteed, since an IIH with the RR bit set will be retransmitted until the CSNP(s) are correctly received.

The CSNPs describe the set of LSPs that are currently held by each neighbor. Synchronization will be complete when all these LSPs have been received.

When (re)starting, a router starts an instance of timer T2 for each LSPDB as described in Section 3.3.1 or Section 3.3.2. In addition to normal processing of the CSNPs, the set of LSPIDs contained in the first complete set of CSNPs received over each interface is recorded, together with their remaining lifetime. In the case of a LAN interface, a complete set of CSNPs MUST consist of CSNPs received from neighbors that are not restarting. If there are multiple interfaces on the (re)starting router, the recorded set of LSPIDs is the union of those received over each interface. LSPs with a remaining lifetime of zero are NOT so recorded.

As LSPs are received (by the normal operation of the update process) over any interface, the corresponding LSPID entry is removed (it is also removed if an LSP arrives before the CSNP containing the reference). When an LSPID has been held in the list for its indicated remaining lifetime, it is removed from the list. When the list of LSPIDs is empty and the timer T1 has been cancelled for all the interfaces that have an adjacency at this level, the timer T2 is cancelled.

At this point, the local database is guaranteed to contain all the LSP(s) (either the same sequence number or a more recent sequence number) that were present in the neighbors' databases at the time of

(re)starting. LSPs that arrived in a neighbor's database after the time of (re)starting may or may not be present, but the normal operation of the update process will guarantee that they will eventually be received. At this point, the local database is deemed to be "synchronized".

Since LSPs mentioned in the CSNP(s) with a zero remaining lifetime are not recorded, and those with a short remaining lifetime are deleted from the list when the lifetime expires, cancellation of the timer T2 will not be prevented by waiting for an LSP that will never arrive.

2.4.1. LSP Generation and Flooding and SPF Computation

The operation of a router starting, as opposed to restarting, is somewhat different. These two cases are dealt with separately below.

2.4.1.1. Restarting

In order to avoid causing unnecessary routing churn in other routers, it is highly desirable that the router's own LSPs generated by the restarting system are the same as those previously present in the network (assuming no other changes have taken place). It is important therefore not to regenerate and flood the LSPs until all the adjacencies have been re-established and any information required for propagation into the local LSPs is fully available. Ideally, the information is loaded into the LSPs in a deterministic way, such that the same information occurs in the same place in the same LSP (and hence the LSPs are identical to their previous versions). If this can be achieved, the new versions may not even cause SPF to be run in other systems. However, provided the same information is included in the set of LSPs (albeit in a different order, and possibly different LSPs), the result of running the SPF will be the same and will not cause churn to the forwarding tables.

In the case of a restarting router, none of the router's own LSPs are transmitted, nor are the router's own forwarding tables updated while the timer T3 is running.

Redistribution of inter-level information MUST be regenerated before this router's LSP is flooded to other nodes. Therefore, the Level-n non-pseudonode LSP(s) MUST NOT be flooded until the other level's T2 timer has expired and its SPF has been run. This ensures that any inter-level information that is to be propagated can be included in the Level-n LSP(s).

During this period, if one of the router's own (including pseudonodes) LSPs is received, which the local router does not

currently have in its own database, it is NOT purged. Under normal operation, such an LSP would be purged, since the LSP clearly should not be present in the global LSP database. However, in the present circumstances, this would be highly undesirable, because it could cause premature removal of a router's own LSP -- and hence churn in remote routers. Even if the local system has one or more of the router's own LSPs (which it has generated, but not yet transmitted), it is still not valid to compare the received LSP against this set, since it may be that as a result of propagation between Level 1 and Level 2 (or vice versa), a further router's own LSP will need to be generated when the LSP databases have synchronized.

During this period, a restarting router SHOULD send CSNPs as it normally would. Information about the router's own LSPs MAY be included, but if it is included it MUST be based on LSPs that have been received, not on versions that have been generated (but not yet transmitted). This restriction is necessary to prevent premature removal of an LSP from the global LSP database.

When the timer T2 expires or is cancelled indicating that synchronization for that level is complete, the SPF for that level is run in order to derive any information that is required to be propagated to another level, but the forwarding tables are not yet updated.

Once the other level's SPF has run and any inter-level propagation has been resolved, the router's own LSPs can be generated and flooded. Any own LSPs that were previously ignored, but that are not part of the current set of own LSPs (including pseudonodes), MUST then be purged. Note that it is possible that a Designated Router change may have taken place, and consequently the router SHOULD purge those pseudonode LSPs that it previously owned, but that are now no longer part of its set of pseudonode LSPs.

When all the T2 timers have expired or been cancelled, the timer T3 is cancelled and the local forwarding tables are updated.

If the timer T3 expires before all the T2 timers have expired or been cancelled, this indicates that the synchronization process is taking longer than the minimum holding time of the neighbors. The router's own LSP(s) for levels that have not yet completed their first SPF computation are then flooded with the overload bit set to indicate that the router's LSPDB is not yet synchronized (and therefore other routers MUST NOT compute routes through this router). Normal operation of the update process resumes, and the local forwarding tables are updated. In order to prevent the neighbor's adjacencies from expiring, IIHs with the normal interface value for the holding time are transmitted over all interfaces with neither RR nor RA set

in the restart TLV. This will cause the neighbors to refresh their adjacencies. The router's own LSP(s) will continue to have the overload bit set until timer T2 has expired or been cancelled.

2.4.1.2. Starting

In the case of a starting router, as soon as each adjacency is established, and before any CSNP exchanges, the router's own zeroth LSP is transmitted with the overload bit set. This prevents other routers from computing routes through the router until it has reliably acquired the complete set of LSPs. The overload bit remains set in subsequent transmissions of the zeroth LSP (such as will occur if a previous copy of the router's own zeroth LSP is still present in the network) while any timer T2 is running.

When all the T2 timers have been cancelled, the router's own LSP(s) MAY be regenerated with the overload bit clear (assuming the router is not in fact overloaded, and there is no other reason, such as incomplete BGP convergence, to keep the overload bit set) and flooded as normal.

Other LSPs owned by this router (including pseudonodes) are generated and flooded as normal, irrespective of the timer T2. The SPF is also run as normal and the Routing Information Base (RIB) and Forwarding Information Base (FIB) updated as routes become available.

To avoid the possible formation of temporary blackholes, the starting router sets the SA bit in the restart TLV (as described in Section 3.3.2) in all IIHs that it sends.

When all T2 timers have been cancelled, the starting router MUST transmit IIHs with the SA bit clear.

3. State Tables

This section presents state tables that summarize the behaviors described in this document. Other behaviors, in particular adjacency state transitions and LSP database update operation, are NOT included in the state tables except where this document modifies the behaviors described in [ISO10589] and [RFC5303].

The states named in the columns of the tables below are a mixture of states that are specific to a single adjacency (ADJ suppressed, ADJ Seen RA, ADJ Seen CSNP) and states that are indicative of the state of the protocol instance (Running, Restarting, Starting, SPF Wait).

Three state tables are presented from the point of view of a running router, a restarting router, and a starting router.

3.1. Running Router

Event	Running	ADJ suppressed
RX PR	Set Planned Restart state. Send PA	
RX PR clr and RR clr	Clear Planned Restart State	
RX RR	Maintain ADJ State Send RA Set SRM, send CSNP (Note 1) Update Hold Time, set Restart Mode (Note 2)	
RX RR clr	Clr Restart mode	
RX SA	Suppress IS neighbor TLV in LSP(s) Goto ADJ Suppressed	
RX SA clr		Unsuppress IS neighbor TLV in LSP(s) Goto Running

Note 1: CSNPs are sent by routers in accordance with Section 2.2.1c

Note 2: If Restart Mode clear

3.2. Restarting Router

Event	Restarting	ADJ Seen RA	ADJ Seen CSNP	SPF Wait
Restart planned	Send PR			
Planned restart canceled	Send PR clr			
Router	Send IIH/RR			

restarts	ADJ Init Start T1,T2,T3			
RX RR	Send RA			
RX RA	Adjust T3 Goto ADJ Seen RA		Cancel T1 Adjust T3	
RX CSNP set	Goto ADJ Seen CSNP	Cancel T1		
RX IIH w/o Restart TLV	Cancel T1 (Point- to-point only)			
T1 expires	Send IIH/RR Restart T1	Send IIH/RR Restart T1	Send IIH/RR Restart T1	
T1 expires nth time	Send IIH/ normal	Send IIH/ normal	Send IIH/ normal	
T2 expires	Trigger SPF Goto SPF Wait			
T3 expires	Set overload bit Flood local LSPs Update fwd plane			
LSP DB Sync	Cancel T2, and T3 Trigger SPF Goto SPF wait			
All SPF done				Clear overload bit Update fwd plane Flood local LSPs Goto Running

3.3. Starting Router

Event	Starting	ADJ Seen RA	ADJ Seen CSNP
Router starts	Send IIH/SA Start T1,T2		
RX RR	Send RA		
RX RA	Goto ADJ Seen RA		Cancel T1
RX CSNP Set	Goto ADJ Seen CSNP	Cancel T1	
RX IIH w no Restart TLV	Cancel T1 (Point-to-Point only)		
ADJ UP	Start T1 Send local LSPs with overload bit set		
T1 expires	Send IIH/RR and SA Restart T1	Send IIH/RR and SA Restart T1	Send IIH/RR and SA Restart T1
T1 expires nth time	Send IIH/SA	Send IIH/SA	Send IIH/SA
T2 expires	Clear overload bit Send IIH normal Goto Running		
LSP DB Sync	Cancel T2 Clear overload bit Send IIH normal		

4. IANA Considerations

This document defines the following IS-IS TLV that is listed in the IS-IS TLV codepoint registry:

Type	Description	IIH	LSP	SNP
211	Restart TLV	y	n	n

5. Security Considerations

Any new security issues raised by the procedures in this document depend upon the ability of an attacker to inject a false but apparently valid IIH, the ease/difficulty of which has not been altered.

If the RR bit is set in a false IIH, neighbors who receive such an IIH will continue to maintain an existing adjacency in the "UP" state and may (re)send a complete set of CSNPs. While the latter action is wasteful, neither action causes any disruption in correct protocol operation.

If the RA bit is set in a false IIH, a (re)starting router that receives such an IIH may falsely believe that there is a neighbor on the corresponding interface that supports the procedures described in this document. In the absence of receipt of a complete set of CSNPs on that interface, this could delay the completion of (re)start procedures by requiring the timer T1 to time out the locally defined maximum number of retries. This behavior is the same as would occur on a LAN where none of the (re)starting router's neighbors support the procedures in this document and is covered in Sections 2.3.1 and 2.3.2.

If an SA bit is set in a false IIH, this could cause suppression of the advertisement of an IS neighbor, which could either continue for an indefinite period or occur intermittently with the result being a possible loss of reachability to some destinations in the network and/or increased frequency of LSP flooding and SPF calculation.

The possibility of IS-IS PDU spoofing can be reduced by the use of authentication as described in [RFC1195] and [ISO10589], and especially the use of cryptographic authentication as described in [RFC5304] and [RFC5310].

6. Manageability Considerations

These extensions that have been designed, developed, and deployed for many years do not have any new impact on management and operation of the IS-IS protocol via this standardization process.

7. Acknowledgements

For RFC 5306 the authors acknowledged contributions made by Jeff Parker, Radia Perlman, Mark Schaefer, Naiming Shen, Nischal Sheth, Russ White, and Rena Yang.

The authors of this updated version acknowledge the contribution of Mike Shand, co-author of RFC 5306.

8. Normative References

- [ISO10589] International Organization for Standardization, "Intermediate system to Intermediate system intra-domain routeing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)", ISO/IEC 10589:2002, Second Edition, Nov 2002.
- [RFC1195] Callon, R., "Use of OSI IS-IS for routing in TCP/IP and dual environments", RFC 1195, DOI 10.17487/RFC1195, December 1990, <<https://www.rfc-editor.org/info/rfc1195>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC5303] Katz, D., Saluja, R., and D. Eastlake 3rd, "Three-Way Handshake for IS-IS Point-to-Point Adjacencies", RFC 5303, DOI 10.17487/RFC5303, October 2008, <<https://www.rfc-editor.org/info/rfc5303>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", RFC 5880, DOI 10.17487/RFC5880, June 2010, <<https://www.rfc-editor.org/info/rfc5880>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

Appendix A. Summary of Changes from RFC 5306

This document extends RFC 5306 by introducing support for signalling the neighbors of a restarting router that a planned restart is about to occur. This allows the neighbors to be aware of the state of the restarting router so that appropriate action may be taken if other topology changes occur while the planned restart is in progress. Since the forwarding plane of the restarting router is maintained based upon the pre-restart state of the network, additional topology changes introduce the possibility that traffic may be lost if paths via the restarting router continue to be used while the restart is in progress.

In support of this new functionality two new flags have been introduced:

- PR - Restart is planned
- PA - Planned restart acknowledgement

No changes to the post restart exchange between the restarting router and its neighbors have been introduced.

Authors' Addresses

Les Ginsberg
Cisco Systems, Inc.

Email: ginsberg@cisco.com

Paul Wells
Cisco Systems, Inc.

Email: pauwells@cisco.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: 13 August 2022

S. Litkowski
Cisco Systems
Y. Qu
Futurewei
P. Sarkar
Individual
I. Chen
The MITRE Corporation
J. Tantsura
Microsoft
9 February 2022

YANG Data Model for IS-IS Segment Routing
draft-ietf-isis-sr-yang-12

Abstract

This document defines a YANG data module that can be used to configure and manage IS-IS Segment Routing, as well as a YANG data module for the management of Signaling Maximum SID Depth (MSD) using IS-IS.

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 <https://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 13 August 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Overview	2
1.1. Requirements Language	3
1.2. Tree Diagrams	3
2. IS-IS MSD	3
2.1. IS-IS MSD YANG Module	3
3. IS-IS Segment Routing	7
3.1. IS-IS Segment Routing configuration	10
3.1.1. Segment Routing activation	10
3.1.2. Advertising mapping server policy	10
3.1.3. IP Fast reroute	11
3.2. IS-IS Segment Routing YANG Module	11
4. Security Considerations	26
5. Contributors	27
6. Acknowledgements	27
7. IANA Considerations	27
8. Normative References	27
Authors' Addresses	29

1. Overview

YANG [RFC7950] is a data definition language used to define the contents of a conceptual data store that allows networked devices to be managed using NETCONF [RFC6241]. YANG is proving relevant beyond its initial confines, as bindings to other interfaces (e.g., ReST) and encodings other than XML (e.g., JSON) are being defined. Furthermore, YANG data models can be used as the basis for implementation of other interfaces, such as CLI and programmatic APIs.

This document defines a YANG data module that can be used to configure and manage IS-IS Segment Routing [RFC8667] and it is an augmentation to the IS-IS YANG data model.

This document also defines a YANG data module for the management of Signaling Maximum SID Depth (MSD) using IS-IS [RFC8491], which augments the base IS-IS YANG data model.

The YANG modules in this document conform to the Network Management Datastore Architecture (NMDA) [RFC8342].

1.1. Requirements Language

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.

1.2. Tree Diagrams

This document uses the graphical representation of data models defined in [RFC8340].

2. IS-IS MSD

This document defines a module for Signaling Maximum SID Depth (MSD) using IS-IS[RFC8667]. It is an augmentation of the IS-IS base model.

The figure below describes the overall structure of the isis-msd YANG module:

```
module: ietf-isis-msd
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis/isis:database
      /isis:levels/isis:lsp/isis:router-capabilities:
        +--ro node-msd-tlv
          +--ro node-msds* [msd-type]
            +--ro msd-type      identityref
            +--ro msd-value?    uint8
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis/isis:database
      /isis:levels/isis:lsp/isis:extended-is-neighbor
        /isis:neighbor:
          +--ro link-msd-sub-tlv
            +--ro link-msds* [msd-type]
              +--ro msd-type      identityref
              +--ro msd-value?    uint8
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis/isis:database
      /isis:levels/isis:lsp/isis:mt-is-neighbor/isis:neighbor:
        +--ro link-msd-sub-tlv
          +--ro link-msds* [msd-type]
            +--ro msd-type      identityref
            +--ro msd-value?    uint8
```

2.1. IS-IS MSD YANG Module

```
<CODE BEGINS> file "ietf-isis-msd@2022-02-09.yang"
module ietf-isis-msd {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:yang:ietf-isis-msd";
  prefix isis-msd;

  import ietf-routing {
    prefix rt;
    reference "RFC 8349: A YANG Data Model for Routing
              Management (NMDA Version)";
  }

  import ietf-isis {
    prefix isis;
  }

  import ietf-mpls-msd {
    prefix mpls-msd;
  }

  organization
    "IETF LSR - LSR Working Group";
  contact
    "WG Web:  <https://tools.ietf.org/wg/mpls/>
    WG List:  <mailto:mpls@ietf.org>

    Author:   Yingzhen Qu
              <mailto:yingzhen.qu@futurewei.com>
    Author:   Acee Lindem
              <mailto:acee@cisco.com>
    Author:   Stephane Litkowski
              <mailto:slitkows.ietf@gmail.com>
    Author:   Jeff Tantsura
              <mailto:jefftant.ietf@gmail.com>

    ";
  description
    "The YANG module augments the base ISIS model to
    manage different types of MSDs.

    This YANG model conforms to the Network Management
    Datastore Architecture (NMDA) as described in RFC 8342.

    Copyright (c) 2022 IETF Trust and the persons identified as
    authors of the code.  All rights reserved.

    Redistribution and use in source and binary forms, with or
    without modification, is permitted pursuant to, and subject
```

to the license terms contained in, the Simplified BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.

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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.";

```
reference "RFC XXXX: YANG Data Model for OSPF MSD.";

revision 2022-02-09 {
  description
    "Initial Version";
  reference "RFC XXXX: YANG Data Model for ISIS MSD.";
}

grouping link-msd-sub-tlv {
  description
    "Link Maximum SID Depth (MSD) grouping for an interface.";
  container link-msd-sub-tlv {
    list link-msds {
      key "msd-type";
      leaf msd-type {
        type identityref {
          base mpls-msd:msd-base-type;
        }
        description
          "MSD-Types";
      }
      leaf msd-value {
        type uint8;
        description
          "MSD value, in the range of 0-255.";
      }
      description
        "List of link MSDs";
    }
    description
      "Link MSD sub-tlvs.";
  }
}
```



```
/* Node MSD TLV */
augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:router-capabilities" {
    when "/rt:routing/rt:control-plane-protocols/"+
        "rt:control-plane-protocol/rt:type = 'isis:isis'" {
        description
            "This augment ISIS routing protocol when used";
    }
    description
        "This augments ISIS protocol LSDB router capability.";
    container node-msd-tlv {
        list node-msds {
            key "msd-type";
            leaf msd-type {
                type identityref {
                    base mpls-msd:msd-base-type;
                }
                description
                    "MSD-Types";
            }
            leaf msd-value {
                type uint8;
                description
                    "MSD value, in the range of 0-255.";
            }
            description
                "Node MSD is the smallest link MSD supported by
                 the node.";
        }
        description
            "Node MSD is the number of SIDs supported by a node.";
        reference
            "RFC 8476: Signaling Maximum SID Depth (MSD) Using OSPF";
    }
}

/* link MSD sub-tlv */
augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:extended-is-neighbor/isis:neighbor" {
    when "/rt:routing/rt:control-plane-protocols/"+
        "rt:control-plane-protocol/rt:type = 'isis:isis'" {
        description
            "This augment ISIS routing protocol when used";
    }
}
```

```

    }
    description
      "This augments ISIS protocol LSDB neighbor with
      Link MSD sub-TLV.";

    uses link-msd-sub-tlv;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:mt-is-neighbor/isis:neighbor" {
    when "/rt:routing/rt:control-plane-protocols/"+
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol LSDB neighbor.";

    uses link-msd-sub-tlv;
  }
}
<CODE ENDS>

```

3. IS-IS Segment Routing

This document defines a model for IS-IS Segment Routing feature. It is an augmentation of the IS-IS base model.

The IS-IS SR YANG module requires support for the base segment routing module [I-D.ietf-spring-sr-yang], which defines the global segment routing configuration independent of any specific routing protocol configuration, and support of IS-IS base model [I-D.ietf-isis-yang-isis-cfg] which defines basic IS-IS configuration and state.

The figure below describes the overall structure of the isis-sr YANG module:

```

module: ietf-isis-sr
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis:
    +--rw segment-routing
    |   +--rw enabled?      boolean
    |   +--rw bindings
    |       +--rw advertise
    |       |   +--rw policies*  string

```

```

|   +--rw receive?      boolean
+--rw protocol-srgb {sr-mpls:protocol-srgb}?
  +--rw srgb* [lower-bound upper-bound]
    +--rw lower-bound   uint32
    +--rw upper-bound   uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:interfaces
  /isis:interface:
    +--rw segment-routing
      +--rw adjacency-sid
        +--rw adj-sids* [value]
          +--rw value-type? enumeration
          +--rw value      uint32
          +--rw protected? boolean
        +--rw advertise-adj-group-sid* [group-id]
          +--rw group-id   uint32
        +--rw advertise-protection? enumeration
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:interfaces
  /isis:interface/isis:fast-reroute:
    +--rw ti-lfa {ti-lfa}?
      +--rw enable?      boolean
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:interfaces
  /isis:interface/isis:fast-reroute/isis:lfa/isis:remote-lfa:
    +--rw use-segment-routing-path? boolean {remote-lfa-sr}?
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:interfaces
  /isis:interface/isis:adjacencies/isis:adjacency:
    +--ro adjacency-sid* [value]
      +--ro af?          iana-rt-types:address-family
      +--ro value        uint32
      +--ro weight?      uint8
      +--ro protection-requested? boolean
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database
  /isis:levels/isis:lsp/isis:router-capabilities:
    +--ro sr-capability
      +--ro sr-capability
        +--ro sr-capability-bits* identityref
      +--ro global-blocks
        +--ro global-block* []
          +--ro range-size?   uint32
          +--ro sid-sub-tlv
            +--ro sid?      uint32
    +--ro sr-algorithms
      +--ro sr-algorithm*   uint8
    +--ro local-blocks

```

```

    | +--ro local-block* []
    |   +--ro range-size?   uint32
    |   +--ro sid-sub-tlv
    |       +--ro sid?   uint32
+--ro srms-preference
  +--ro preference?   uint8
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database/isis:levels
  /isis:lsp/isis:extended-is-neighbor/isis:neighbor:
+--ro sid-list* [value]
  +--ro adj-sid-flags
  | +--ro bits*   identityref
  +--ro weight?   uint8
  +--ro neighbor-id?   isis:system-id
  +--ro value      uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database
  /isis:levels/isis:lsp/isis:mt-is-neighbor/isis:neighbor:
+--ro sid-list* [value]
  +--ro adj-sid-flags
  | +--ro bits*   identityref
  +--ro weight?   uint8
  +--ro neighbor-id?   isis:system-id
  +--ro value      uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database
  /isis:levels/isis:lsp/isis:extended-ipv4-reachability
  /isis:prefixes:
+--ro sid-list* [value]
  +--ro prefix-sid-flags
  | +--ro bits*   identityref
  +--ro algorithm?   uint8
  +--ro value      uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database
  /isis:levels/isis:lsp/isis:mt-extended-ipv4-reachability
  /isis:prefixes:
+--ro sid-list* [value]
  +--ro prefix-sid-flags
  | +--ro bits*   identityref
  +--ro algorithm?   uint8
  +--ro value      uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/isis:isis/isis:database
  /isis:levels/isis:lsp/isis:ipv6-reachability/isis:prefixes:
+--ro sid-list* [value]
  +--ro prefix-sid-flags
  | +--ro bits*   identityref

```

```

    +--ro algorithm?          uint8
    +--ro value                uint32
augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis/isis:database
    /isis:levels/isis:lsp/isis:mt-ipv6-reachability/isis:prefixes:
    +--ro sid-list* [value]
    +--ro prefix-sid-flags
    |   +--ro bits* identityref
    +--ro algorithm?          uint8
    +--ro value                uint32
augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/isis:isis/isis:database
    /isis:levels/isis:lsp:
    +--ro segment-routing-bindings* [fec range]
    +--ro fec                  string
    +--ro range                uint16
    +--ro sid-binding-flags
    |   +--ro bits* identityref
    +--ro binding
    +--ro prefix-sid
    +--ro sid-list* [value]
    +--ro prefix-sid-flags
    |   +--ro bits* identityref
    +--ro algorithm?          uint8
    +--ro value                uint32

```

3.1. IS-IS Segment Routing configuration

3.1.1. Segment Routing activation

Activation of segment-routing IS-IS is done by setting the "enable" leaf to true. This triggers advertisement of segment-routing extensions based on the configuration parameters that have been setup using the base segment routing module.

3.1.2. Advertising mapping server policy

The base segment routing module defines mapping server policies. By default, IS-IS will not advertise nor receive any mapping server entry. The IS-IS segment-routing module allows to advertise one or multiple mapping server policies through the "bindings/advertise/policies" leaf-list. The "bindings/receive" leaf allows to enable the reception of mapping server entries.

3.1.3. IP Fast reroute

IS-IS SR model augments the fast-reroute container under interface. It brings the ability to activate TI-LFA (topology independent LFA) and also enhances remote LFA to use segment-routing tunneling instead of LDP.

3.2. IS-IS Segment Routing YANG Module

```
<CODE BEGINS> file "ietf-isis-sr@2022-02-09.yang"
module ietf-isis-sr {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:"
    + "yang:ietf-isis-sr";
  prefix isis-sr;

  import ietf-routing {
    prefix "rt";
    reference
      "RFC 8349 - A YANG Data Model for Routing
        Management (NMDA Version)";
  }

  import ietf-segment-routing-common {
    prefix "sr-cmn";
    reference
      "RFC 9020 - YANG Data Model for Segment Routing";
  }

  import ietf-segment-routing-mpls {
    prefix "sr-mpls";
    reference
      "RFC 9020 - YANG Data Model for Segment Routing";
  }

  import ietf-isis {
    prefix "isis";
  }

  import iana-routing-types {
    prefix "iana-rt-types";
    reference "RFC 8294 - Common YANG Data Types for the
      Routing Area";
  }

  organization
    "IETF LSR - LSR Working Group";
```

contact

"WG List: <<mailto:lsr@ietf.org>>

Editor: Stephane Litkowski
<<mailto:stephane.litkowski@orange.com>>

Author: Acee Lindem
<<mailto:acee@cisco.com>>

Author: Yingzhen Qu
<<mailto:yingzhen.qu@futurewei.com>>

Author: Pushpasis Sarkar
<<mailto:pushpasis.ietf@gmail.com>>

Author: Ing-Wher Chen
<<mailto:ingwherchen@mitre.org>>

Author: Jeff Tantsura
<<mailto:jefftant.ietf@gmail.com>>

";

description

"The YANG module defines a generic configuration model for Segment routing ISIS extensions common across all of the vendor implementations.

This YANG model conforms to the Network Management Datastore Architecture (NMDA) as described in RFC 8342.

Copyright (c) 2022 IETF Trust and the persons identified as authors of the code. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, is permitted pursuant to, and subject to the license terms contained in, the Simplified BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX (<https://www.rfc-editor.org/info/rfcXXXX>); see the RFC itself for full legal notices.

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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.";

```
reference "RFC XXXX";

revision 2022-02-09 {
  description
    "Initial revision.";
  reference "RFC XXXX";
}

/* Identities */
identity sr-capability {
  description
    "Base identity for ISIS SR-Capabilities sub-TLV flgs";
}

identity mpls-ipv4 {
  base sr-capability;
  description
    "If set, then the router is capable of
    processing SR MPLS encapsulated IPv4 packets
    on all interfaces.";
}

identity mpls-ipv6 {
  base sr-capability;
  description
    "If set, then the router is capable of
    processing SR MPLS encapsulated IPv6 packets
    on all interfaces.";
}

identity prefix-sid-bit {
  description
    "Base identity for prefix sid sub-tlv bits.";
}

identity r-bit {
  base prefix-sid-bit;
  description
    "Re-advertisement Flag.";
}

identity n-bit {
  base prefix-sid-bit;
  description
    "Node-SID Flag.";
}
```



```
identity p-bit {
  base prefix-sid-bit;
  description
    "No-PHP (No Penultimate Hop-Popping) Flag.";
}

identity e-bit {
  base prefix-sid-bit;
  description
    "Explicit NULL Flag.";
}

identity v-bit {
  base prefix-sid-bit;
  description
    "Value Flag.";
}

identity l-bit {
  base prefix-sid-bit;
  description
    "Local Flag.";
}

identity adj-sid-bit {
  description
    "Base identity for adj sid sub-tlv bits.";
}

identity f-bit {
  base adj-sid-bit;
  description
    "Address-Family flag.";
}

identity b-bit {
  base adj-sid-bit;
  description
    "Backup flag.";
}

identity vi-bit {
  base adj-sid-bit;
  description
    "Value/Index flag.";
}

identity lo-bit {
```

```
    base adj-sid-bit;
    description
        "Local flag.";
}

identity s-bit {
    base adj-sid-bit;
    description
        "Group flag.";
}

identity pe-bit {
    base adj-sid-bit;
    description
        "Persistent flag.";
}

identity sid-binding-bit {
    description
        "Base identity for sid binding tlv bits.";
}

identity af-bit {
    base sid-binding-bit;
    description
        "Address-Family flag.";
}

identity m-bit {
    base sid-binding-bit;
    description
        "Mirror Context flag.";
}

identity sf-bit {
    base sid-binding-bit;
    description
        "S flag. If set, the binding label tlv should be flooded
        across the entire routing domain.";
}

identity d-bit {
    base sid-binding-bit;
    description
        "Leaking flag.";
}

identity a-bit {
```

```
    base sid-binding-bit;
    description
        "Attached flag.";
}

/* Features */

feature remote-lfa-sr {
    description
        "Enhance rLFA to use SR path.";
}

feature ti-lfa {
    description
        "Enhance IPFRR with ti-lfa
        support";
}

/* Groupings */

grouping sid-sub-tlv {
    description "SID/Label sub-TLV grouping.";
    container sid-sub-tlv {
        description
            "Used to advertise the SID/Label associated with a
            prefix or adjacency.";
        leaf sid {
            type uint32;
            description
                "Segment Identifier (SID) - A 20 bit label or
                32 bit SID.";
        }
    }
}

grouping sr-capability {
    description
        "SR capability grouping.";
    container sr-capability {
        description
            "Segment Routing capability.";
        container sr-capability {
            leaf-list sr-capability-bits {
                type identityref {
                    base sr-capability;
                }
            }
            description "SR Capability sub-tlv flags list.";
        }
    }
}
```

```
    }
    description
      "SR Capability Flags.";
  }
  container global-blocks {
    description
      "Segment Routing Global Blocks.";
    list global-block {
      description "Segment Routing Global Block.";
      leaf range-size {
        type uint32;
        description "The SID range.";
      }
      uses sid-sub-tlv;
    }
  }
}

grouping sr-algorithm {
  description
    "SR algorithm grouping.";
  container sr-algorithms {
    description "All SR algorithms.";
    leaf-list sr-algorithm {
      type uint8;
      description
        "The Segment Routing (SR) algorithms that the router is
        currently using.";
    }
  }
}

grouping srlb {
  description
    "SR Local Block grouping.";
  container local-blocks {
    description "List of SRLBs.";
    list local-block {
      description "Segment Routing Local Block.";
      leaf range-size {
        type uint32;
        description "The SID range.";
      }
      uses sid-sub-tlv;
    }
  }
}
```

```
grouping srms-preference {
  description "The SRMS preference TLV is used to advertise
              a preference associated with the node that acts
              as an SR Mapping Server.";
  container srms-preference {
    description "SRMS Preference TLV.";
    leaf preference {
      type uint8 {
        range "0 .. 255";
      }
      description "SRMS preference TLV, vlaue from 0 to 255.";
    }
  }
}

grouping adjacency-state {
  description
    "This group will extend adjacency state.";
  list adjacency-sid {
    key value;
    config false;
    leaf af {
      type iana-rt-types:address-family;
      description
        "Address-family associated with the
        segment ID";
    }
    leaf value {
      type uint32;
      description
        "Value of the Adj-SID.";
    }
    leaf weight {
      type uint8;
      description
        "Weight associated with
        the adjacency SID.";
    }
    leaf protection-requested {
      type boolean;
      description
        "Describe if the adjacency SID
        must be protected.";
    }
  }
  description
    "List of adjacency Segment IDs.";
}
}
```

```
grouping prefix-segment-id {
  description
    "This group defines segment routing extensions
    for prefixes.";

  list sid-list {
    key value;

    container prefix-sid-flags {
      leaf-list bits {
        type identityref {
          base prefix-sid-bit;
        }
        description
          "Prefix SID Sub-TLV flag bits list.";
      }
      description
        "Describes flags associated with the
        segment ID.";
    }

    leaf algorithm {
      type uint8;
      description
        "Algorithm to be used for path computation.";
    }
    leaf value {
      type uint32;
      description
        "Value of the prefix-SID.";
    }
    description
      "List of segments.";
  }
}

grouping adjacency-segment-id {
  description
    "This group defines segment routing extensions
    for adjacencies.";

  list sid-list {
    key value;

    container adj-sid-flags {
      leaf-list bits {
        type identityref {
          base adj-sid-bit;
        }
      }
    }
  }
}
```

```
    }
    description "Adj sid sub-tlv flags list.";
  }
  description "Adj-sid sub-tlv flags.";
}

leaf weight {
  type uint8;
  description
    "The value represents the weight of the Adj-SID
    for the purpose of load balancing.";
}
leaf neighbor-id {
  type isis:system-id;
  description
    "Describes the system ID of the neighbor
    associated with the SID value. This is only
    used on LAN adjacencies.";
}
leaf value {
  type uint32;
  description
    "Value of the Adj-SID.";
}
description
  "List of segments.";
}
}

grouping segment-routing-binding-tlv {
  list segment-routing-bindings {
    key "fec range";

    leaf fec {
      type string;
      description
        "IP (v4 or v6) range to be bound to SIDs.";
    }

    leaf range {
      type uint16;
      description
        "Describes number of elements to assign
        a binding to.";
    }

    container sid-binding-flags {
      leaf-list bits {
```

```
        type identityref {
            base sid-binding-bit;
        }
        description
            "SID Binding TLV flag bits list.";
    }
    description
        "Binding flags.";
}

container binding {
    container prefix-sid {
        uses prefix-segment-id;
        description
            "Binding prefix SID to the range.";
    }
    description
        "Bindings associated with the range.";
}

description
    "This container describes list of SID/Label bindings.
    ISIS reference is TLV 149.";
}
description
    "Defines binding TLV for database.";
}

/* Cfg */

augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis" {
    when "/rt:routing/rt:control-plane-protocols/" +
        "rt:control-plane-protocol/rt:type = 'isis:isis'" {
        description
            "This augment ISIS routing protocol when used";
    }
    description
        "This augments ISIS protocol configuration
        with segment routing.";

    uses sr-mpls:sr-control-plane;
    container protocol-srgb {
        if-feature sr-mpls:protocol-srgb;
        uses sr-cmn:srgb;
        description
            "Per-protocol SRGB.";
    }
}
```



```
    }
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:interfaces/isis:interface" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol configuration
        with segment routing.";

    uses sr-mpls:igp-interface;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:interfaces/isis:interface"+
    "/isis:fast-reroute" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS IP FRR with TILFA.";

    container ti-lfa {
      if-feature ti-lfa;
      leaf enable {
        type boolean;
        description
          "Enables TI-LFA computation.";
      }
      description
        "TILFA configuration.";
    }
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:interfaces/isis:interface"+
    "/isis:fast-reroute/isis:lfa/isis:remote-lfa" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
```

```
        description
          "This augment ISIS routing protocol when used";
      }
      description
        "This augments ISIS remoteLFA config with
         use of segment-routing path.";

      leaf use-segment-routing-path {
        if-feature remote-lfa-sr;
        type boolean;
        description
          "force remote LFA to use segment routing
           path instead of LDP path.";
      }
    }
  }

  /* Operational states */

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:interfaces/isis:interface" +
    "/isis:adjacencies/isis:adjacency" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol configuration
       with segment routing.";

    uses adjacency-state;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:router-capabilities" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol LSDB router capability.";

    uses sr-capability;
    uses sr-algorithm;
  }
```

```
    uses srlb;
    uses srms-preference;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:extended-is-neighbor/isis:neighbor" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol LSDB neighbor.";
    uses adjacency-segment-id;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:mt-is-neighbor/isis:neighbor" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol LSDB neighbor.";
    uses adjacency-segment-id;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:extended-ipv4-reachability/isis:prefixes" {
    when "/rt:routing/rt:control-plane-protocols/" +
      "rt:control-plane-protocol/rt:type = 'isis:isis'" {
      description
        "This augment ISIS routing protocol when used";
    }
    description
      "This augments ISIS protocol LSDB prefix.";
    uses prefix-segment-id;
  }

  augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
```

```
        "/isis:isis/isis:database/isis:levels/isis:lsp"+
        "/isis:mt-extended-ipv4-reachability/isis:prefixes" {
when "/rt:routing/rt:control-plane-protocols/"+
    "rt:control-plane-protocol/rt:type = 'isis:isis'" {
    description
        "This augment ISIS routing protocol when used";
    }
    description
        "This augments ISIS protocol LSDB prefix.";
    uses prefix-segment-id;
}

augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:ipv6-reachability/isis:prefixes" {
when "/rt:routing/rt:control-plane-protocols/"+
    "rt:control-plane-protocol/rt:type = 'isis:isis'" {
    description
        "This augment ISIS routing protocol when used";
    }
    description
        "This augments ISIS protocol LSDB prefix.";
    uses prefix-segment-id;
}

augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp"+
    "/isis:mt-ipv6-reachability/isis:prefixes" {
when "/rt:routing/rt:control-plane-protocols/"+
    "rt:control-plane-protocol/rt:type = 'isis:isis'" {
    description
        "This augment ISIS routing protocol when used";
    }
    description
        "This augments ISIS protocol LSDB prefix.";
    uses prefix-segment-id;
}

augment "/rt:routing/" +
    "rt:control-plane-protocols/rt:control-plane-protocol"+
    "/isis:isis/isis:database/isis:levels/isis:lsp" {
when "/rt:routing/rt:control-plane-protocols/"+
    "rt:control-plane-protocol/rt:type = 'isis:isis'" {
    description
        "This augment ISIS routing protocol when used";
    }
}
```

```
    description
      "This augments ISIS protocol LSDB.";
      uses segment-routing-binding-tlv;
  }

  /* Notifications */
}
<CODE ENDS>
```

4. Security Considerations

The YANG module specified in this document defines a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC8446].

The NETCONF Access Control Model (NACM) [RFC8341] provides the means to restrict access for particular NETCONF or RESTCONF users to a pre-configured subset of all available NETCONF or RESTCONF protocol operations and content.

There are a number of data nodes defined in this YANG module that are writable/creatable/deletable (i.e., config true, which is the default). These data nodes may be considered sensitive or vulnerable in some network environments. Write operations (e.g., edit-config) to these data nodes without proper protection can have a negative effect on network operations. These are the subtrees and data nodes and their sensitivity/vulnerability:

```
/isis:isis/segment-routing
```

```
/isis:isis/protocol-srgb
```

```
/isis:isis/isis:interfaces/isis:interface/segment-routing
```

Some of the readable data nodes in the modules may be considered sensitive or vulnerable in some network environments. It is thus important to control read access (e.g., via get, get-config, or notification) to these data nodes.

```
/isis:router-capabilities/sr-capability
```

```
/isis:router-capabilities/sr-algorithms
```

```
/isis:router-capabilities/local-blocks  
  
/isis:router-capabilities/srms-preference  
  
/isis:router-capabilities/node-msd-tlv
```

And the augmentations to the ISIS link state database.

Unauthorized access to any data node of these subtrees can disclose the operational state information of IS-IS protocol on this device.

5. Contributors

Authors would like to thank Derek Yeung, Acee Lindem, Yi Yang for their major contributions to the draft.

6. Acknowledgements

MITRE has approved this document for Public Release, Distribution Unlimited, with Public Release Case Number 19-3033.

7. IANA Considerations

The IANA is requested to assign two new URIs from the IETF XML registry ([RFC3688]). Authors are suggesting the following URI:

```
URI: urn:ietf:params:xml:ns:yang:ietf-isis-sr  
Registrant Contact: The IESG.  
XML: N/A, the requested URI is an XML namespace
```

```
URI: urn:ietf:params:xml:ns:yang:ietf-isis-msd  
Registrant Contact: The IESG.  
XML: N/A, the requested URI is an XML namespace
```

This document also requests one new YANG module name in the YANG Module Names registry ([RFC6020]) with the following suggestion :

```
name: ietf-isis-sr  
namespace: urn:ietf:params:xml:ns:yang:ietf-isis-sr  
prefix: isis-sr  
reference: RFC XXXX
```

```
name: ietf-isis-msd  
namespace: urn:ietf:params:xml:ns:yang:ietf-isis-msd  
prefix: isis-msd  
reference: RFC XXXX
```

8. Normative References

- [I-D.ietf-isis-yang-isis-cfg]
Litkowski, S., Yeung, D., Lindem, A., Zhang, J., and L. Lhotka, "YANG Data Model for IS-IS Protocol", Work in Progress, Internet-Draft, draft-ietf-isis-yang-isis-cfg-42, 15 October 2019, <<https://www.ietf.org/archive/id/draft-ietf-isis-yang-isis-cfg-42.txt>>.
- [I-D.ietf-spring-sr-yang]
Litkowski, S., Qu, Y., Sarkar, P., and J. Tantsura, "YANG Data Model for Segment Routing", Work in Progress, Internet-Draft, draft-ietf-spring-sr-yang-15, 28 December 2017, <<http://www.ietf.org/internet-drafts/draft-ietf-spring-sr-yang-15.txt>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.
- [RFC6536] Bierman, A. and M. Bjorklund, "Network Configuration Protocol (NETCONF) Access Control Model", RFC 6536, DOI 10.17487/RFC6536, March 2012, <<https://www.rfc-editor.org/info/rfc6536>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.

- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.
- [RFC8341] Bierman, A. and M. Bjorklund, "Network Configuration Access Control Model", STD 91, RFC 8341, DOI 10.17487/RFC8341, March 2018, <<https://www.rfc-editor.org/info/rfc8341>>.
- [RFC8342] Bjorklund, M., Schoenwaelder, J., Shafer, P., Watsen, K., and R. Wilton, "Network Management Datastore Architecture (NMDA)", RFC 8342, DOI 10.17487/RFC8342, March 2018, <<https://www.rfc-editor.org/info/rfc8342>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.
- [RFC8491] Tantsura, J., Chunduri, U., Aldrin, S., and L. Ginsberg, "Signaling Maximum SID Depth (MSD) Using IS-IS", RFC 8491, DOI 10.17487/RFC8491, November 2018, <<https://www.rfc-editor.org/info/rfc8491>>.
- [RFC8667] Previdi, S., Ed., Ginsberg, L., Ed., Filsfils, C., Bashandy, A., Gredler, H., and B. Decraene, "IS-IS Extensions for Segment Routing", RFC 8667, DOI 10.17487/RFC8667, December 2019, <<https://www.rfc-editor.org/info/rfc8667>>.

Authors' Addresses

Stephane Litkowski
Cisco Systems

Email: slitkows.ietf@gmail.com

Yingzhen Qu
Futurewei

Email: yingzhen.qu@futurewei.com

Pushpasis Sarkar
Individual

Email: pushpasis.ietf@gmail.com

Ing-Wher Chen
The MITRE Corporation

Email: ingwherchen@mitre.org

Jeff Tantsura
Microsoft

Email: jefftant.ietf@gmail.com

IS-IS Working Group
Internet-Draft
Intended status: Standards Track
Expires: April 17, 2020

S. Litkowski
Cisco Systems
D. Yeung
Arrcus, Inc
A. Lindem
Cisco Systems
J. Zhang
Juniper Networks
L. Lhotka
CZ.NIC
October 15, 2019

YANG Data Model for IS-IS Protocol
draft-ietf-isis-yang-isis-cfg-42

Abstract

This document defines a YANG data model that can be used to configure and manage the IS-IS protocol on network elements.

Requirements Language

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.

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 <https://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 April 17, 2020.

Copyright Notice

Copyright (c) 2019 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 (<https://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	2
2. Design of the Data Model	3
2.1. IS-IS Configuration	9
2.2. Multi-topology Parameters	10
2.3. Per-Level Parameters	10
2.4. Per-Interface Parameters	12
2.5. Authentication Parameters	19
2.6. IGP/LDP synchronization	19
2.7. ISO parameters	20
2.8. IP FRR	20
2.9. Operational States	20
3. RPC Operations	21
4. Notifications	21
5. Interaction with Other YANG Modules	22
6. IS-IS YANG Module	23
7. Security Considerations	108
8. Contributors	110
9. Acknowledgements	110
10. IANA Considerations	110
11. References	110
11.1. Normative References	110
11.2. Informative References	115
Appendix A. Example of IS-IS configuration in XML	115
Authors' Addresses	117

1. Introduction

This document defines a YANG [RFC7950] data model for IS-IS routing protocol.

The data model covers configuration of an IS-IS routing protocol instance, as well as, the retrieval of IS-IS operational states.

A simplified tree representation of the data model is presented in Section 2. Tree diagrams used in this document follow the notation defined in [RFC8340].

The module is designed as per the NMDA (Network Management Datastore Architecture) [RFC8342].

2. Design of the Data Model

The IS-IS YANG module augments the "control-plane-protocol" list in the ietf-routing module [RFC8349] with specific IS-IS parameters.

The figure below describes the overall structure of the ietf-isis YANG module:

```

module: ietf-isis
augment /rt:routing/rt:ribs/rt:rib/rt:routes/rt:route:
  +--ro metric?          uint32
  +--ro tag*             uint64
  +--ro route-type?     enumeration
augment /if:interfaces/if:interface:
  +--rw clns-mtu?      uint16 {osi-interface}?
augment /rt:routing/rt:control-plane-protocols/rt:
  control-plane-protocol:
  +--rw isis
    +--rw enable?          boolean {admin-control}?
    +--rw level-type?      level
    +--rw system-id?       system-id
    +--rw maximum-area-addresses? uint8 {maximum-area-addresses}?
    +--rw area-address*    area-address
    +--rw lsp-mtu?         uint16
    +--rw lsp-lifetime?    uint16
    +--rw lsp-refresh?     rt-types:timer-value-seconds16
    |                       {lsp-refresh}?
    +--rw poi-tlv?         boolean {poi-tlv}?
    +--rw graceful-restart {graceful-restart}?
    |   +--rw enable?      boolean
    |   +--rw restart-interval? rt-types:timer-value-seconds16
    |   +--rw helper-enable? boolean
    +--rw nsr {nsr}?
    |   +--rw enable?      boolean
    +--rw node-tags {node-tag}?
    |   +--rw node-tag* [tag]
    |   ...

```

```
+--rw metric-type
|   +--rw value?      enumeration
|   +--rw level-1
|   |   ...
|   +--rw level-2
|   |   ...
+--rw default-metric
|   +--rw value?      wide-metric
|   +--rw level-1
|   |   ...
|   +--rw level-2
|   |   ...
+--rw auto-cost {auto-cost}?
|   +--rw enable?      boolean
|   +--rw reference-bandwidth? uint32
+--rw authentication
|   +--rw (authentication-type)?
|   |   ...
|   +--rw level-1
|   |   ...
|   +--rw level-2
|   |   ...
+--rw address-families {nlpid-control}?
|   +--rw address-family-list* [address-family]
|   |   ...
+--rw mpls
|   +--rw te-rid {te-rid}?
|   |   ...
|   +--rw ldp
|   |   ...
+--rw spf-control
|   +--rw paths?      uint16 {max-ecmp}?
|   +--rw ietf-spf-delay {ietf-spf-delay}?
|   |   ...
+--rw fast-reroute {fast-reroute}?
|   +--rw lfa {lfa}?
+--rw preference
|   +--rw (granularity)?
|   |   ...
+--rw overload
|   +--rw status?      boolean
+--rw overload-max-metric {overload-max-metric}?
|   +--rw timeout?     rt-types:timer-value-seconds16
+--ro spf-log
|   +--ro event* [id]
|   |   ...
+--ro lsp-log
|   +--ro event* [id]
```

```

|      ...
+---ro hostnames
|   +---ro hostname* [system-id]
|      ...
+---ro database
|   +---ro levels* [level]
|      ...
+---ro local-rib
|   +---ro route* [prefix]
|      ...
+---ro system-counters
|   +---ro level* [level]
|      ...
+---ro protected-routes
|   +---ro address-family-stats* [address-family prefix alternate]
|      ...
+---ro unprotected-routes
|   +---ro prefixes* [address-family prefix]
|      ...
+---ro protection-statistics* [frr-protection-method]
|   +---ro frr-protection-method    identityref
|   +---ro address-family-stats* [address-family]
|      ...
+---rw discontinuity-time?          yang:date-and-time
+---rw topologies {multi-topology}?
|   +---rw topology* [name]
|      ...
+---rw interfaces
|   +---rw interface* [name]
|      ...

rpcs:
+---x clear-adjacency
|   +---w input
|       +---w routing-protocol-instance-name -> /rt:routing/
|           control-plane-protocols/
|           control-plane-protocol/name
|       +---w level?                          level
|       +---w interface?                      if:interface-ref
+---x clear-database
|   +---w input
|       +---w routing-protocol-instance-name -> /rt:routing/
|           control-plane-protocols/
|           control-plane-protocol/name
|       +---w level?                          level

notifications:
+---n database-overload

```

```

|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro overload?               enumeration
+---n lsp-too-large
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro interface-name?         if:interface-ref
|   +--ro interface-level?       level
|   +--ro extended-circuit-id?    extended-circuit-id
|   +--ro pdu-size?               uint32
|   +--ro lsp-id?                 lsp-id
+---n if-state-change
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro interface-name?         if:interface-ref
|   +--ro interface-level?       level
|   +--ro extended-circuit-id?    extended-circuit-id
|   +--ro state?                  if-state-type
+---n corrupted-lsp-detected
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro lsp-id?                 lsp-id
+---n attempt-to-exceed-max-sequence
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro lsp-id?                 lsp-id
+---n id-len-mismatch
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?             level
|   +--ro interface-name?         if:interface-ref
|   +--ro interface-level?       level
|   +--ro extended-circuit-id?    extended-circuit-id
|   +--ro pdu-field-len?          uint8
|   +--ro raw-pdu?                binary
+---n max-area-addresses-mismatch
|   +--ro routing-protocol-name?  -> /rt:routing/

```

```

| | | | | control-plane-protocols/
| | | | | control-plane-protocol/name
| | | | |
| | | | | +---ro isis-level? level
| | | | | +---ro interface-name? if:interface-ref
| | | | | +---ro interface-level? level
| | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | +---ro max-area-addresses? uint8
| | | | | +---ro raw-pdu? binary
| | | | |
| | | | | +---n own-lsp-purge
| | | | | | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | | | | | control-plane-protocols/
| | | | | | | | | | | | | control-plane-protocol/name
| | | | | | | | | | | | |
| | | | | | | | | | | | | +---ro isis-level? level
| | | | | | | | | | | | | +---ro interface-name? if:interface-ref
| | | | | | | | | | | | | +---ro interface-level? level
| | | | | | | | | | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | | | | | | | | | +---ro lsp-id? lsp-id
| | | | | | | | | | | | |
| | | | | | | | | | | | | +---n sequence-number-skipped
| | | | | | | | | | | | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | | | | | | | | | | | control-plane-protocols/
| | | | | | | | | | | | | | | | | | | control-plane-protocol/name
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | +---ro isis-level? level
| | | | | | | | | | | | | | | | | | | +---ro interface-name? if:interface-ref
| | | | | | | | | | | | | | | | | | | +---ro interface-level? level
| | | | | | | | | | | | | | | | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | | | | | | | | | | | | | | | +---ro lsp-id? lsp-id
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | +---n authentication-type-failure
| | | | | | | | | | | | | | | | | | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocols/
| | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocol/name
| | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | +---ro isis-level? level
| | | | | | | | | | | | | | | | | | | | | | | | | +---ro interface-name? if:interface-ref
| | | | | | | | | | | | | | | | | | | | | | | | | +---ro interface-level? level
| | | | | | | | | | | | | | | | | | | | | | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | | | | | | | | | | | | | | | | | | | | | +---ro raw-pdu? binary
| | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | +---n authentication-failure
| | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocols/
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocol/name
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro isis-level? level
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro interface-name? if:interface-ref
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro interface-level? level
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro raw-pdu? binary
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---n version-skew
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocols/
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | control-plane-protocol/name

```



```

| +--ro isis-level?                level
| +--ro interface-name?           if:interface-ref
| +--ro interface-level?         level
| +--ro extended-circuit-id?     extended-circuit-id
| +--ro protocol-version?        uint8
| +--ro raw-pdu?                 binary
+---n area-mismatch
| +--ro routing-protocol-name?    -> /rt:routing/
| |                               control-plane-protocols/
| |                               control-plane-protocol/name
| +--ro isis-level?              level
| +--ro interface-name?          if:interface-ref
| +--ro interface-level?         level
| +--ro extended-circuit-id?     extended-circuit-id
| +--ro raw-pdu?                 binary
+---n rejected-adjacency
| +--ro routing-protocol-name?    -> /rt:routing/
| |                               control-plane-protocols/
| |                               control-plane-protocol/name
| +--ro isis-level?              level
| +--ro interface-name?          if:interface-ref
| +--ro interface-level?         level
| +--ro extended-circuit-id?     extended-circuit-id
| +--ro raw-pdu?                 binary
| +--ro reason?                  string
+---n protocols-supported-mismatch
| +--ro routing-protocol-name?    -> /rt:routing/
| |                               control-plane-protocols/
| |                               control-plane-protocol/name
| +--ro isis-level?              level
| +--ro interface-name?          if:interface-ref
| +--ro interface-level?         level
| +--ro extended-circuit-id?     extended-circuit-id
| +--ro raw-pdu?                 binary
| +--ro protocols*               uint8
+---n lsp-error-detected
| +--ro routing-protocol-name?    -> /rt:routing/
| |                               control-plane-protocols/
| |                               control-plane-protocol/name
| +--ro isis-level?              level
| +--ro interface-name?          if:interface-ref
| +--ro interface-level?         level
| +--ro extended-circuit-id?     extended-circuit-id
| +--ro lsp-id?                  lsp-id
| +--ro raw-pdu?                 binary
| +--ro error-offset?            uint32
| +--ro tlv-type?                uint8
+---n adjacency-state-change

```

```

|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?            level
|   +--ro interface-name?        if:interface-ref
|   +--ro interface-level?       level
|   +--ro extended-circuit-id?    extended-circuit-id
|   +--ro neighbor?              string
|   +--ro neighbor-system-id?     system-id
|   +--ro state?                  adj-state-type
|   +--ro reason?                 string
+---n lsp-received
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?            level
|   +--ro interface-name?        if:interface-ref
|   +--ro interface-level?       level
|   +--ro extended-circuit-id?    extended-circuit-id
|   +--ro lsp-id?                 lsp-id
|   +--ro sequence?               uint32
|   +--ro received-timestamp?     yang:timestamp
|   +--ro neighbor-system-id?     system-id
+---n lsp-generation
|   +--ro routing-protocol-name?  -> /rt:routing/
|   |                             control-plane-protocols/
|   |                             control-plane-protocol/name
|   +--ro isis-level?            level
|   +--ro lsp-id?                 lsp-id
|   +--ro sequence?               uint32
|   +--ro send-timestamp?         yang:timestamp

```

2.1. IS-IS Configuration

The IS-IS configuration is divided into:

- o Global parameters.
- o Per-interface configuration (see Section 2.4).

Additional modules may be created to support additional parameters. These additional modules **MUST** augment the `ietf-isis` module.

The model includes optional features, for which the corresponding configuration data nodes are also optional. As an example, the ability to control the administrative state of a particular IS-IS instance is optional. By advertising the feature "admin-control", a

device communicates to the client that it supports the ability to shutdown a particular IS-IS instance.

The global configuration contains usual IS-IS parameters, such as, `lsp-mtu`, `lsp-lifetime`, `lsp-refresh`, `default-metric`, etc.

2.2. Multi-topology Parameters

The model supports multi-topology (MT) IS-IS as defined in [RFC5120].

The "topologies" container is used to enable support of the MT extensions.

The "name" used in the topology list should refer to an existing Routing Information Base (RIB) defined for the device [RFC8349].

Some specific parameters can be defined on a per-topology basis, both at the global level and at the interface level: for example, an interface metric can be defined per topology.

Multiple address families (such as, IPv4 or IPv6) can also be enabled within the default topology. This can be achieved using the address-families container (requiring the "nlpid-control" feature to be supported).

2.3. Per-Level Parameters

Some parameters allow a per-level configuration. For such parameters, the parameter is modeled as a container with three configuration locations:

- o a Top-level container: Corresponds to level-1-2, so the configuration applies to both levels.
- o a Level-1 container: Corresponds to level-1 specific parameters.
- o a Level-2 container: Corresponds to level-2 specific parameters.

```

+--rw priority
|   +--rw value?      uint8
|   +--rw level-1
|   |   +--rw value?  uint8
|   +--rw level-2
|       +--rw value?  uint8
```

Example:

```
<priority>
  <value>250</value>
  <level-1>
    <value>100</value>
  </level-1>
</priority>
```

An implementation MUST prefer a level-specific parameter over a top-level parameter. For example, if the priority is 100 for the level-1 and 250 for the top-level configuration, the implementation must use 100 for the level-1 priority and 250 for the level-2 priority.

Some parameters, such as, "overload bit" and "route preference", are not modeled to support a per-level configuration. If an implementation supports per-level configuration for such parameter, this implementation MUST augment the current model by adding both level-1 and level-2 containers and MUST reuse existing configuration groupings.

Example of augmentation:

```
augment "/rt:routing/" +
  "rt:control-plane-protocols/rt:control-plane-protocol"+
  "/isis:isis/isis:overload" {
  when "rt:type = 'isis:isis'" {
    description
      "This augment IS-IS routing protocol when used";
  }
  description
    "This augments IS-IS overload configuration
    with per-level configuration.";

  container level-1 {
    uses isis:overload-global-cfg;
    description
      "Level 1 configuration.";
  }
  container level-2 {
    uses isis:overload-global-cfg;
    description
      "Level 2 configuration.";
  }
}
```

If an implementation does not support per-level configuration for a parameter modeled with per-level configuration, the implementation should advertise a deviation to announce the non-support of the level-1 and level-2 containers.

Finally, if an implementation supports per-level configuration but does not support the level-1-2 configuration, it should also advertise a deviation.

2.4. Per-Interface Parameters

The per-interface section of the IS-IS instance describes the interface-specific parameters.

The interface is modeled as a reference to an existing interface defined in the "ietf-interfaces" YANG model ([RFC8343]).

Each interface has some interface-specific parameters that may have a different per-level value as described in the previous section. An interface-specific parameter MUST be preferred over an IS-IS global parameter.

Some parameters, such as, hello-padding are defined as containers to allow easy extension by vendor-specific modules.

```

+--rw interfaces
  +--rw interface* [name]
    +--rw name                               if:interface-ref
    +--rw enable?                             boolean {admin-control}?
    +--rw level-type?                          level
    +--rw lsp-pacing-interval?                 rt-types:
    |                                           timer-value-milliseconds
    +--rw lsp-retransmit-interval?             rt-types:
    |                                           timer-value-seconds16
    +--rw passive?                             boolean
    +--rw csnp-interval?                       rt-types:
    |                                           timer-value-seconds16
    +--rw hello-padding
    |   +--rw enable?                          boolean
    +--rw mesh-group-enable?                   mesh-group-state
    +--rw mesh-group?                          uint8
    +--rw interface-type?                      interface-type
    +--rw tag*                                uint32 {prefix-tag}?
    +--rw tag64*                              uint64 {prefix-tag64}?
    +--rw node-flag?                          boolean {node-flag}?
    +--rw hello-authentication
    |   +--rw (authentication-type)?
    |   |   +--:(key-chain) {key-chain}?
    |   |   |   +--rw key-chain?              key-chain:key-chain-ref
    |   |   +--:(password)
    |   |   |   +--rw key?                      string
    |   |   |   +--rw crypto-algorithm?        identityref
    |   +--rw level-1

```

```

+--rw (authentication-type)?
+--:(key-chain) {key-chain}?
|   +--rw key-chain?          key-chain:key-chain-ref
+--:(password)
|   +--rw key?                string
|   +--rw crypto-algorithm?   identityref
+--rw level-2
+--rw (authentication-type)?
+--:(key-chain) {key-chain}?
|   +--rw key-chain?          key-chain:key-chain-ref
+--:(password)
|   +--rw key?                string
|   +--rw crypto-algorithm?   identityref
+--rw hello-interval
+--rw value?                  rt-types:timer-value-seconds16
+--rw level-1
|   +--rw value?              rt-types:timer-value-seconds16
+--rw level-2
|   +--rw value?              rt-types:timer-value-seconds16
+--rw hello-multiplier
+--rw value?                  uint16
+--rw level-1
|   +--rw value?              uint16
+--rw level-2
|   +--rw value?              uint16
+--rw priority
+--rw value?                  uint8
+--rw level-1
|   +--rw value?              uint8
+--rw level-2
|   +--rw value?              uint8
+--rw metric
+--rw value?                  wide-metric
+--rw level-1
|   +--rw value?              wide-metric
+--rw level-2
|   +--rw value?              wide-metric
+--rw bfd {bfd}?
+--rw enable?                  boolean
+--rw local-multiplier?        multiplier
+--rw (interval-config-type)?
+--:(tx-rx-intervals)
|   +--rw desired-min-tx-interval?  uint32
|   +--rw required-min-rx-interval? uint32
+--:(single-interval) {single-minimum-interval}?
|   +--rw min-interval?            uint32
+--rw address-families {nlpid-control}?
+--rw address-family-list* [address-family]

```

```

|      +--rw address-family      iana-rt-types:address-family
+--rw mpls
|   +--rw ldp
|       +--rw igp-sync?    boolean {ldp-igp-sync}?
+--rw fast-reroute {fast-reroute}?
|   +--rw lfa {lfa}?
|       +--rw candidate-enable?    boolean
|       +--rw enable?              boolean
|       +--rw remote-lfa {remote-lfa}?
|       |   +--rw enable?    boolean
|       +--rw level-1
|       |   +--rw candidate-enable?    boolean
|       |   +--rw enable?              boolean
|       |   +--rw remote-lfa {remote-lfa}?
|       |   |   +--rw enable?    boolean
|       +--rw level-2
|       |   +--rw candidate-enable?    boolean
|       |   +--rw enable?              boolean
|       |   +--rw remote-lfa {remote-lfa}?
|       |   |   +--rw enable?    boolean
+--ro adjacencies
|   +--ro adjacency* []
|       +--ro neighbor-sys-type?          level
|       +--ro neighbor-sysid?             system-id
|       +--ro neighbor-extended-circuit-id? extended-circuit-id
|       +--ro neighbor-snpa?              snpa
|       +--ro usage?                      level
|       +--ro hold-timer?                 rt-types:
|       |                               timer-value-seconds16
|       +--ro neighbor-priority?          uint8
|       +--ro lastuptime?                 yang:timestamp
|       +--ro state?                      adj-state-type
+--ro event-counters
|   +--ro adjacency-changes?              uint32
|   +--ro adjacency-number?              uint32
|   +--ro init-fails?                    uint32
|   +--ro adjacency-rejects?             uint32
|   +--ro id-len-mismatch?               uint32
|   +--ro max-area-addresses-mismatch?   uint32
|   +--ro authentication-type-fails?     uint32
|   +--ro authentication-fails?         uint32
|   +--ro lan-dis-changes?               uint32
+--ro packet-counters
|   +--ro level* [level]
|       +--ro level          level-number
|       +--ro iih
|       |   +--ro in?        uint32
|       |   +--ro out?       uint32

```

```

    |
    | +--ro ish
    | |   +--ro in?      uint32
    | |   +--ro out?     uint32
    | +--ro esh
    | |   +--ro in?      uint32
    | |   +--ro out?     uint32
    | +--ro lsp
    | |   +--ro in?      uint32
    | |   +--ro out?     uint32
    | +--ro psnp
    | |   +--ro in?      uint32
    | |   +--ro out?     uint32
    | +--ro csnp
    | |   +--ro in?      uint32
    | |   +--ro out?     uint32
    | +--ro unknown
    | |   +--ro in?      uint32
    +--rw discontinuity-time?      yang:date-and-time
    +--rw topologies {multi-topology}?
    |   +--rw topology* [name]
    |   |   +--rw name      ->
    |   |   |   ../.../rt:ribs/rib/name
    |   +--rw metric
    |   |   +--rw value?      wide-metric
    |   |   +--rw level-1
    |   |   |   +--rw value?  wide-metric
    |   |   +--rw level-2
    |   |   |   +--rw value?  wide-metric
    +--rw rpcs:
    |   +---x clear-adjacency
    |   |   +---w input
    |   |   |   +---w routing-protocol-instance-name  -> /rt:routing/
    |   |   |   |   control-plane-protocols/
    |   |   |   |   control-plane-protocol/name
    |   |   |   +---w level?                          level
    |   |   |   +---w interface?                       if:interface-ref
    |   +---x clear-database
    |   |   +---w input
    |   |   |   +---w routing-protocol-instance-name  -> /rt:routing/
    |   |   |   |   control-plane-protocols/
    |   |   |   |   control-plane-protocol/name
    |   |   |   +---w level?                          level
    +--rw notifications:
    |   +---n database-overload
    |   |   +---ro routing-protocol-name?  -> /rt:routing/
    |   |   |   control-plane-protocols/

```



```

| | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro overload? enumeration
+---n lsp-too-large
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro interface-name? if:interface-ref
| | | | | +---ro interface-level? level
| | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | +---ro pdu-size? uint32
| | | | | +---ro lsp-id? lsp-id
+---n if-state-change
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro interface-name? if:interface-ref
| | | | | +---ro interface-level? level
| | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | +---ro state? if-state-type
+---n corrupted-lsp-detected
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro lsp-id? lsp-id
+---n attempt-to-exceed-max-sequence
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro lsp-id? lsp-id
+---n id-len-mismatch
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name
| | | | | +---ro isis-level? level
| | | | | +---ro interface-name? if:interface-ref
| | | | | +---ro interface-level? level
| | | | | +---ro extended-circuit-id? extended-circuit-id
| | | | | +---ro pdu-field-len? uint8
| | | | | +---ro raw-pdu? binary
+---n max-area-addresses-mismatch
| | | | | +---ro routing-protocol-name? -> /rt:routing/
| | | | | | | | | control-plane-protocols/
| | | | | | | | | control-plane-protocol/name

```

```

+---ro isis-level?                level
+---ro interface-name?            if:interface-ref
+---ro interface-level?          level
+---ro extended-circuit-id?       extended-circuit-id
+---ro max-area-addresses?        uint8
+---ro raw-pdu?                   binary
+---n own-lsp-purge
+---ro routing-protocol-name?     -> /rt:routing/
                                   control-plane-protocols/
                                   control-plane-protocol/name
+---ro isis-level?                level
+---ro interface-name?            if:interface-ref
+---ro interface-level?          level
+---ro extended-circuit-id?       extended-circuit-id
+---ro lsp-id?                    lsp-id
+---n sequence-number-skipped
+---ro routing-protocol-name?     -> /rt:routing/
                                   control-plane-protocols/
                                   control-plane-protocol/name
+---ro isis-level?                level
+---ro interface-name?            if:interface-ref
+---ro interface-level?          level
+---ro extended-circuit-id?       extended-circuit-id
+---ro lsp-id?                    lsp-id
+---n authentication-type-failure
+---ro routing-protocol-name?     -> /rt:routing/
                                   control-plane-protocols/
                                   control-plane-protocol/name
+---ro isis-level?                level
+---ro interface-name?            if:interface-ref
+---ro interface-level?          level
+---ro extended-circuit-id?       extended-circuit-id
+---ro raw-pdu?                   binary
+---n authentication-failure
+---ro routing-protocol-name?     -> /rt:routing/
                                   control-plane-protocols/
                                   control-plane-protocol/name
+---ro isis-level?                level
+---ro interface-name?            if:interface-ref
+---ro interface-level?          level
+---ro extended-circuit-id?       extended-circuit-id
+---ro raw-pdu?                   binary
+---n version-skew
+---ro routing-protocol-name?     -> /rt:routing/
                                   control-plane-protocols/
                                   control-plane-protocol/name
+---ro isis-level?                level
+---ro interface-name?            if:interface-ref

```

```

| +---ro interface-level?          level
| +---ro extended-circuit-id?     extended-circuit-id
| +---ro protocol-version?        uint8
| +---ro raw-pdu?                 binary
+---n area-mismatch
| +---ro routing-protocol-name?   -> /rt:routing/
|                               control-plane-protocols/
|                               control-plane-protocol/name
| +---ro isis-level?              level
| +---ro interface-name?          if:interface-ref
| +---ro interface-level?         level
| +---ro extended-circuit-id?     extended-circuit-id
| +---ro raw-pdu?                 binary
+---n rejected-adjacency
| +---ro routing-protocol-name?   -> /rt:routing/
|                               control-plane-protocols/
|                               control-plane-protocol/name
| +---ro isis-level?              level
| +---ro interface-name?          if:interface-ref
| +---ro interface-level?         level
| +---ro extended-circuit-id?     extended-circuit-id
| +---ro raw-pdu?                 binary
| +---ro reason?                  string
+---n protocols-supported-mismatch
| +---ro routing-protocol-name?   -> /rt:routing/
|                               control-plane-protocols/
|                               control-plane-protocol/name
| +---ro isis-level?              level
| +---ro interface-name?          if:interface-ref
| +---ro interface-level?         level
| +---ro extended-circuit-id?     extended-circuit-id
| +---ro raw-pdu?                 binary
| +---ro protocols*               uint8
+---n lsp-error-detected
| +---ro routing-protocol-name?   -> /rt:routing/
|                               control-plane-protocols/
|                               control-plane-protocol/name
| +---ro isis-level?              level
| +---ro interface-name?          if:interface-ref
| +---ro interface-level?         level
| +---ro extended-circuit-id?     extended-circuit-id
| +---ro lsp-id?                  lsp-id
| +---ro raw-pdu?                 binary
| +---ro error-offset?            uint32
| +---ro tlv-type?                uint8
+---n adjacency-state-change
| +---ro routing-protocol-name?   -> /rt:routing/
|                               control-plane-protocols/

```

```

|
|                                     control-plane-protocol/name
+--ro isis-level?                    level
+--ro interface-name?               if:interface-ref
+--ro interface-level?              level
+--ro extended-circuit-id?          extended-circuit-id
+--ro neighbor?                     string
+--ro neighbor-system-id?           system-id
+--ro state?                         adj-state-type
+--ro reason?                       string
+---n lsp-received
|   +--ro routing-protocol-name?    -> /rt:routing/
|   |                               control-plane-protocols/
|   |                               control-plane-protocol/name
|   +--ro isis-level?               level
|   +--ro interface-name?           if:interface-ref
|   +--ro interface-level?          level
|   +--ro extended-circuit-id?      extended-circuit-id
|   +--ro lsp-id?                   lsp-id
|   +--ro sequence?                 uint32
|   +--ro received-timestamp?       yang:timestamp
|   +--ro neighbor-system-id?       system-id
+---n lsp-generation
|   +--ro routing-protocol-name?    -> /rt:routing/
|   |                               control-plane-protocols/
|   |                               control-plane-protocol/name
|   +--ro isis-level?               level
|   +--ro lsp-id?                   lsp-id
|   +--ro sequence?                 uint32
|   +--ro send-timestamp?           yang:timestamp

```

2.5. Authentication Parameters

The module enables authentication configuration through the IETF key-chain module [RFC8177]. The IS-IS module imports the "ietf-key-chain" module and reuses some groupings to allow global and per-interface configuration of authentication. If global authentication is configured, an implementation SHOULD authenticate PSNPs (Partial Sequence Number Packets), CSNPs (Complete Sequence Number Packets) and LSPs (Link State Packets) with the authentication parameters supplied. The authentication of HELLO PDUs (Protocol Data Units) can be activated on a per-interface basis.

2.6. IGP/LDP synchronization

[RFC5443] defines a mechanism where IGP (Interior Gateway Protocol) needs to be synchronized with LDP (Label Distribution Protocol). An "ldp-igp-sync" feature has been defined in the model to support this functionality. The "mpls/ldp/igp-sync" leaf under "interface" allows

activation of the functionality on a per-interface basis. The "mpls/ldp/igp-sync" container in the global configuration is intentionally empty and is not required for feature activation. The goal of this empty container is to facilitate augmentation with additional parameters, e.g., timers.

2.7. ISO parameters

As the IS-IS protocol is based on the ISO protocol suite, some ISO parameters may be required.

This module augments interface configuration model to support selected ISO configuration parameters.

The clns-mtu can be configured for an interface.

2.8. IP FRR

This YANG module supports LFA (Loop Free Alternates) [RFC5286] and remote LFA [RFC7490] as IP Fast Re-Route (FRR) techniques. The "fast-reroute" container may be augmented by other models to support other IP FRR flavors (MRT as defined in [RFC7812], TI-LFA as defined in [I-D.ietf-rtgwg-segment-routing-ti-lfa], etc.).

The current version of the model supports activation of LFA and remote LFA at the interface-level only. The global "lfa" container is present but kept empty to allow augmentation with vendor-specific properties, e.g., policies.

Remote LFA is considered as an extension of LFA. Remote LFA cannot be enabled if LFA is not enabled.

The "candidate-enable" data leaf designates that an interface can be used as a backup.

2.9. Operational States

Operational state is defined in module in various containers at various levels:

- o system-counters: Provides statistical information about the global system.
- o interface: Provides configuration state information for each interface.
- o adjacencies: Provides state information about current IS-IS adjacencies.

- o `spf-log`: Provides information about SPF events for an IS-IS instance. This SHOULD be implemented as a wrapping buffer.
- o `lsp-log`: Provides information about LSP events for an IS-IS instance (reception of an LSP or modification of a local LSP). This SHOULD be implemented as a wrapping buffer and the implementation MAY optionally log LSP refreshes.
- o `local-rib`: Provides the IS-IS internal routing table.
- o `database`: Provides contents of the current Link State Database.
- o `hostnames`: Provides the system-id to hostname mappings [RFC5301].
- o `fast-reroute`: Provides IP FRR state information.

3. RPC Operations

The "ietf-isis" module defines two RPC operations:

- o `clear-database`: Reset the content of a particular IS-IS database and restart database synchronization with all neighbors.
- o `clear-adjacency`: Restart a particular set of IS-IS adjacencies.

4. Notifications

The "ietf-isis" module defines the following notifications:

`database-overload`: This notification is sent when the IS-IS Node overload condition changes.

`lsp-too-large`: This notification is sent when the system tries to propagate a PDU that is too large.

`if-state-change`: This notification is sent when an interface's state changes.

`corrupted-lsp-detected`: This notification is sent when the IS-IS node discovers that an LSP that was previously stored in the Link State Database, i.e., local memory, has become corrupted.

`attempt-to-exceed-max-sequence`: This notification is sent when the system wraps the 32-bit sequence counter of an LSP.

`id-len-mismatch`: This notification is sent when we receive a PDU with a different value for the System ID length.

max-area-addresses-mismatch: This notification is sent when we receive a PDU with a different value for the Maximum Area Addresses.

own-lsp-purge: This notification is sent when the system receives a PDU with its own system ID and zero age.

sequence-number-skipped: This notification is sent when the system receives a PDU with its own system ID and different contents. The system has to reissue the LSP with a higher sequence number.

authentication-type-failure: This notification is sent when the system receives a PDU with the wrong authentication type field.

authentication-failure: This notification is sent when the system receives a PDU with the wrong authentication information.

version-skew: This notification is sent when the system receives a PDU with a different protocol version number.

area-mismatch: This notification is sent when the system receives a Hello PDU from an IS that does not share any area address.

rejected-adjacency: This notification is sent when the system receives a Hello PDU from an IS but does not establish an adjacency for some reason.

protocols-supported-mismatch: This notification is sent when the system receives a non-pseudonode LSP that has no matching protocol supported.

lsp-error-detected: This notification is sent when the system receives an LSP with a parse error.

adjacency-state-change: This notification is sent when an IS-IS adjacency moves to Up state or to Down state.

lsp-received: This notification is sent when an LSP is received.

lsp-generation: This notification is sent when an LSP is regenerated.

5. Interaction with Other YANG Modules

The "isis" container augments the "/rt:routing/rt:control-plane-protocols/control-plane-protocol" container of the ietf-routing [RFC8349] module with IS-IS-specific parameters.

The "isis" module augments "/if:interfaces/if:interface" defined by [RFC8343] with ISO specific parameters.

The "isis" operational state container augments the "/rt:routing-state/rt:control-plane-protocols/control-plane-protocol" container of the ietf-routing module with IS-IS-specific operational states.

Some IS-IS-specific route attributes are added to route objects in the ietf-routing module by augmenting "/rt:routing-state/rt:ribs/rt:rib/rt:routes/rt:route".

The modules defined in this document uses some groupings from ietf-keychain [RFC8177].

The module reuses types from [RFC6991] and [RFC8294].

To support BFD for fast detection, the module relies on [I-D.ietf-bfd-yang].

6. IS-IS YANG Module

The following RFCs, drafts and external standards are not referenced in the document text but are referenced in the ietf-isis.yang module: [ISO-10589], [RFC1195], [RFC4090], [RFC5029], [RFC5130], [RFC5302], [RFC5305], [RFC5306], [RFC5307], [RFC5308], [RFC5880], [RFC5881], [RFC6119], [RFC6232], [RFC7794], [RFC7981], [RFC8570], [RFC7917], [RFC8405].

```
<CODE BEGINS> file "ietf-isis@2019-10-15.yang"
module ietf-isis {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:yang:ietf-isis";

  prefix isis;

  import ietf-routing {
    prefix "rt";
    reference "RFC 8349 - A YANG Data Model for Routing
              Management (NMDA Version)";
  }

  import ietf-inet-types {
    prefix inet;
    reference "RFC 6991 - Common YANG Data Types";
  }

  import ietf-yang-types {
```



```
    prefix yang;
    reference "RFC 6991 - Common YANG Data Types";
}

import ietf-interfaces {
    prefix "if";
    reference "RFC 8343 - A YANG Data Model for Interface
              Management (NDMA Version)";
}

import ietf-key-chain {
    prefix "key-chain";
    reference "RFC 8177 - YANG Data Model for Key Chains";
}

import ietf-routing-types {
    prefix "rt-types";
    reference "RFC 8294 - Common YANG Data Types for the
              Routing Area";
}

import iana-routing-types {
    prefix "iana-rt-types";
    reference "RFC 8294 - Common YANG Data Types for the
              Routing Area";
}

import ietf-bfd-types {
    prefix "bfd-types";
    reference "RFC YYYY - YANG Data Model for Bidirectional
              Forwarding Detection (BFD)".
}

-- Note to RFC Editor Please replace YYYY with published RFC
   number for draft-ietf-bfd-yang.";

}

organization
    "IETF LSR Working Group";

contact
    "WG Web:    <https://datatracker.ietf.org/group/lsr/>
    WG List:    <mailto:lsr@ietf.org>

    Editor:     Stephane Litkowski
                <mailto:slitkows.ietf@gmail.com>
    Author:     Derek Yeung
                <mailto:derek@arrcus.com>
```

Author: Acee Lindem
<mailto:acee@cisco.com>
Author: Jeffrey Zhang
<mailto:zzhang@juniper.net>
Author: Ladislav Lhotka
<mailto:llhotka@nic.cz>;

description

"This YANG module defines the generic configuration and operational state for the IS-IS protocol common to all vendor implementations. It is intended that the module will be extended by vendors to define vendor-specific IS-IS configuration parameters and policies, for example, route maps or route policies.

This YANG model conforms to the Network Management Datastore Architecture (NMDA) as described in RFC 8242.

Copyright (c) 2018 IETF Trust and the persons identified as authors of the code. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, is permitted pursuant to, and subject to the license terms contained in, the Simplified BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX (<https://www.rfc-editor.org/info/rfcXXXX>); see the RFC itself for full legal notices.

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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.";

```
revision 2019-10-15 {  
  description  
    "Initial revision."  
  reference "RFC XXXX";  
}
```

```
/* Identities */
```

```
identity isis {
  base rt:routing-protocol;
  description "Identity for the IS-IS routing protocol.";
}

identity lsp-log-reason {
  description "Base identity for an LSP change log reason.";
}

identity refresh {
  base lsp-log-reason;
  description
    "Identity used when the LSP log reason is
    a refresh LSP received.";
}

identity content-change {
  base lsp-log-reason;
  description
    "Identity used when the LSP log reason is
    a change in the content of the LSP.";
}

identity frr-protection-method {
  description
    "Base identity for a Fast Reroute protection method.";
}

identity frr-protection-method-lfa {
  base frr-protection-method;
  description "Loop Free Alternate as defined in RFC5286.";
}

identity frr-protection-method-rlfa {
  base frr-protection-method;
  description "Remote Loop Free Alternate as defined in RFC7490.";
}

identity frr-protection-method-rsvpte {
  base frr-protection-method;
  description "RSVP-TE as defined in RFC4090.";
}

identity frr-protection-available-type {
  description "Base identity for Fast Reroute protection types
    provided by an alternate path.";
}

identity frr-protection-available-node-type {
  base frr-protection-available-type;
  description "Node protection is provided by the alternate.";
}
```

```
identity frr-protection-available-link-type {
  base frr-protection-available-type;
  description "Link protection is provided by the alternate.";
}
identity frr-protection-available-srlg-type {
  base frr-protection-available-type;
  description "SRLG protection is provided by the alternate.";
}
identity frr-protection-available-downstream-type {
  base frr-protection-available-type;
  description "The alternate is downstream of node in the path.";
}
identity frr-protection-available-other-type {
  base frr-protection-available-type;
  description "The level of protection is unknown.";
}

identity frr-alternate-type {
  description "Base identity for IP Fast Reroute alternate type.";
}
identity frr-alternate-type-equal-cost {
  base frr-alternate-type;
  description "ECMP alternate.";
}
identity frr-alternate-type-lfa {
  base frr-alternate-type;
  description "LFA alternate.";
}
identity frr-alternate-type-remote-lfa {
  base frr-alternate-type;
  description "Remote LFA alternate.";
}
identity frr-alternate-type-tunnel {
  base frr-alternate-type;
  description "Tunnel based alternate (such as,
    RSVP-TE or GRE).";
}
identity frr-alternate-mrt {
  base frr-alternate-type;
  description "MRT alternate.";
}
identity frr-alternate-tilfa {
  base frr-alternate-type;
  description "TILFA alternate.";
}
identity frr-alternate-other {
  base frr-alternate-type;
  description "Other alternate.";
```

```
}

identity unidirectional-link-delay-subtlv-flag {
    description "Base identity for unidirectional-link-delay
                subTLV flags. Flags are defined in RFC8570.";
}
identity unidirectional-link-delay-subtlv-a-flag {
    base unidirectional-link-delay-subtlv-flag;
    description
        "The A bit represents the Anomalous (A) bit.
        The A bit is set when the measured value of
        this parameter exceeds its configured
        maximum threshold.
        The A bit is cleared when the measured value
        falls below its configured reuse threshold.
        If the A bit is clear,
        the value represents steady-state link performance.";
}
identity min-max-unidirectional-link-delay-subtlv-flag {
    description
        "Base identity for min-max-unidirectional-link-delay
        subTLV flags. Flags are defined in RFC8570.";
}
identity min-max-unidirectional-link-delay-subtlv-a-flag {
    base min-max-unidirectional-link-delay-subtlv-flag;
    description
        "The A bit represents the Anomalous (A) bit.
        The A bit is set when the measured value of
        this parameter exceeds its configured
        maximum threshold.
        The A bit is cleared when the measured value
        falls below its configured reuse threshold.
        If the A bit is clear,
        the value represents steady-state link performance.";
}
identity unidirectional-link-loss-subtlv-flag {
    description "Base identity for unidirectional-link-loss
                subTLV flags. Flags are defined in RFC8570.";
}

identity unidirectional-link-loss-subtlv-a-flag {
    base unidirectional-link-loss-subtlv-flag;
    description
        "The A bit represents the Anomalous (A) bit.
        The A bit is set when the measured value of
        this parameter exceeds its configured
        maximum threshold.
```

```
        The A bit is cleared when the measured value
        falls below its configured reuse threshold.
        If the A bit is clear,
        the value represents steady-state link performance.";
    }
    identity tlv229-flag {
        description "Base identity for TLV229 flags. Flags are defined
            in RFC5120.";
    }
    identity tlv229-overload-flag {
        base tlv229-flag;
        description
            "If set, the originator is overloaded,
            and must be avoided in path calculation.";
    }
    identity tlv229-attached-flag {
        base tlv229-flag;
        description
            "If set, the originator is attached to
            another area using the referred metric.";
    }
    identity router-capability-flag {
        description "Base identity for router capability flags.
            Flags are defined in RFC7981.";
    }
    identity router-capability-flooding-flag {
        base router-capability-flag;
        description
            "Quote from RFC7981: 'If the S bit is set,
            the IS-IS Router CAPABILITY
            TLV MUST be flooded across the entire routing
            domain. If the S bit is clear, the TLV MUST NOT
            be leaked between levels. This bit MUST NOT
            be altered during the TLV leaking'.";
    }
    identity router-capability-down-flag {
        base router-capability-flag;
        description
            "Quote from RFC7981: 'When the IS-IS Router CAPABILITY TLV
            is leaked from level-2 to level-1, the D bit MUST be set.
            Otherwise, this bit MUST be clear. IS-IS Router
            capability TLVs with the D bit set MUST NOT be
            leaked from level-1 to level-2 in to prevent
            TLV looping'.";
    }

    identity lsp-flag {
        description "Base identity for LSP attributes.
```

```

        Attributes are defined in ISO 10589";
    }
    identity lsp-partitioned-flag {
        base lsp-flag;
        description "Originator partition repair supported";
    }
    identity lsp-attached-error-metric-flag {
        base lsp-flag;
        description "Set when originator is attached to
            another area using the error metric.";
    }
    identity lsp-attached-delay-metric-flag {
        base lsp-flag;
        description "Set when originator is attached to
            another area using the delay metric.";
    }
    identity lsp-attached-expense-metric-flag {
        base lsp-flag;
        description "Set when originator is attached to
            another area using the expense metric.";
    }
    identity lsp-attached-default-metric-flag {
        base lsp-flag;
        description "Set when originator is attached to
            another area using the default metric.";
    }
    identity lsp-overload-flag {
        base lsp-flag;
        description
            "If set, the originator is overloaded,
            and must be avoided in path calculation.";
    }
    identity lsp-l1system-flag {
        base lsp-flag;
        description
            "Set when the Intermediate System has an L1 type.";
    }
    identity lsp-l2system-flag {
        base lsp-flag;
        description
            "Set when the Intermediate System has an L2 type.";
    }
}

/* Feature definitions */

feature osi-interface {
    description "Support of OSI specific parameters on an
```

```
        interface.";
    }
    feature poi-tlv {
        description "Support of Purge Originator Identification.";
        reference "RFC 6232 - Purge Originator Identification TLV
            for IS-IS";
    }
    feature ietf-spf-delay {
        description
            "Support for IETF SPF delay algorithm.";
        reference "RFC 8405 - SPF Back-off algorithm for link
            state IGP";
    }
    feature bfd {
        description
            "Support for BFD detection of IS-IS neighbor reachability.";
        reference "RFC 5880 - Bidirectional Forwarding Detection (BFD)
            RFC 5881 - Bidirectional Forwarding Detection
            (BFD) for IPv4 and IPv6 (Single Hop)";
    }
    feature key-chain {
        description
            "Support of keychain for authentication.";
        reference "RFC8177 - YANG Data Model for Key Chains";
    }
    feature node-flag {
        description
            "Support for node-flag for IS-IS prefixes.";
        reference "RFC7794 - IS-IS Prefix Attributes for
            Extended IP and IPv6 Reachability";
    }
    feature node-tag {
        description
            "Support for node admin tag for IS-IS routing instances.";
        reference "RFC7917 - Advertising Node Administrative Tags
            in IS-IS";
    }
    feature ldp-igp-sync {
        description
            "Support for LDP IGP synchronization.";
        reference "RFC5443 - LDP IGP Synchronization.";
    }
    feature fast-reroute {
        description
            "Support for IP Fast Reroute (IP-FRR).";
    }
    feature nsr {
        description
```



```
    "Support for Non-Stop-Routing (NSR). The IS-IS NSR feature
      allows a router with redundant control-plane capability
      (e.g., dual Route-Processor (RP) cards) to maintain its
      state and adjacencies during planned and unplanned
      IS-IS instance restarts. It differs from graceful-restart
      or Non-Stop Forwarding (NSF) in that no protocol signaling
      or assistance from adjacent IS-IS neighbors is required to
      recover control-plane state.";
  }
  feature lfa {
    description
      "Support for Loop-Free Alternates (LFAs).";
    reference "RFC5286 - Basic Specification of IP Fast-Reroute:
      Loop-free Alternates";
  }
  feature remote-lfa {
    description
      "Support for Remote Loop-Free Alternates (R-LFAs).";
    reference "RFC7490 - Remote Loop-Free Alternate Fast Reroute";
  }

  feature overload-max-metric {
    description
      "Support of overload by setting all links to max metric.
      In IS-IS, the overload bit is usually used to signal that
      a node cannot be used as a transit. The overload-max-metric
      feature brings a similar behavior leveraging on setting all
      the link metrics to MAX_METRIC.";
  }
  feature prefix-tag {
    description
      "Support for 32-bit prefix tags";
    reference "RFC5130 - A Policy Control Mechanism in
      IS-IS Using Administrative Tags";
  }
  feature prefix-tag64 {
    description
      "Support for 64-bit prefix tags";
    reference "RFC5130 - A Policy Control Mechanism in
      IS-IS Using Administrative Tags";
  }
  feature auto-cost {
    description
      "Support for IS-IS interface metric computation
      according to a reference bandwidth.";
  }

  feature te-rid {
```

```
    description
        "Traffic-Engineering Router-ID.";
    reference "RFC5305 - IS-IS Extensions for Traffic Engineering
        RFC6119 - IPv6 Traffic Engineering in IS-IS";
}
feature max-ecmp {
    description
        "Setting maximum number of ECMP paths.";
}
feature multi-topology {
    description
        "Support for Multiple-Topology Routing (MTR).";
    reference "RFC5120 - M-IS-IS: Multi Topology Routing in IS-IS";
}
feature nlpid-control {
    description
        "Support for the advertisement
        of a Network Layer Protocol Identifier within IS-IS
        configuration.";
}
feature graceful-restart {
    description
        "IS-IS Graceful restart support.";
    reference "RFC5306 - Restart Signaling in IS-IS";
}

feature lsp-refresh {
    description
        "Configuration of LSP refresh interval.";
}

feature maximum-area-addresses {
    description
        "Support for maximum-area-addresses configuration.";
}

feature admin-control {
    description
        "Administrative control of the protocol state.";
}

/* Type definitions */

typedef circuit-id {
    type uint8;
    description
        "This type defines the circuit ID
        associated with an interface.";
```

```
}

typedef extended-circuit-id {
    type uint32;
    description
        "This type defines the extended circuit ID
        associated with an interface.";
}

typedef interface-type {
    type enumeration {
        enum broadcast {
            description
                "Broadcast interface type.";
        }
        enum point-to-point {
            description
                "Point-to-point interface type.";
        }
    }
    description
        "This type defines the type of adjacency
        to be established for the interface.
        The interface-type determines the type
        of hello message that is used.";
}

typedef level {
    type enumeration {
        enum "level-1" {
            description
                "This enum indicates L1-only capability.";
        }
        enum "level-2" {
            description
                "This enum indicates L2-only capability.";
        }
        enum "level-all" {
            description
                "This enum indicates capability for both levels.";
        }
    }
    default "level-all";
    description
        "This type defines IS-IS level of an object.";
}
```

```
typedef adj-state-type {
    type enumeration {
        enum "up" {
            description
                "State indicates the adjacency is established.";
        }
        enum "down" {
            description
                "State indicates the adjacency is NOT established.";
        }
        enum "init" {
            description
                "State indicates the adjacency is establishing.";
        }
        enum "failed" {
            description
                "State indicates the adjacency is failed.";
        }
    }
    description
        "This type defines states of an adjacency";
}

typedef if-state-type {
    type enumeration {
        enum "up" {
            description "Up state.";
        }
        enum "down" {
            description "Down state";
        }
    }
    description
        "This type defines the state of an interface";
}

typedef level-number {
    type uint8 {
        range "1 .. 2";
    }
    description
        "This type defines the current IS-IS level.";
}

typedef lsp-id {
    type string {
        pattern
```

```
        '[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.[0-9A-Fa-f]'
        +'{4}\.[0-9][0-9]-[0-9][0-9]';
    }
    description
        "This type defines the IS-IS LSP ID format using a
        pattern. An example LSP ID is 0143.0438.AEF0.02-01";
}

typedef area-address {
    type string {
        pattern '[0-9A-Fa-f]{2}(\.[0-9A-Fa-f]{4}){0,6}';
    }
    description
        "This type defines the area address format.";
}

typedef snpa {
    type string {
        length "0 .. 20";
    }
    description
        "This type defines the Subnetwork Point
        of Attachment (SNPA) format.
        The SNPA should be encoded according to the rules
        specified for the particular type of subnetwork
        being used. As an example, for an ethernet subnetwork,
        the SNPA is encoded as a MAC address, such as,
        '00aa.bbcc.ddee'.";
}

typedef system-id {
    type string {
        pattern
            '[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}';
    }
    description
        "This type defines IS-IS system-id using pattern,
        An example system-id is 0143.0438.AEF0";
}

typedef extended-system-id {
    type string {
        pattern
            '[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.[0-9A-Fa-f]{4}\.'
            +'[0-9][0-9]';
    }
    description
        "This type defines IS-IS system-id using pattern. The extended
        system-id contains the pseudonode number in addition to the
```

```
        system-id.  
        An example system-id is 0143.0438.AEF0.00";  
    }  
  
    typedef wide-metric {  
        type uint32 {  
            range "0 .. 16777215";  
        }  
        description  
            "This type defines wide style format of IS-IS metric.";  
    }  
  
    typedef std-metric {  
        type uint8 {  
            range "0 .. 63";  
        }  
        description  
            "This type defines old style format of IS-IS metric.";  
    }  
  
    typedef mesh-group-state {  
        type enumeration {  
            enum "mesh-inactive" {  
                description  
                    "Interface is not part of a mesh group.";  
            }  
            enum "mesh-set" {  
                description  
                    "Interface is part of a mesh group.";  
            }  
            enum "mesh-blocked" {  
                description  
                    "LSPs must not be flooded over this interface.";  
            }  
        }  
        description  
            "This type describes mesh group state of an interface";  
    }  
  
    /* Grouping for notifications */  
  
    grouping notification-instance-hdr {  
        description  
            "Instance specific IS-IS notification data grouping";  
        leaf routing-protocol-name {  
            type leafref {  
                path "/rt:routing/rt:control-plane-protocols/"  
                    + "rt:control-plane-protocol/rt:name";  
            }  
        }  
    }
```

```
    }
    description "Name of the IS-IS instance.";
  }
  leaf isis-level {
    type level;
    description "IS-IS level of the instance.";
  }
}

grouping notification-interface-hdr {
  description
    "Interface specific IS-IS notification data grouping";
  leaf interface-name {
    type if:interface-ref;
    description "IS-IS interface name";
  }
  leaf interface-level {
    type level;
    description "IS-IS level of the interface.";
  }
  leaf extended-circuit-id {
    type extended-circuit-id;
    description "Extended circuit-id of the interface.";
  }
}

/* Groupings for IP Fast Reroute */

grouping instance-fast-reroute-config {
  description
    "This group defines global configuration of IP
    Fast ReRoute (FRR).";
  container fast-reroute {
    if-feature fast-reroute;
    description
      "This container may be augmented with global
      parameters for IP-FRR.";
    container lfa {
      if-feature lfa;
      description
        "This container may be augmented with
        global parameters for Loop-Free Alternatives (LFA).
        Container creation has no effect on LFA activation.";
    }
  }
}
```

```
grouping interface-lfa-config {
  leaf candidate-enable {
    type boolean;
    default "true";
    description
      "Enable the interface to be used as backup.";
  }
  leaf enable {
    type boolean;
    default false;
    description
      "Activates LFA - Per-prefix LFA computation
       is assumed.";
  }
  container remote-lfa {
    if-feature remote-lfa;
    leaf enable {
      type boolean;
      default false;
      description
        "Activates Remote LFA (R-LFA).";
    }
    description
      "Remote LFA configuration.";
  }
  description "Grouping for LFA interface configuration";
}
grouping interface-fast-reroute-config {
  description
    "This group defines interface configuration of IP-FRR.";
  container fast-reroute {
    if-feature fast-reroute;
    container lfa {
      if-feature lfa;
      uses interface-lfa-config;
      container level-1 {
        uses interface-lfa-config;
        description
          "LFA level 1 config";
      }
      container level-2 {
        uses interface-lfa-config;
        description
          "LFA level 2 config";
      }
    }
    description
      "LFA configuration.";
  }
}
```



```
        description
            "Interface IP Fast-reroute configuration.";
    }
}
grouping instance-fast-reroute-state {
    description "IPFRR state data grouping";
    container protected-routes {
        config false;
        list address-family-stats {
            key "address-family prefix alternate";

            leaf address-family {
                type iana-rt-types:address-family;
                description
                    "Address-family";
            }
            leaf prefix {
                type inet:ip-prefix;
                description
                    "Protected prefix.";
            }
            leaf alternate {
                type inet:ip-address;
                description
                    "Alternate next hop for the prefix.";
            }
            leaf alternate-type {
                type identityref {
                    base frr-alternate-type;
                }
                description
                    "Type of alternate.";
            }
            leaf best {
                type boolean;
                description
                    "Is set when the alternate is the preferred one,
                     is clear otherwise.";
            }
            leaf non-best-reason {
                type string {
                    length "1..255";
                }
                description
                    "Information field to describe why the alternate
                     is not best. The length should be limited to 255
                     unicode characters. The expected format is a single
                     line text.";
```

```
}
container protection-available {
  leaf-list protection-types {
    type identityref {
      base frr-protection-available-type;
    }
    description "This list contains a set of protection
                  types defined as identities.
                  An identity must be added for each type of
                  protection provided by the alternate.
                  As an example, if an alternate provides
                  SRLG, node and link protection, three
                  identities must be added in this list:
                  one for SRLG protection, one for node
                  protection, one for link protection.";
  }
  description "Protection types provided by the alternate.";
}
leaf alternate-metric1 {
  type uint32;
  description
    "Metric from Point of Local Repair (PLR) to
     destination through the alternate path.";
}
leaf alternate-metric2 {
  type uint32;
  description
    "Metric from PLR to the alternate node";
}
leaf alternate-metric3 {
  type uint32;
  description
    "Metric from alternate node to the destination";
}
description
  "Per-AF protected prefix statistics.";
}
description
  "List of prefixes that are protected.";
}

container unprotected-routes {
  config false;
  list prefixes {
    key "address-family prefix";

    leaf address-family {
      type iana-rt-types:address-family;
    }
  }
}
```

```
        description "Address-family";
    }
    leaf prefix {
        type inet:ip-prefix;
        description "Unprotected prefix.";
    }
    description
        "Per-AF unprotected prefix statistics.";
}
description
    "List of prefixes that are not protected.";
}

list protection-statistics {
    key frr-protection-method;
    config false;
    leaf frr-protection-method {
        type identityref {
            base frr-protection-method;
        }
        description "Protection method used.";
    }
}
list address-family-stats {
    key address-family;

    leaf address-family {
        type iana-rt-types:address-family;

        description "Address-family";
    }
    leaf total-routes {
        type yang:gauge32;
        description "Total prefixes.";
    }
    leaf unprotected-routes {
        type yang:gauge32;
        description
            "Total prefixes that are not protected.";
    }
    leaf protected-routes {
        type yang:gauge32;
        description
            "Total prefixes that are protected.";
    }
    leaf link-protected-routes {
        type yang:gauge32;
        description
            "Total prefixes that are link protected.";
    }
}
```

```
    }
    leaf node-protected-routes {
        type yang:gauge32;
        description
            "Total prefixes that are node protected.";
    }
    description
        "Per-AF protected prefix statistics.";
}

description "Global protection statistics.";
}
}

/* Route table and local RIB groupings */

grouping local-rib {
    description "Local-rib - RIB for Routes computed by the local
        IS-IS routing instance.";
    container local-rib {
        config false;
        description "Local-rib.";
        list route {
            key "prefix";
            description "Routes";
            leaf prefix {
                type inet:ip-prefix;
                description "Destination prefix.";
            }
            container next-hops {
                description "Next hops for the route.";
                list next-hop {
                    key "next-hop";
                    description "List of next hops for the route";
                    leaf outgoing-interface {
                        type if:interface-ref;
                        description
                            "Name of the outgoing interface.";
                    }
                    leaf next-hop {
                        type inet:ip-address;
                        description "Next hop address.";
                    }
                }
            }
        }
        leaf metric {
            type uint32;
            description "Metric for this route.";
        }
    }
}
```

```
    }
    leaf level {
        type level-number;
        description "Level number for this route.";
    }
    leaf route-tag {
        type uint32;
        description "Route tag for this route.";
    }
}
}
}

grouping route-content {
    description
        "IS-IS protocol-specific route properties grouping.";
    leaf metric {
        type uint32;
        description "IS-IS metric of a route.";
    }
    leaf-list tag {
        type uint64;
        description
            "List of tags associated with the route.
             This list provides a consolidated view of both
             32-bit and 64-bit tags (RFC5130) available for the prefix.";
    }
    leaf route-type {
        type enumeration {
            enum l2-intra-area {
                description "Level 2 internal route. As per RFC5302,
                             the prefix is directly connected to the
                             advertising router. It cannot be
                             distinguished from an L1->L2 inter-area
                             route.";
            }
            enum l1-intra-area {
                description "Level 1 internal route. As per RFC5302,
                             the prefix is directly connected to the
                             advertising router.";
            }
            enum l2-external {
                description "Level 2 external route. As per RFC5302,
                             such a route is learned from other IGPs.
                             It cannot be distinguished from an L1->L2
                             inter-area external route.";
            }
            enum l1-external {
```

```
        description "Level 1 external route. As per RFC5302,
                    such a route is learned from other IGPs.";
    }
    enum l1-inter-area {
        description "These prefixes are learned via L2 routing.";
    }
    enum l1-inter-area-external {
        description "These prefixes are learned via L2 routing
                    towards an l2-external route.";
    }
}
description "IS-IS route type.";
}
```

```
/* Grouping definitions for configuration and ops state */
```

```
grouping adjacency-state {
    container adjacencies {
        config false;
        list adjacency {
            leaf neighbor-sys-type {
                type level;
                description
                    "Level capability of neighboring system";
            }
            leaf neighbor-sysid {
                type system-id;
                description
                    "The system-id of the neighbor";
            }
            leaf neighbor-extended-circuit-id {
                type extended-circuit-id;
                description
                    "Circuit ID of the neighbor";
            }
            leaf neighbor-snpa {
                type snpa;
                description
                    "SNPA of the neighbor";
            }
            leaf usage {
                type level;
                description
                    "Define the level(s) activated for the adjacency.
                     On a p2p link this might be level 1 and 2,
```

```
        but on a LAN, the usage will be level 1
        between neighbors at level 1 or level 2 between
        neighbors at level 2.";
    }
    leaf hold-timer {
        type rt-types:timer-value-seconds16;
        units seconds;
        description
            "The holding time in seconds for this
            adjacency. This value is based on
            received hello PDUs and the elapsed
            time since receipt.";
    }
    leaf neighbor-priority {
        type uint8 {
            range "0 .. 127";
        }
        description
            "Priority of the neighboring IS for becoming
            the DIS.";
    }
    leaf lastuptime {
        type yang:timestamp;
        description
            "When the adjacency most recently entered
            state 'up', measured in hundredths of a
            second since the last reinitialization of
            the network management subsystem.
            The value is 0 if the adjacency has never
            been in state 'up'.";
    }
    leaf state {
        type adj-state-type;
        description
            "This leaf describes the state of the interface.";
    }
    description
        "List of operational adjacencies.";
}
description
    "This container lists the adjacencies of
    the local node.";
}
description
    "Adjacency state";
}
```

```
grouping admin-control {
  leaf enable {
    if-feature admin-control;
    type boolean;
    default "true";
    description
      "Enable/Disable the protocol.";
  }
  description
    "Grouping for admin control.";
}

grouping ietf-spf-delay {
  leaf initial-delay {
    type rt-types:timer-value-milliseconds;
    units msec;
    description
      "Delay used while in QUIET state (milliseconds).";
  }
  leaf short-delay {
    type rt-types:timer-value-milliseconds;
    units msec;
    description
      "Delay used while in SHORT_WAIT state (milliseconds).";
  }
  leaf long-delay {
    type rt-types:timer-value-milliseconds;
    units msec;
    description
      "Delay used while in LONG_WAIT state (milliseconds).";
  }

  leaf hold-down {
    type rt-types:timer-value-milliseconds;
    units msec;
    description
      "Timer used to consider an IGP stability period
        (milliseconds).";
  }
  leaf time-to-learn {
    type rt-types:timer-value-milliseconds;
    units msec;
    description
      "Duration used to learn all the IGP events
        related to a single component failure (milliseconds).";
  }
  leaf current-state {
    type enumeration {
```



```
        enum "quiet" {
            description "QUIET state";
        }
        enum "short-wait" {
            description "SHORT_WAIT state";
        }
        enum "long-wait" {
            description "LONG_WAIT state";
        }
    }
    config false;
    description
        "Current SPF back-off algorithm state.";
}
leaf remaining-time-to-learn {
    type rt-types:timer-value-milliseconds;
    units "msec";
    config false;
    description
        "Remaining time until time-to-learn timer fires.";
}
leaf remaining-hold-down {
    type rt-types:timer-value-milliseconds;
    units "msec";
    config false;
    description
        "Remaining time until hold-down timer fires.";
}
leaf last-event-received {
    type yang:timestamp;
    config false;
    description
        "Time of last IGP event received";
}
leaf next-spf-time {
    type yang:timestamp;
    config false;
    description
        "Time when next SPF has been scheduled.";
}
leaf last-spf-time {
    type yang:timestamp;
    config false;
    description
        "Time of last SPF computation.";
}
description
    "Grouping for IETF SPF delay configuration and state.";
```

```
}

grouping node-tag-config {
  description
    "IS-IS node tag config state.";
  container node-tags {
    if-feature node-tag;
    list node-tag {
      key tag;
      leaf tag {
        type uint32;
        description
          "Node tag value.";
      }
      description
        "List of tags.";
    }
    description
      "Container for node admin tags.";
  }
}

grouping authentication-global-cfg {
  choice authentication-type {
    case key-chain {
      if-feature key-chain;
      leaf key-chain {
        type key-chain:key-chain-ref;
        description
          "Reference to a key-chain.";
      }
    }
    case password {
      leaf key {
        type string;
        description
          "This leaf specifies the authentication key. The
          length of the key may be dependent on the
          cryptographic algorithm.";
      }
      leaf crypto-algorithm {
        type identityref {
          base key-chain:crypto-algorithm;
        }
        description
          "Cryptographic algorithm associated with key.";
      }
    }
  }
}
```

```
    }
  }
  description "Choice of authentication.";
}
description "Grouping for global authentication config.";
}

grouping metric-type-global-cfg {
  leaf value {
    type enumeration {
      enum wide-only {
        description
          "Advertise new metric style only (RFC5305)";
      }
      enum old-only {
        description
          "Advertise old metric style only (RFC1195)";
      }
      enum both {
        description "Advertise both metric styles";
      }
    }
  }
  description
    "Type of metric to be generated:
    - wide-only means only new metric style
      is generated,
    - old-only means that only old-style metric
      is generated,
    - both means that both are advertised.
    This leaf is only affecting IPv4 metrics.";
}
description
  "Grouping for global metric style config.";
}

grouping metric-type-global-cfg-with-default {
  leaf value {
    type enumeration {
      enum wide-only {
        description
          "Advertise new metric style only (RFC5305)";
      }
      enum old-only {
        description
          "Advertise old metric style only (RFC1195)";
      }
      enum both {
        description "Advertise both metric styles";
      }
    }
  }
}
```

```
    }
  }
  default wide-only;
  description
    "Type of metric to be generated:
    - wide-only means only new metric style
      is generated,
    - old-only means that only old-style metric
      is generated,
    - both means that both are advertised.
    This leaf is only affecting IPv4 metrics.";
}
description
  "Grouping for global metric style config.";
}

grouping default-metric-global-cfg {
  leaf value {
    type wide-metric;
    description "Value of the metric";
  }
  description
    "Global default metric config grouping.";
}

grouping default-metric-global-cfg-with-default {
  leaf value {
    type wide-metric;
    default "10";
    description "Value of the metric";
  }
  description
    "Global default metric config grouping.";
}

grouping overload-global-cfg {
  leaf status {
    type boolean;
    default false;
    description
      "This leaf specifies the overload status.";
  }
  description "Grouping for overload bit config.";
}

grouping overload-max-metric-global-cfg {
  leaf timeout {
    type rt-types:timer-value-seconds16;
```

```
        units "seconds";
        description
            "Timeout (in seconds) of the overload condition.";
    }
    description
        "Overload maximum metric configuration grouping";
}

grouping route-preference-global-cfg {
    choice granularity {
        case detail {
            leaf internal {
                type uint8;
                description
                    "Protocol preference for internal routes.";
            }
            leaf external {
                type uint8;
                description
                    "Protocol preference for external routes.";
            }
        }
        case coarse {
            leaf default {
                type uint8;
                description
                    "Protocol preference for all IS-IS routes.";
            }
        }
    }
    description
        "Choice for implementation of route preference.";
}
description
    "Global route preference grouping";
}

grouping hello-authentication-cfg {
    choice authentication-type {
        case key-chain {
            if-feature key-chain;
            leaf key-chain {
                type key-chain:key-chain-ref;
                description "Reference to a key-chain.";
            }
        }
        case password {
            leaf key {
                type string;
            }
        }
    }
}
```

```
        description "Authentication key specification - The
                    length of the key may be dependent on the
                    cryptographic algorithm.";
    }
    leaf crypto-algorithm {
        type identityref {
            base key-chain:crypto-algorithm;
        }
        description
            "Cryptographic algorithm associated with key.";
    }
    }
    description "Choice of authentication.";
}
description "Grouping for hello authentication.";
}

grouping hello-interval-cfg {
    leaf value {
        type rt-types:timer-value-seconds16;
        units "seconds";
        description
            "Interval (in seconds) between successive hello
            messages.";
    }

    description "Interval between hello messages.";
}
grouping hello-interval-cfg-with-default {
    leaf value {
        type rt-types:timer-value-seconds16;
        units "seconds";
        default 10;
        description
            "Interval (in seconds) between successive hello
            messages.";
    }

    description "Interval between hello messages.";
}

grouping hello-multiplier-cfg {
    leaf value {
        type uint16;
        description
            "Number of missed hello messages prior to
            declaring the adjacency down.";
    }
}
```

```
        description
            "Number of missed hello messages prior to
            adjacency down grouping.";
    }
    grouping hello-multiplier-cfg-with-default {
        leaf value {
            type uint16;
            default 3;
            description
                "Number of missed hello messages prior to
                declaring the adjacency down.";
        }
        description
            "Number of missed hello messages prior to
            adjacency down grouping.";
    }

    grouping priority-cfg {
        leaf value {
            type uint8 {
                range "0 .. 127";
            }
            description
                "Priority of interface for DIS election.";
        }

        description "Interface DIS election priority grouping";
    }
    grouping priority-cfg-with-default {
        leaf value {
            type uint8 {
                range "0 .. 127";
            }
            default 64;
            description
                "Priority of interface for DIS election.";
        }

        description "Interface DIS election priority grouping";
    }

    grouping metric-cfg {
        leaf value {
            type wide-metric;
            description "Metric value.";
        }
        description "Interface metric grouping";
    }
}
```

```
grouping metric-cfg-with-default {
  leaf value {
    type wide-metric;
    default "10";
    description "Metric value.";
  }
  description "Interface metric grouping";
}

grouping metric-parameters {
  container metric-type {
    uses metric-type-global-cfg-with-default;
    container level-1 {
      uses metric-type-global-cfg;
      description "level-1 specific configuration";
    }
    container level-2 {
      uses metric-type-global-cfg;
      description "level-2 specific configuration";
    }
    description "Metric style global configuration";
  }

  container default-metric {
    uses default-metric-global-cfg-with-default;
    container level-1 {
      uses default-metric-global-cfg;
      description "level-1 specific configuration";
    }
    container level-2 {
      uses default-metric-global-cfg;
      description "level-2 specific configuration";
    }
    description "Default metric global configuration";
  }

  container auto-cost {
    if-feature auto-cost;
    description
      "Interface Auto-cost configuration state.";
    leaf enable {
      type boolean;
      description
        "Enable/Disable interface auto-cost.";
    }
    leaf reference-bandwidth {
      when "../enable = 'true'" {
        description "Only when auto cost is enabled";
      }
    }
  }
}
```



```
    }
    type uint32 {
        range "1..4294967";
    }
    units Mbits;
    description
        "Configure reference bandwidth used to automatically
        determine interface cost (Mbits). The cost is the
        reference bandwidth divided by the interface speed
        with 1 being the minimum cost.";
    }
}

description "Grouping for global metric parameters.";
}

grouping high-availability-parameters {
    container graceful-restart {
        if-feature graceful-restart;
        leaf enable {
            type boolean;
            default false;
            description "Enable graceful restart.";
        }
        leaf restart-interval {
            type rt-types:timer-value-seconds16;
            units "seconds";
            description
                "Interval (in seconds) to attempt graceful restart prior
                to failure.";
        }
        leaf helper-enable {
            type boolean;
            default "true";
            description
                "Enable local IS-IS router as graceful restart helper.";
        }
        description "Graceful-Restart Configuration.";
    }
    container nsr {
        if-feature nsr;
        description "Non-Stop Routing (NSR) configuration.";
        leaf enable {
            type boolean;
            default false;
            description "Enable/Disable Non-Stop Routing (NSR).";
        }
    }
}
```

```
    description "Grouping for High Availability parameters.";
}

grouping authentication-parameters {
    container authentication {
        uses authentication-global-cfg;

        container level-1 {
            uses authentication-global-cfg;
            description "level-1 specific configuration";
        }
        container level-2 {
            uses authentication-global-cfg;
            description "level-2 specific configuration";
        }
        description "Authentication global configuration for
            both LSPs and SNPs.";
    }
    description "Grouping for authentication parameters";
}

grouping address-family-parameters {
    container address-families {
        if-feature nlpid-control;
        list address-family-list {
            key address-family;
            leaf address-family {
                type iana-rt-types:address-family;
                description "Address-family";
            }
            leaf enable {
                type boolean;
                description "Activate the address family.";
            }
            description
                "List of address families and whether or not they
                are activated.";
        }
        description "Address Family configuration";
    }
    description "Grouping for address family parameters.";
}

grouping mpls-parameters {
    container mpls {
        container te-rid {
            if-feature te-rid;
            description
                "Stable ISIS Router IP Address used for Traffic
```

```
        Engineering";
    leaf ipv4-router-id {
        type inet:ipv4-address;
        description
            "Router ID value that would be used in TLV 134.";
    }
    leaf ipv6-router-id {
        type inet:ipv6-address;
        description
            "Router ID value that would be used in TLV 140.";
    }
}
container ldp {
    container igp-sync {
        if-feature ldp-igp-sync;
        description
            "This container may be augmented with global
            parameters for igp-ldp-sync.";
    }
    description "LDP configuration.";
}
description "MPLS configuration";
}
description "Grouping for MPLS global parameters.";
}

grouping lsp-parameters {
    leaf lsp-mtu {
        type uint16;
        units "bytes";
        default 1492;
        description
            "Maximum size of an LSP PDU in bytes.";
    }
    leaf lsp-lifetime {
        type uint16 {
            range "1..65535";
        }
        units "seconds";
        description
            "Lifetime of the router's LSPs in seconds.";
    }
    leaf lsp-refresh {
        if-feature lsp-refresh;
        type rt-types:timer-value-seconds16;
        units "seconds";
        description
            "Refresh interval of the router's LSPs in seconds.";
    }
}
```

```
    }
    leaf poi-tlv {
        if-feature poi-tlv;
        type boolean;
        default false;
        description
            "Enable advertisement of IS-IS Purge Originator
             Identification TLV.";
    }
    description "Grouping for LSP global parameters.";
}
grouping spf-parameters {
    container spf-control {
        leaf paths {
            if-feature max-ecmp;
            type uint16 {
                range "1..65535";
            }
            description
                "Maximum number of Equal-Cost Multi-Path (ECMP) paths.";
        }
        container ietf-spf-delay {
            if-feature ietf-spf-delay;
            uses ietf-spf-delay;
            description "IETF SPF delay algorithm configuration.";
        }
        description
            "SPF calculation control.";
    }
    description "Grouping for SPF global parameters.";
}
grouping instance-config {
    description "IS-IS global configuration grouping";

    uses admin-control;

    leaf level-type {
        type level;
        default "level-all";
        description
            "Level of an IS-IS node - can be level-1,
             level-2 or level-all.";
    }

    leaf system-id {
        type system-id;
        description "system-id of the node.";
    }
}
```

```
leaf maximum-area-addresses {
  if-feature maximum-area-addresses;
  type uint8;
  default 3;
  description "Maximum areas supported.";
}

leaf-list area-address {
  type area-address;
  description
    "List of areas supported by the protocol instance.";
}

uses lsp-parameters;
uses high-availability-parameters;
uses node-tag-config;
uses metric-parameters;
uses authentication-parameters;
uses address-family-parameters;
uses mpls-parameters;
uses spf-parameters;
uses instance-fast-reroute-config;

container preference {
  uses route-preference-global-cfg;
  description "Router preference configuration for IS-IS
    protocol instance route installation";
}

container overload {
  uses overload-global-cfg;
  description "Router protocol instance overload state
    configuration";
}

container overload-max-metric {
  if-feature overload-max-metric;
  uses overload-max-metric-global-cfg;
  description
    "Router protocol instance overload maximum
    metric advertisement configuration.";
}
}

grouping instance-state {
  description
    "IS-IS instance operational state.";
  uses spf-log;
}
```

```
    uses lsp-log;
    uses hostname-db;
    uses lsdb;
    uses local-rib;
    uses system-counters;
    uses instance-fast-reroute-state;
    leaf discontinuity-time {
        type yang:date-and-time;
        description
            "The time of the most recent occasion at which any one
            or more of this IS-IS instance's counters suffered a
            discontinuity. If no such discontinuities have occurred
            since the IS-IS instance was last re-initialized, then
            this node contains the time the IS-IS instance was
            re-initialized which normally occurs when it was
            created.";
    }
}

grouping multi-topology-config {
    description "Per-topology configuration";
    container default-metric {
        uses default-metric-global-cfg;
        container level-1 {
            uses default-metric-global-cfg;
            description "level-1 specific configuration";
        }
        container level-2 {
            uses default-metric-global-cfg;
            description "level-2 specific configuration";
        }
        description "Default metric per-topology configuration";
    }
    uses node-tag-config;
}

grouping interface-config {
    description "Interface configuration grouping";

    uses admin-control;

    leaf level-type {
        type level;
        default "level-all";
        description "IS-IS level of the interface.";
    }
    leaf lsp-pacing-interval {
        type rt-types:timer-value-milliseconds;
    }
}
```

```
    units "milliseconds";
    default 33;
    description
        "Interval (in milli-seconds) between LSP
        transmissions.";
}
leaf lsp-retransmit-interval {
    type rt-types:timer-value-seconds16;
    units "seconds";
    description
        "Interval (in seconds) between LSP
        retransmissions.";
}
leaf passive {
    type boolean;
    default "false";
    description
        "Indicates whether the interface is in passive mode (IS-IS
        not running but network is advertised).";
}
leaf csnp-interval {
    type rt-types:timer-value-seconds16;
    units "seconds";
    default 10;
    description
        "Interval (in seconds) between CSNP messages.";
}
container hello-padding {
    leaf enable {
        type boolean;
        default "true";
        description
            "IS-IS Hello-padding activation - enabled by default.";
    }
    description "IS-IS hello padding configuration.";
}
leaf mesh-group-enable {
    type mesh-group-state;
    description "IS-IS interface mesh-group state";
}
leaf mesh-group {
    when "../mesh-group-enable = 'mesh-set'" {
        description
            "Only valid when mesh-group-enable equals mesh-set";
    }
    type uint8;
    description "IS-IS interface mesh-group ID.";
}
```

```
leaf interface-type {
    type interface-type;
    default "broadcast";
    description
        "Type of adjacency to be established for the interface. This
        dictates the type of hello messages that are used.";
}

leaf-list tag {
    if-feature prefix-tag;
    type uint32;
    description
        "List of tags associated with the interface.";
}

leaf-list tag64 {
    if-feature prefix-tag64;
    type uint64;
    description
        "List of 64-bit tags associated with the interface.";
}

leaf node-flag {
    if-feature node-flag;
    type boolean;
    default false;
    description
        "Set prefix as a node representative prefix.";
}

container hello-authentication {
    uses hello-authentication-cfg;
    container level-1 {
        uses hello-authentication-cfg;
        description "level-1 specific configuration";
    }
    container level-2 {
        uses hello-authentication-cfg;
        description "level-2 specific configuration";
    }
    description
        "Authentication type to be used in hello messages.";
}

container hello-interval {
    uses hello-interval-cfg-with-default;
    container level-1 {
        uses hello-interval-cfg;
        description "level-1 specific configuration";
    }
    container level-2 {
        uses hello-interval-cfg;
    }
}
```



```
        description "level-2 specific configuration";
    }
    description "Interval between hello messages.";
}
container hello-multiplier {
    uses hello-multiplier-cfg-with-default;
    container level-1 {
        uses hello-multiplier-cfg;
        description "level-1 specific configuration";
    }
    container level-2 {
        uses hello-multiplier-cfg;
        description "level-2 specific configuration";
    }
    description "Hello multiplier configuration.";
}
container priority {
    must '../interface-type = "broadcast"' {
        error-message
            "Priority only applies to broadcast interfaces.";
        description "Check for broadcast interface.";
    }
    uses priority-cfg-with-default;
    container level-1 {
        uses priority-cfg;
        description "level-1 specific configuration";
    }
    container level-2 {
        uses priority-cfg;
        description "level-2 specific configuration";
    }
    description "Priority for DIS election.";
}
container metric {
    uses metric-cfg-with-default;
    container level-1 {
        uses metric-cfg;
        description "level-1 specific configuration";
    }
    container level-2 {
        uses metric-cfg;
        description "level-2 specific configuration";
    }
    description "Metric configuration.";
}
container bfd {
    if-feature bfd;
    description "BFD Client Configuration.";
```

```
    uses bfd-types:client-cfg-parms;

    reference "RFC YYYY - YANG Data Model for Bidirectional
        Forwarding Detection (BFD).

-- Note to RFC Editor Please replace YYYY with published FC
    number for draft-ietf-bfd-yang.";

}
container address-families {
    if-feature nlpid-control;
    list address-family-list {
        key address-family;
        leaf address-family {
            type iana-rt-types:address-family;
            description "Address-family";
        }
        description "List of AFs.";
    }
    description "Interface address-families";
}
container mpls {
    container ldp {
        leaf igp-sync {
            if-feature ldp-igp-sync;
            type boolean;
            default false;
            description "Enables IGP/LDP synchronization";
        }
        description "LDP protocol related configuration.";
    }
    description "MPLS configuration for IS-IS interfaces";
}
uses interface-fast-reroute-config;
}

grouping multi-topology-interface-config {
    description "IS-IS interface topology configuration.";
    container metric {
        uses metric-cfg;
        container level-1 {
            uses metric-cfg;
            description "level-1 specific configuration";
        }
        container level-2 {
            uses metric-cfg;
            description "level-2 specific configuration";
        }
    }
}
```

```
        description "Metric IS-IS interface configuration.";
    }
}
grouping interface-state {
    description
        "IS-IS interface operational state.";
    uses adjacency-state;
    uses event-counters;
    uses packet-counters;
    leaf discontinuity-time {
        type yang:date-and-time;
        description
            "The time of the most recent occasion at which any one
            or more of this IS-IS interface's counters suffered a
            discontinuity.  If no such discontinuities have occurred
            since the IS-IS interface was last re-initialized, then
            this node contains the time the IS-IS interface was
            re-initialized which normally occurs when it was
            created.";
    }
}

/* Grouping for the hostname database */

grouping hostname-db {
    container hostnames {
        config false;
        list hostname {
            key system-id;
            leaf system-id {
                type system-id;
                description
                    "system-id associated with the hostname.";
            }
            leaf hostname {
                type string {
                    length "1..255";
                }
                description
                    "Hostname associated with the system-id
                    as defined in RFC5301.";
            }
        }
        description
            "List of system-id/hostname associations.";
    }
    description
        "Hostname to system-id mapping database.";
}
```

```
    description
      "Grouping for hostname to system-id mapping database.";
  }

/* Groupings for counters */

grouping system-counters {
  container system-counters {
    config false;
    list level {
      key level;

      leaf level {
        type level-number;
        description "IS-IS level.";
      }
      leaf corrupted-lsps {
        type uint32;
        description
          "Number of corrupted in-memory LSPs detected.
          LSPs received from the wire with a bad
          checksum are silently dropped and not counted.
          LSPs received from the wire with parse errors
          are counted by lsp-errors.";
      }
      leaf authentication-type-fails {
        type uint32;
        description
          "Number of authentication type mismatches.";
      }
      leaf authentication-fails {
        type uint32;
        description
          "Number of authentication key failures.";
      }
      leaf database-overload {
        type uint32;
        description
          "Number of times the database has become
          overloaded.";
      }
      leaf own-lsp-purge {
        type uint32;
        description
          "Number of times a zero-aged copy of the system's
          own LSP is received from some other IS-IS node.";
      }
      leaf manual-address-drop-from-area {
```

```
        type uint32;
        description
            "Number of times a manual address
             has been dropped from the area.";
    }
    leaf max-sequence {
        type uint32;
        description
            "Number of times the system has attempted
             to exceed the maximum sequence number.";
    }
    leaf sequence-number-skipped {
        type uint32;
        description
            "Number of times a sequence number skip has
             occurred.";
    }
    leaf id-len-mismatch {
        type uint32;
        description
            "Number of times a PDU is received with a
             different value for the ID field length
             than that of the receiving system.";
    }
    leaf partition-changes {
        type uint32;
        description
            "Number of partition changes detected.";
    }
    leaf lsp-errors {
        type uint32;
        description
            "Number of LSPs with errors we have received.";
    }
    leaf spf-runs {
        type uint32;
        description
            "Number of times we ran SPF at this level.";
    }
    description
        "List of supported levels.";
}
description
    "List counters for the IS-IS protocol instance";
}
description
    "Grouping for IS-IS system counters";
}
```

```
grouping event-counters {
  container event-counters {
    config false;
    leaf adjacency-changes {
      type uint32;
      description
        "The number of times an adjacency state change has
        occurred on this interface.";
    }
    leaf adjacency-number {
      type uint32;
      description
        "The number of adjacencies on this interface.";
    }
    leaf init-fails {
      type uint32;
      description
        "The number of times initialization of this
        interface has failed. This counts events such
        as PPP NCP failures. Failures to form an
        adjacency are counted by adjacency-rejects.";
    }
    leaf adjacency-rejects {
      type uint32;
      description
        "The number of times an adjacency has been
        rejected on this interface.";
    }
    leaf id-len-mismatch {
      type uint32;
      description
        "The number of times an IS-IS PDU with an ID
        field length different from that for this
        system has been received on this interface.";
    }
    leaf max-area-addresses-mismatch {
      type uint32;
      description
        "The number of times an IS-IS PDU has been
        received on this interface with the
        max area address field differing from that of
        this system.";
    }
    leaf authentication-type-fails {
      type uint32;
      description
        "Number of authentication type mismatches.";
    }
  }
}
```

```
    leaf authentication-fails {
        type uint32;
        description
            "Number of authentication key failures.";
    }
    leaf lan-dis-changes {
        type uint32;
        description
            "The number of times the DIS has changed on this
            interface at this level. If the interface type is
            point-to-point, the count is zero.";
    }
    description "IS-IS interface event counters.";
}
description
    "Grouping for IS-IS interface event counters";
}

grouping packet-counters {
    container packet-counters {
        config false;
        list level {
            key level;

            leaf level {
                type level-number;
                description "IS-IS level.";
            }
        }
        container iih {
            leaf in {
                type uint32;
                description "Received IIH PDUs.";
            }
            leaf out {
                type uint32;
                description "Sent IIH PDUs.";
            }
            description "Number of IIH PDUs received/sent.";
        }
        container ish {
            leaf in {
                type uint32;
                description "Received ISH PDUs.";
            }
            leaf out {
                type uint32;
                description "Sent ISH PDUs.";
            }
        }
    }
}
```

```
        description
            "ISH PDUs received/sent.";
    }
    container esh {
        leaf in {
            type uint32;
            description "Received ESH PDUs.";
        }
        leaf out {
            type uint32;
            description "Sent ESH PDUs.";
        }
        description "Number of ESH PDUs received/sent.";
    }
    container lsp {
        leaf in {
            type uint32;
            description "Received LSP PDUs.";
        }
        leaf out {
            type uint32;
            description "Sent LSP PDUs.";
        }
        description "Number of LSP PDUs received/sent.";
    }
    container psnp {
        leaf in {
            type uint32;
            description "Received PSNP PDUs.";
        }
        leaf out {
            type uint32;
            description "Sent PSNP PDUs.";
        }
        description "Number of PSNP PDUs received/sent.";
    }
    container csnp {
        leaf in {
            type uint32;
            description "Received CSNP PDUs.";
        }
        leaf out {
            type uint32;
            description "Sent CSNP PDUs.";
        }
        description "Number of CSNP PDUs received/sent.";
    }
    container unknown {
```



```
        leaf in {
            type uint32;
            description "Received unknown PDUs.";
        }
        description "Number of unknown PDUs received/sent.";
    }
    description
        "List of packet counter for supported levels.";
}
description "Packet counters per IS-IS level.";
}
description
    "Grouping for per IS-IS Level packet counters.";
}

/* Groupings for various log buffers */
grouping spf-log {
    container spf-log {
        config false;
        list event {
            key id;

            leaf id {
                type yang:counter32;
                description
                    "Event identifier - purely internal value.
                     It is expected the most recent events to have the bigger
                     id number.";
            }
            leaf spf-type {
                type enumeration {
                    enum full {
                        description "Full SPF computation.";
                    }
                    enum route-only {
                        description
                            "Route reachability only SPF computation";
                    }
                }
                description "Type of SPF computation performed.";
            }
            leaf level {
                type level-number;
                description
                    "IS-IS level number for SPF computation";
            }
            leaf schedule-timestamp {
                type yang:timestamp;
            }
        }
    }
}
```

```
        description
            "Timestamp of when the SPF computation was
            scheduled.";
    }
    leaf start-timestamp {
        type yang:timestamp;
        description
            "Timestamp of when the SPF computation started.";
    }
    leaf end-timestamp {
        type yang:timestamp;
        description
            "Timestamp of when the SPF computation ended.";
    }
    list trigger-lsp {
        key "lsp";
        leaf lsp {
            type lsp-id;
            description
                "LSP ID of the LSP triggering SPF computation.";
        }
        leaf sequence {
            type uint32;
            description
                "Sequence number of the LSP triggering SPF
                computation";
        }
        description
            "This list includes the LSPs that triggered the
            SPF computation.";
    }
    description
        "List of computation events - implemented as a
        wrapping buffer.";
}

description
    "This container lists the SPF computation events.";
}
description "Grouping for spf-log events.";
}

grouping lsp-log {
    container lsp-log {
        config false;
        list event {
            key id;
```

```
leaf id {
  type yang:counter32;
  description
    "Event identifier - purely internal value.
    It is expected the most recent events to have the bigger
    id number.";
}
leaf level {
  type level-number;
  description
    "IS-IS level number for LSP";
}
container lsp {
  leaf lsp {
    type lsp-id;
    description
      "LSP ID of the LSP.";
  }
  leaf sequence {
    type uint32;
    description
      "Sequence number of the LSP.";
  }
  description
    "LSP identification container - either the received
    LSP or the locally generated LSP.";
}

leaf received-timestamp {
  type yang:timestamp;
  description
    "This is the timestamp when the LSA was received.
    In case of local LSA update, the timestamp refers
    to the LSA origination time.";
}

leaf reason {
  type identityref {
    base lsp-log-reason;
  }
  description "Type of LSP change.";
}

description
  "List of LSP events - implemented as a
  wrapping buffer.";
```

```
        description
            "This container lists the LSP log.
            Local LSP modifications are also included
            in the list.";
    } description "Grouping for LSP log.";
}

/* Groupings for the LSDB description */

/* Unknown TLV and sub-TLV description */
grouping tlv {
    description
        "Type-Length-Value (TLV)";
    leaf type {
        type uint16;
        description "TLV type.";
    }
    leaf length {
        type uint16;
        description "TLV length (octets).";
    }
    leaf value {
        type yang:hex-string;
        description "TLV value.";
    }
}

grouping unknown-tlvs {
    description
        "Unknown TLVs grouping - Used for unknown TLVs or
        unknown sub-TLVs.";
    container unknown-tlvs {
        description "All unknown TLVs.";
        list unknown-tlv {
            description "Unknown TLV.";
            uses tlv;
        }
    }
}

/* TLVs and sub-TLVs for prefixes */

grouping prefix-reachability-attributes {
    description
        "Grouping for extended reachability attributes of an
```

```
        IPv4 or IPv6 prefix.";

    leaf external-prefix-flag {
        type boolean;
        description "External prefix flag.";
    }
    leaf readvertisement-flag {
        type boolean;
        description "Re-advertisement flag.";
    }
    leaf node-flag {
        type boolean;
        description "Node flag.";
    }
}

grouping prefix-ipv4-source-router-id {
    description
        "Grouping for the IPv4 source router ID of a prefix
        advertisement.";

    leaf ipv4-source-router-id {
        type inet:ipv4-address;
        description "IPv4 Source router ID address.";
    }
}

grouping prefix-ipv6-source-router-id {
    description
        "Grouping for the IPv6 source router ID of a prefix
        advertisement.";

    leaf ipv6-source-router-id {
        type inet:ipv6-address;
        description "IPv6 Source router ID address.";
    }
}

grouping prefix-attributes-extension {
    description "Prefix extended attributes
        as defined in RFC7794.";

    uses prefix-reachability-attributes;
    uses prefix-ipv4-source-router-id;
    uses prefix-ipv6-source-router-id;
}

grouping prefix-ipv4-std {
```

```
description
  "Grouping for attributes of an IPv4 standard prefix
   as defined in RFC1195.";
leaf ip-prefix {
  type inet:ipv4-address;
  description "IPv4 prefix address";
}
leaf prefix-len {
  type uint8;
  description "IPv4 prefix length (in bits)";
}
leaf i-e {
  type boolean;
  description
    "Internal or External (I/E) Metric bit value.
     Set to 'false' to indicate an internal metric.";
}
container default-metric {
  leaf metric {
    type std-metric;
    description "Default IS-IS metric for IPv4 prefix";
  }
  description "IS-IS default metric container.";
}
container delay-metric {
  leaf metric {
    type std-metric;
    description "IS-IS delay metric for IPv4 prefix";
  }
  leaf supported {
    type boolean;
    default "false";
    description
      "Indicates whether IS-IS delay metric is supported.";
  }
  description "IS-IS delay metric container.";
}
container expense-metric {
  leaf metric {
    type std-metric;
    description "IS-IS expense metric for IPv4 prefix";
  }
  leaf supported {
    type boolean;
    default "false";
    description
      "Indicates whether IS-IS expense metric is supported.";
  }
}
```

```
        description "IS-IS expense metric container.";
    }
    container error-metric {
        leaf metric {
            type std-metric;
            description
                "This leaf describes the IS-IS error metric value";
        }
        leaf supported {
            type boolean;
            default "false";
            description
                "Indicates whether IS-IS error metric is supported.";
        }
        description "IS-IS error metric container.";
    }
}

grouping prefix-ipv4-extended {
    description
        "Grouping for attributes of an IPv4 extended prefix
        as defined in RFC5305.";
    leaf up-down {
        type boolean;
        description "Value of up/down bit.
            Set to true when the prefix has been advertised down
            the hierarchy.";
    }
    leaf ip-prefix {
        type inet:ipv4-address;
        description "IPv4 prefix address";
    }
    leaf prefix-len {
        type uint8;
        description "IPv4 prefix length (in bits)";
    }
    leaf metric {
        type wide-metric;
        description "IS-IS wide metric value";
    }
    leaf-list tag {
        type uint32;
        description
            "List of 32-bit tags associated with the IPv4 prefix.";
    }
    leaf-list tag64 {
        type uint64;
        description
```

```
        "List of 64-bit tags associated with the IPv4 prefix.";
    }
    uses prefix-attributes-extension;
}

grouping prefix-ipv6-extended {
    description "Grouping for attributes of an IPv6 prefix
                as defined in RFC5308.";
    leaf up-down {
        type boolean;
        description "Value of up/down bit.
                    Set to true when the prefix has been advertised down
                    the hierarchy.";
    }
    leaf ip-prefix {
        type inet:ipv6-address;
        description "IPv6 prefix address";
    }
    leaf prefix-len {
        type uint8;
        description "IPv6 prefix length (in bits)";
    }
    leaf metric {
        type wide-metric;
        description "IS-IS wide metric value";
    }
    leaf-list tag {
        type uint32;
        description
            "List of 32-bit tags associated with the IPv4 prefix.";
    }
    leaf-list tag64 {
        type uint64;
        description
            "List of 64-bit tags associated with the IPv4 prefix.";
    }
    uses prefix-attributes-extension;
}

/* TLVs and sub-TLVs for neighbors */

grouping neighbor-link-attributes {
    description
        "Grouping for link attributes as defined
        in RFC5029";
    leaf link-attributes-flags {
        type uint16;
        description
```



```
        "Flags for the link attributes";
    }
}
grouping neighbor-gmpls-extensions {
    description
        "Grouping for GMPLS attributes of a neighbor as defined
        in RFC5307";
    leaf link-local-id {
        type uint32;
        description
            "Local identifier of the link.";
    }
    leaf remote-local-id {
        type uint32;
        description
            "Remote identifier of the link.";
    }
    leaf protection-capability {
        type uint8;
        description
            "Describes the protection capabilities
            of the link. This is the value of the
            first octet of the sub-TLV type 20 value.";
    }
    container interface-switching-capability {
        description
            "Interface switching capabilities of the link.";
        leaf switching-capability {
            type uint8;
            description
                "Switching capability of the link.";
        }
    }
    leaf encoding {
        type uint8;
        description
            "Type of encoding of the LSP being used.";
    }
    container max-lsp-bandwidths {
        description "Per-priority max LSP bandwidths.";
        list max-lsp-bandwidth {
            leaf priority {
                type uint8 {
                    range "0 .. 7";
                }
                description "Priority from 0 to 7.";
            }
            leaf bandwidth {
                type rt-types:bandwidth-ieee-float32;
            }
        }
    }
}
```

```
        description "max LSP bandwidth.";
    }
    description
        "List of max LSP bandwidths for different
        priorities.";
    }
}
container tdm-specific {
    when "../switching-capability = 100";
    description
        "Switching Capability-specific information applicable
        when switching type is TDM.";

    leaf minimum-lsp-bandwidth {
        type rt-types:bandwidth-ieee-float32;
        description "minimum LSP bandwidth.";
    }
    leaf indication {
        type uint8;
        description
            "The indication whether the interface supports Standard
            or Arbitrary SONET/SDH.";
    }
}
container psc-specific {
    when "../switching-capability >= 1 and
        ../switching-capability <= 4";
    description
        "Switching Capability-specific information applicable
        when switching type is PSC1,PSC2,PSC3 or PSC4.";

    leaf minimum-lsp-bandwidth {
        type rt-types:bandwidth-ieee-float32;
        description "minimum LSP bandwidth.";
    }
    leaf mtu {
        type uint16;
        units bytes;
        description
            "Interface MTU";
    }
}
}
}

grouping neighbor-extended-te-extensions {
    description
        "Grouping for TE attributes of a neighbor as defined
```

```
    in RFC8570";

container unidirectional-link-delay {
  description
    "Container for the average delay
    from the local neighbor to the remote one.";
  container flags {
    leaf-list unidirectional-link-delay-subtlv-flags {
      type identityref {
        base unidirectional-link-delay-subtlv-flag;
      }
      description
        "This list contains identities for the bits
        which are set.";
    }
    description
      "unidirectional-link-delay subTLV flags.";
  }
  leaf value {
    type uint32;
    units usec;
    description
      "Delay value expressed in microseconds.";
  }
}

container min-max-unidirectional-link-delay {
  description
    "Container for the min and max delay
    from the local neighbor to the remote one.";
  container flags {
    leaf-list min-max-unidirectional-link-delay-subtlv-flags {
      type identityref {
        base min-max-unidirectional-link-delay-subtlv-flag;
      }
      description
        "This list contains identities for the bits which are
        set.";
    }
    description
      "min-max-unidirectional-link-delay subTLV flags.";
  }
  leaf min-value {
    type uint32;
    units usec;
    description
      "Minimum delay value expressed in microseconds.";
  }
  leaf max-value {
```

```
        type uint32;
        units usec;
        description
            "Maximum delay value expressed in microseconds.";
    }
}
container unidirectional-link-delay-variation {
    description
        "Container for the average delay variation
        from the local neighbor to the remote one.";
    leaf value {
        type uint32;
        units usec;
        description
            "Delay variation value expressed in microseconds.";
    }
}
container unidirectional-link-loss {
    description
        "Container for the packet loss
        from the local neighbor to the remote one.";
    container flags {
        leaf-list unidirectional-link-loss-subtlv-flags {
            type identityref {
                base unidirectional-link-loss-subtlv-flag;
            }
            description
                "This list contains identities for the bits which are
                set.";
        }
        description
            "unidirectional-link-loss subTLV flags.";
    }
    leaf value {
        type uint32;
        units percent;
        description
            "Link packet loss expressed as a percentage
            of the total traffic sent over a configurable interval.";
    }
}
container unidirectional-link-residual-bandwidth {
    description
        "Container for the residual bandwidth
        from the local neighbor to the remote one.";
    leaf value {
        type rt-types:bandwidth-ieee-float32;
        units Bps;
    }
}
```

```
        description
            "Residual bandwidth.";
    }
}
container unidirectional-link-available-bandwidth {
    description
        "Container for the available bandwidth
        from the local neighbor to the remote one.";
    leaf value {
        type rt-types:bandwidth-ieee-float32;
        units Bps;
        description
            "Available bandwidth.";
    }
}
container unidirectional-link-utilized-bandwidth {
    description
        "Container for the utilized bandwidth
        from the local neighbor to the remote one.";
    leaf value {
        type rt-types:bandwidth-ieee-float32;
        units Bps;
        description
            "Utilized bandwidth.";
    }
}
}

grouping neighbor-te-extensions {
    description
        "Grouping for TE attributes of a neighbor as defined
        in RFC5305";
    leaf admin-group {
        type uint32;
        description
            "Administrative group/Resource Class/Color.";
    }
    container local-if-ipv4-addrs {
        description "All local interface IPv4 addresses.";
        leaf-list local-if-ipv4-addr {
            type inet:ipv4-address;
            description
                "List of local interface IPv4 addresses.";
        }
    }
    container remote-if-ipv4-addrs {
        description "All remote interface IPv4 addresses.";
        leaf-list remote-if-ipv4-addr {
```

```
        type inet:ipv4-address;
        description
            "List of remote interface IPv4 addresses.";
    }
}
leaf te-metric {
    type uint32;
    description "TE metric.";
}
leaf max-bandwidth {
    type rt-types:bandwidth-ieee-float32;
    description "Maximum bandwidth.";
}
leaf max-reservable-bandwidth {
    type rt-types:bandwidth-ieee-float32;
    description "Maximum reservable bandwidth.";
}
container unreserved-bandwidths {
    description "All unreserved bandwidths.";
    list unreserved-bandwidth {
        leaf priority {
            type uint8 {
                range "0 .. 7";
            }
            description "Priority from 0 to 7.";
        }
        leaf unreserved-bandwidth {
            type rt-types:bandwidth-ieee-float32;
            description "Unreserved bandwidth.";
        }
    }
    description
        "List of unreserved bandwidths for different
        priorities.";
}
}
}

grouping neighbor-extended {
    description
        "Grouping for attributes of an IS-IS extended neighbor.";
    leaf neighbor-id {
        type extended-system-id;
        description "system-id of the extended neighbor.";
    }
}
container instances {
    description "List of all adjacencies between the local
        system and the neighbor system-id.";
    list instance {
```

```
    key id;

    leaf id {
        type uint32;
        description "Unique identifier of an instance of a
            particular neighbor.";
    }
    leaf metric {
        type wide-metric;
        description "IS-IS wide metric for extended neighbor";
    }
    uses neighbor-gmpls-extensions;
    uses neighbor-te-extensions;
    uses neighbor-extended-te-extensions;
    uses neighbor-link-attributes;
    uses unknown-tlvs;
    description "Instance of a particular adjacency.";
}
}
}

grouping neighbor {
    description "IS-IS standard neighbor grouping.";
    leaf neighbor-id {
        type extended-system-id;
        description "IS-IS neighbor system-id";
    }
    container instances {
        description "List of all adjacencies between the local
            system and the neighbor system-id.";
        list instance {
            key id;

            leaf id {
                type uint32;
                description "Unique identifier of an instance of a
                    particular neighbor.";
            }
            leaf i-e {
                type boolean;
                description
                    "Internal or External (I/E) Metric bit value.
                    Set to 'false' to indicate an internal metric.";
            }
            container default-metric {
                leaf metric {
                    type std-metric;
                    description "IS-IS default metric value";
                }
            }
        }
    }
}
```

```
    }
    description "IS-IS default metric container";
  }
  container delay-metric {
    leaf metric {
      type std-metric;
      description "IS-IS delay metric value";
    }
    leaf supported {
      type boolean;
      default "false";
      description "IS-IS delay metric supported";
    }
    description "IS-IS delay metric container";
  }
  container expense-metric {
    leaf metric {
      type std-metric;
      description "IS-IS expense metric value";
    }
    leaf supported {
      type boolean;
      default "false";
      description "IS-IS expense metric supported";
    }
    description "IS-IS expense metric container";
  }
  container error-metric {
    leaf metric {
      type std-metric;
      description "IS-IS error metric value";
    }
    leaf supported {
      type boolean;
      default "false";
      description "IS-IS error metric supported";
    }
    description "IS-IS error metric container";
  }
  description "Instance of a particular adjacency
    as defined in ISO10589.";
}
}
}

/* Top-level TLVs */

grouping tlv132-ipv4-addresses {
```



```
    leaf-list ipv4-addresses {
      type inet:ipv4-address;
      description
        "List of IPv4 addresses of the IS-IS node - IS-IS
        reference is TLV 132.";
    }
    description "Grouping for TLV132.";
  }
  grouping tlv232-ipv6-addresses {
    leaf-list ipv6-addresses {
      type inet:ipv6-address;
      description
        "List of IPv6 addresses of the IS-IS node - IS-IS
        reference is TLV 232.";
    }
    description "Grouping for TLV232.";
  }
  grouping tlv134-ipv4-te-rid {
    leaf ipv4-te-routerid {
      type inet:ipv4-address;
      description
        "IPv4 Traffic Engineering router ID of the IS-IS node -
        IS-IS reference is TLV 134.";
    }
    description "Grouping for TLV134.";
  }
  grouping tlv140-ipv6-te-rid {
    leaf ipv6-te-routerid {
      type inet:ipv6-address;
      description
        "IPv6 Traffic Engineering router ID of the IS-IS node -
        IS-IS reference is TLV 140.";
    }
    description "Grouping for TLV140.";
  }
  grouping tlv129-protocols {
    leaf-list protocol-supported {
      type uint8;
      description
        "List of supported protocols of the IS-IS node -
        IS-IS reference is TLV 129.";
    }
    description "Grouping for TLV129.";
  }
  grouping tlv137-hostname {
    leaf dynamic-hostname {
      type string;
      description
```

```
        "Host Name of the IS-IS node - IS-IS reference
        is TLV 137.";
    }
    description "Grouping for TLV137.";
}
grouping tlv10-authentication {
    container authentication {
        leaf authentication-type {
            type identityref {
                base key-chain:crypto-algorithm;
            }
            description
                "Authentication type to be used with IS-IS node.";
        }
        leaf authentication-key {
            type string;
            description
                "Authentication key to be used. For security reasons,
                the authentication key MUST NOT be presented in
                a clear text format in response to any request
                (e.g., via get, get-config).";
        }
        description
            "IS-IS node authentication information container -
            IS-IS reference is TLV 10.";
    }
    description "Grouping for TLV10.";
}
grouping tlv229-mt {
    container mt-entries {
        list topology {
            description
                "List of topologies supported";

            leaf mt-id {
                type uint16 {
                    range "0 .. 4095";
                }
                description
                    "Multi-Topology identifier of topology.";
            }
        }
        container attributes {
            leaf-list flags {
                type identityref {
                    base tlv229-flag;
                }
                description
                    "This list contains identities for the bits which are
```

```
        set.";
    }
    description
        "TLV 229 flags.";
}
}
description
    "IS-IS node topology information container -
    IS-IS reference is TLV 229.";
}
description "Grouping for TLV229.";
}

grouping tlv242-router-capabilities {
    container router-capabilities {
        list router-capability {
            container flags {
                leaf-list router-capability-flags {
                    type identityref {
                        base router-capability-flag;
                    }
                    description
                        "This list contains identities for the bits which are
                        set.";
                }
                description
                    "Router capability flags.";
            }
            container node-tags {
                if-feature node-tag;
                list node-tag {
                    leaf tag {
                        type uint32;
                        description "Node tag value.";
                    }
                    description "List of tags.";
                }
                description "Container for node admin tags";
            }
        }
        description "List of router capability TLVs.";
    }
}
```

```
    }
    description "Grouping for TLV242.";
}

grouping tlv138-srlg {
  description
    "Grouping for TLV138.";
  container links-srlgs {
    list links {
      leaf neighbor-id {
        type extended-system-id;
        description "system-id of the extended neighbor.";
      }
      leaf flags {
        type uint8;
        description
          "Flags associated with the link.";
      }
      leaf link-local-id {
        type union {
          type inet:ip-address;
          type uint32;
        }
        description
          "Local identifier of the link.
          It could be an IPv4 address or a local identifier.";
      }
      leaf link-remote-id {
        type union {
          type inet:ip-address;
          type uint32;
        }
        description
          "Remote identifier of the link.
          It could be an IPv4 address or a remotely learned
          identifier.";
      }
    }
    container srlgs {
      description "List of SRLGs.";
      leaf-list srlg {
        type uint32;
        description
          "SRLG value of the link.";
      }
    }
    description
      "SRLG attribute of a link.";
  }
}
```

```
        description
            "List of links with SRLGs";
    }
}

/* Grouping for LSDB description */

grouping lsp-entry {
    description "IS-IS LSP database entry grouping";

    leaf decoded-completed {
        type boolean;
        description "IS-IS LSP body fully decoded.";
    }
    leaf raw-data {
        type yang:hex-string;
        description
            "The hexadecimal representation of the complete LSP in
            network-byte order (NBO) as received or originated.";
    }
    leaf lsp-id {
        type lsp-id;
        description "LSP ID of the LSP";
    }
    leaf checksum {
        type uint16;
        description "LSP checksum";
    }
    leaf remaining-lifetime {
        type uint16;
        units "seconds";
        description
            "Remaining lifetime (in seconds) until LSP expiration.";
    }
    leaf sequence {
        type uint32;
        description
            "This leaf describes the sequence number of the LSP.";
    }
    container attributes {
        leaf-list lsp-flags {
            type identityref {
                base lsp-flag;
            }
            description
                "This list contains identities for the bits which are
                set.";
        }
    }
}
```

```
        description "LSP attributes.";
    }

    uses tlv132-ipv4-addresses;
    uses tlv232-ipv6-addresses;
    uses tlv134-ipv4-te-rid;
    uses tlv140-ipv6-te-rid;
    uses tlv129-protocols;
    uses tlv137-hostname;
    uses tlv10-authentication;
    uses tlv229-mt;
    uses tlv242-router-capabilities;
    uses tlv138-srlg;
    uses unknown-tlvs;

    container is-neighbor {
        list neighbor {
            key neighbor-id;

            uses neighbor;
            description "List of neighbors.";
        }
        description
            "Standard IS neighbors container - IS-IS reference is
             TLV 2.";
    }

    container extended-is-neighbor {
        list neighbor {
            key neighbor-id;

            uses neighbor-extended;
            description
                "List of extended IS neighbors";
        }
        description
            "Standard IS extended neighbors container - IS-IS
             reference is TLV 22";
    }

    container ipv4-internal-reachability {
        list prefixes {
            uses prefix-ipv4-std;
            description "List of prefixes.";
        }
        description
            "IPv4 internal reachability information container - IS-IS
             reference is TLV 128.";
    }
```

```
}

container ipv4-external-reachability {
  list prefixes {
    uses prefix-ipv4-std;
    description "List of prefixes.";
  }
  description
    "IPv4 external reachability information container -
    IS-IS reference is TLV 130.";
}

container extended-ipv4-reachability {
  list prefixes {
    uses prefix-ipv4-extended;
    uses unknown-tlvs;
    description "List of prefixes.";
  }
  description
    "IPv4 extended reachability information container -
    IS-IS reference is TLV 135.";
}

container mt-is-neighbor {
  list neighbor {
    leaf mt-id {
      type uint16 {
        range "0 .. 4095";
      }
      description "Multi-topology (MT) identifier";
    }
    uses neighbor-extended;
    description "List of neighbors.";
  }
  description
    "IS-IS multi-topology neighbor container - IS-IS
    reference is TLV 223.";
}

container mt-extended-ipv4-reachability {
  list prefixes {
    leaf mt-id {
      type uint16 {
        range "0 .. 4095";
      }
      description "Multi-topology (MT) identifier";
    }
    uses prefix-ipv4-extended;
  }
}
```

```
        uses unknown-tlvs;
        description "List of extended prefixes.";
    }
    description
        "IPv4 multi-topology (MT) extended reachability
        information container - IS-IS reference is TLV 235.";
}

container mt-ipv6-reachability {
    list prefixes {
        leaf MT-ID {
            type uint16 {
                range "0 .. 4095";
            }
            description "Multi-topology (MT) identifier";
        }
        uses prefix-ipv6-extended;
        uses unknown-tlvs;
        description "List of IPv6 extended prefixes.";
    }
    description
        "IPv6 multi-topology (MT) extended reachability
        information container - IS-IS reference is TLV 237.";
}

container ipv6-reachability {
    list prefixes {
        uses prefix-ipv6-extended;
        uses unknown-tlvs;
        description "List of IPv6 prefixes.";
    }
    description
        "IPv6 reachability information container - IS-IS
        reference is TLV 236.";
}
}

grouping lsdb {
    description "Link State Database (LSDB) grouping";
    container database {
        config false;
        list levels {
            key level;

            leaf level {
                type level-number;
                description "LSDB level number (1 or 2)";
            }
        }
    }
}
```



```
    list lsp {
      key lsp-id;
      uses lsp-entry;
      description "List of LSPs in LSDB";
    }
    description "List of LSPs for the LSDB level container";
  }
  description "IS-IS Link State database container";
}
}
```

```
/* Augmentations */
```

```
augment "/rt:routing/"
+ "rt:ribs/rt:rib/rt:routes/rt:route" {
  when "rt:source-protocol = 'isis:isis'" {
    description "IS-IS-specific route attributes.";
  }
  uses route-content;
  description
    "This augments route object in RIB with IS-IS-specific
    attributes.";
}
```

```
augment "/if:interfaces/if:interface" {
  leaf clns-mtu {
    if-feature osi-interface;
    type uint16;
    description "CLNS MTU of the interface";
  }
  description "ISO specific interface parameters.";
}
```

```
augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol" {
  when "rt:type = 'isis:isis'" {
    description
      "This augment is only valid when routing protocol
      instance type is 'isis'";
  }
  description
    "This augments a routing protocol instance with IS-IS
    specific parameters.";
  container isis {
```

```
must "count(area-address) > 0" {
  error-message
    "At least one area-address must be configured.";
  description
    "Enforce configuration of at least one area.";
}

uses instance-config;
uses instance-state;

container topologies {
  if-feature multi-topology;
  list topology {
    key "name";
    leaf enable {
      type boolean;
      description "Topology enable configuration";
    }
    leaf name {
      type leafref {
        path "../.../.../.../rt:ribs/rt:rib/rt:name";
      }
      description
        "Routing Information Base (RIB) corresponding
        to topology.";
    }
  }

  uses multi-topology-config;

  description "List of topologies";
}
description "Multi-topology container";
}

container interfaces {
  list interface {
    key "name";
    leaf name {
      type if:interface-ref;

      description
        "Reference to the interface within
        the routing-instance.";
    }
  }
  uses interface-config;
  uses interface-state;
  container topologies {
    if-feature multi-topology;
    list topology {
```

```
        key name;

        leaf name {
            type leafref {
                path "../../../../../../../../../"+
                    "rt:ribs/rt:rib/rt:name";
            }

            description
                "Routing Information Base (RIB) corresponding
                to topology.";
        }
        uses multi-topology-interface-config;
        description "List of interface topologies";
    }
    description "Multi-topology container";
}
description "List of IS-IS interfaces.";
}
description
    "IS-IS interface specific configuration container";
}

description
    "IS-IS configuration/state top-level container";
}

/* RPC methods */

rpc clear-adjacency {
    description
        "This RPC request clears a particular set of IS-IS
        adjacencies. If the operation fails due to an internal
        reason, then the error-tag and error-app-tag should be
        set indicating the reason for the failure.";
    input {

        leaf routing-protocol-instance-name {
            type leafref {
                path "/rt:routing/rt:control-plane-protocols/"
                    + "rt:control-plane-protocol/rt:name";
            }
            mandatory "true";
            description
                "Name of the IS-IS protocol instance whose IS-IS
                adjacency is being cleared."
        }
    }
}
```

```
        If the corresponding IS-IS instance doesn't exist,
        then the operation will fail with an error-tag of
        'data-missing' and an error-app-tag of
        'routing-protocol-instance-not-found'.";
    }
    leaf level {
        type level;
        description
            "IS-IS level of the adjacency to be cleared. If the
            IS-IS level is level-1-2, both level 1 and level 2
            adjacencies would be cleared.

            If the value provided is different from the one
            authorized in the enum type, then the operation
            SHALL fail with an error-tag of 'data-missing' and
            an error-app-tag of 'bad-isis-level'.";
    }
    leaf interface {
        type if:interface-ref;
        description
            "IS-IS interface name.

            If the corresponding IS-IS interface doesn't exist,
            then the operation SHALL fail with an error-tag of
            'data-missing' and an error-app-tag of
            'isis-interface-not-found'.";
    }
}

rpc clear-database {
    description
        "This RPC request clears a particular IS-IS database. If
        the operation fails for an IS-IS internal reason, then
        the error-tag and error-app-tag should be set
        indicating the reason for the failure.";
    input {
        leaf routing-protocol-instance-name {
            type leafref {
                path "/rt:routing/rt:control-plane-protocols/"
                    + "rt:control-plane-protocol/rt:name";
            }
            mandatory "true";
            description
                "Name of the IS-IS protocol instance whose IS-IS
                database(s) is/are being cleared.

                If the corresponding IS-IS instance doesn't exist,
```

```
        then the operation will fail with an error-tag of
        'data-missing' and an error-app-tag of
        'routing-protocol-instance-not-found'.";
    }
    leaf level {
        type level;
        description
            "IS-IS level of the adjacency to be cleared. If the
            IS-IS level is level-1-2, both level 1 and level 2
            databases would be cleared.

            If the value provided is different from the one
            authorized in the enum type, then the operation
            SHALL fail with an error-tag of 'data-missing' and
            an error-app-tag of 'bad-isis-level'.";
    }
}

/* Notifications */

notification database-overload {
    uses notification-instance-hdr;

    leaf overload {
        type enumeration {
            enum off {
                description
                    "Indicates IS-IS instance has left overload state";
            }
            enum on {
                description
                    "Indicates IS-IS instance has entered overload state";
            }
        }
        description "New overload state of the IS-IS instance";
    }
    description
        "This notification is sent when an IS-IS instance
        overload state changes.";
}

notification lsp-too-large {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
```

```
    leaf pdu-size {
      type uint32;
      description "Size of the LSP PDU";
    }
    leaf lsp-id {
      type lsp-id;
      description "LSP ID";
    }
    description
      "This notification is sent when we attempt to propagate
      an LSP that is larger than the dataLinkBlockSize (ISO10589)
      for the circuit. The notification generation must be
      throttled with at least 5 seconds between successive
      notifications.";
  }

  notification if-state-change {
    uses notification-instance-hdr;
    uses notification-interface-hdr;

    leaf state {
      type if-state-type;
      description "Interface state.";
    }
    description
      "This notification is sent when an interface
      state change is detected.";
  }

  notification corrupted-lsp-detected {
    uses notification-instance-hdr;
    leaf lsp-id {
      type lsp-id;
      description "LSP ID";
    }
    description
      "This notification is sent when we find that
      an LSP that was stored in memory has become
      corrupted.";
  }

  notification attempt-to-exceed-max-sequence {
    uses notification-instance-hdr;
    leaf lsp-id {
      type lsp-id;
      description "LSP ID";
    }
    description
```

```
        "This notification is sent when the system
        wraps the 32-bit sequence counter of an LSP.";
    }

notification id-len-mismatch {
    uses notification-instance-hdr;
    uses notification-interface-hdr;

    leaf pdu-field-len {
        type uint8;
        description "Size of the ID length in the received PDU";
    }
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when we receive a PDU
        with a different value for the system-id length.
        The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification max-area-addresses-mismatch {
    uses notification-instance-hdr;
    uses notification-interface-hdr;

    leaf max-area-addresses {
        type uint8;
        description "Received number of supported areas";
    }
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when we receive a PDU
        with a different value for the Maximum Area Addresses.
        The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification own-lsp-purge {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf lsp-id {
```

```
        type lsp-id;
        description "LSP ID";
    }
    description
        "This notification is sent when the system receives
        a PDU with its own system-id and zero age.";
}

notification sequence-number-skipped {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf lsp-id {
        type lsp-id;
        description "LSP ID";
    }
    description
        "This notification is sent when the system receives a
        PDU with its own system-id and different contents. The
        system has to originate the LSP with a higher sequence
        number.";
}

notification authentication-type-failure {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when the system receives a
        PDU with the wrong authentication type field.
        The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification authentication-failure {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when the system receives
        a PDU with the wrong authentication information.
        The notification generation must be throttled
```



```
        with at least 5 seconds between successive
        notifications.";
    }

notification version-skew {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf protocol-version {
        type uint8;
        description "Protocol version received in the PDU.";
    }
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when the system receives a
        PDU with a different protocol version number.
        The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification area-mismatch {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    description
        "This notification is sent when the system receives a
        Hello PDU from an IS that does not share any area
        address. The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification rejected-adjacency {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf raw-pdu {
        type binary;
        description
            "Received raw PDU.";
    }
    leaf reason {
        type string {
```

```
        length "0..255";
    }
    description
        "The system may provide a reason to reject the
        adjacency. If the reason is not available,
        the reason string will not be returned.
        The expected format is a single line text.";
    }
    description
        "This notification is sent when the system receives a
        Hello PDU from an IS but does not establish an adjacency
        for some reason. The notification generation must be
        throttled with at least 5 seconds between successive
        notifications.";
}

notification protocols-supported-mismatch {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
    leaf-list protocols {
        type uint8;
        description
            "List of protocols supported by the remote system.";
    }
    description
        "This notification is sent when the system receives a
        non-pseudonode LSP that has no matching protocols
        supported. The notification generation must be throttled
        with at least 5 seconds between successive
        notifications.";
}

notification lsp-error-detected {
    uses notification-instance-hdr;
    uses notification-interface-hdr;
    leaf lsp-id {
        type lsp-id;
        description "LSP ID.";
    }
    leaf raw-pdu {
        type binary;
        description "Received raw PDU.";
    }
}
```

```
leaf error-offset {
  type uint32;
  description
    "If the problem is a malformed TLV, the error-offset
    points to the start of the TLV. If the problem is with
    the LSP header, the error-offset points to the errant
    byte";
}
leaf tlv-type {
  type uint8;
  description
    "If the problem is a malformed TLV, the tlv-type is set
    to the type value of the suspicious TLV. Otherwise,
    this leaf is not present.";
}
description
  "This notification is sent when the system receives an
  LSP with a parse error. The notification generation must
  be throttled with at least 5 seconds between successive
  notifications.";
}

notification adjacency-state-change {
  uses notification-instance-hdr;
  uses notification-interface-hdr;
  leaf neighbor {
    type string {
      length "1..255";
    }
    description
      "Name of the neighbor.
      It corresponds to the hostname associated
      with the system-id of the neighbor in the
      mapping database (RFC5301).
      If the name of the neighbor is
      not available, it is not returned.";
  }
  leaf neighbor-system-id {
    type system-id;
    description "Neighbor system-id";
  }
  leaf state {
    type adj-state-type;

    description "New state of the IS-IS adjacency.";
  }
  leaf reason {
    type string {
```

```
        length "1..255";
    }
    description
        "If the adjacency is going to DOWN, this leaf provides
        a reason for the adjacency going down. The reason is
        provided as a text. If the adjacency is going to UP, no
        reason is provided. The expected format is a single line
        text.";
    }
    description
        "This notification is sent when an IS-IS adjacency
        moves to Up state or to Down state.";
}

notification lsp-received {
    uses notification-instance-hdr;
    uses notification-interface-hdr;

    leaf lsp-id {
        type lsp-id;
        description "LSP ID";
    }
    leaf sequence {
        type uint32;
        description "Sequence number of the received LSP.";
    }
    leaf received-timestamp {
        type yang:timestamp;

        description "Timestamp when the LSP was received.";
    }
    leaf neighbor-system-id {
        type system-id;
        description "Neighbor system-id of LSP sender";
    }
    description
        "This notification is sent when an LSP is received.
        The notification generation must be throttled with at
        least 5 seconds between successive notifications.";
}

notification lsp-generation {
    uses notification-instance-hdr;

    leaf lsp-id {
        type lsp-id;
        description "LSP ID";
    }
}
```

```
    leaf sequence {
      type uint32;
      description "Sequence number of the received LSP.";
    }
    leaf send-timestamp {
      type yang:timestamp;

      description "Timestamp when our LSP was regenerated.";
    }
    description
      "This notification is sent when an LSP is regenerated.
      The notification generation must be throttled with at
      least 5 seconds between successive notifications.";
  }
}
<CODE ENDS>
```

7. Security Considerations

The YANG modules specified in this document define a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC8446].

The NETCONF Access Control Model (NACM) [RFC8341] provides the means to restrict access for particular NETCONF or RESTCONF users to a pre-configured subset of all available NETCONF or RESTCONF protocol operations and content.

There are a number of data nodes defined in ietf-isis.yang module that are writable/creatable/deletable (i.e., config true, which is the default). These data nodes may be considered sensitive or vulnerable in some network environments. Write operations (e.g., edit-config) to these data nodes without proper protection can have a negative effect on network operations. Writable data node represent configuration of each instance and interface. These correspond to the following schema nodes:

```
/isis
```

```
/isis/interfaces/interface[name]
```

For IS-IS, the ability to modify IS-IS configuration will allow the entire IS-IS domain to be compromised including forming adjacencies with unauthorized routers to misroute traffic or mount a massive

Denial-of-Service (DoS) attack. For example, adding IS-IS on any unprotected interface could allow an IS-IS adjacency to be formed with an unauthorized and malicious neighbor. Once an adjacency is formed, traffic could be hijacked. As a simpler example, a Denial-Of-Service attack could be mounted by changing the cost of an IS-IS interface to be asymmetric such that a hard routing loop ensues. In general, unauthorized modification of most IS-IS features will pose their own set of security risks and the "Security Considerations" in the respective reference RFCs should be consulted.

Some of the readable data nodes in the `ietf-isis.yang` module may be considered sensitive or vulnerable in some network environments. It is thus important to control read access (e.g., via `get`, `get-config`, or `notification`) to these data nodes. The exposure of the Link State Database (LSDB) will expose the detailed topology of the network. Similarly, the IS-IS local RIB exposes the reachable prefixes in the IS-IS routing domain. The Link State Database (LSDB) and local RIB are represented by the following schema nodes:

```
/isis/database
```

```
/isis/local-rib
```

Exposure of the Link State Database and local RIB include information beyond the scope of the IS-IS router and this may be undesirable since exposure may facilitate other attacks. Additionally, the complete IP network topology and, if deployed, the traffic engineering topology of the IS-IS domain can be reconstructed from the Link State Database. Though not as straightforward, the IS-IS local RIB can also be discover topological information. Network operators may consider their topologies to be sensitive confidential data.

For IS-IS authentication, configuration is supported via the specification of `key-chain` [RFC8177] or the direct specification of key and authentication algorithm. Hence, authentication configuration using the `"auth-table-trailer"` case in the `"authentication"` container inherits the security considerations of [RFC8177]. This includes the considerations with respect to the local storage and handling of authentication keys.

Some of the RPC operations in this YANG module may be considered sensitive or vulnerable in some network environments. It is thus important to control access to these operations. The IS-IS YANG module support the `"clear-adjacency"` and `"clear-database"` RPCs. If access to either of these is compromised, they can result in temporary network outages be employed to mount DoS attacks.

The actual authentication key data (whether locally specified or part of a key-chain) is sensitive and needs to be kept secret from unauthorized parties; compromise of the key data would allow an attacker to forge IS-IS traffic that would be accepted as authentic, potentially compromising the entirety IS-IS domain.

The model describes several notifications, implementations must rate-limit the generation of these notifications to avoid creating significant notification load. Otherwise, this notification load may have some side effects on the system stability and may be exploited as an attack vector.

8. Contributors

The authors would like to thank Kiran Agrahara Sreenivasa, Dean Bogdanovic, Yingzhen Qu, Yi Yang, Jeff Tanstura for their major contributions to the draft.

9. Acknowledgements

The authors would like to thank Tom Petch, Alvaro Retana, Stewart Bryant, Barry Leiba, Benjamin Kaduk and Adam Roach, and Roman Danyliw for their review and comments.

10. IANA Considerations

The IANA is requested to assign two new URIs from the IETF XML registry [RFC3688]. Authors are suggesting the following URI:

```
URI: urn:ietf:params:xml:ns:yang:ietf-isis
Registrant Contact: The IESG
XML: N/A, the requested URI is an XML namespace
```

This document also requests one new YANG module name in the YANG Module Names registry [RFC6020] with the following suggestion:

```
name: ietf-isis
namespace: urn:ietf:params:xml:ns:yang:ietf-isis
prefix: isis
reference: RFC XXXX
```

11. References

11.1. Normative References

- [I-D.ietf-bfd-yang]
Rahman, R., Zheng, L., Jethanandani, M., Networks, J., and G. Mirsky, "YANG Data Model for Bidirectional Forwarding Detection (BFD)", draft-ietf-bfd-yang-17 (work in progress), August 2018.
- [ISO-10589]
"Intermediate System to Intermediate System intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode network service (ISO 8473)", International Standard 10589: 2002, Second Edition, 2002.
- [RFC1195] Callon, R., "Use of OSI IS-IS for routing in TCP/IP and dual environments", RFC 1195, DOI 10.17487/RFC1195, December 1990, <<https://www.rfc-editor.org/info/rfc1195>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC4090] Pan, P., Ed., Swallow, G., Ed., and A. Atlas, Ed., "Fast Reroute Extensions to RSVP-TE for LSP Tunnels", RFC 4090, DOI 10.17487/RFC4090, May 2005, <<https://www.rfc-editor.org/info/rfc4090>>.
- [RFC5029] Vasseur, JP. and S. Previdi, "Definition of an IS-IS Link Attribute Sub-TLV", RFC 5029, DOI 10.17487/RFC5029, September 2007, <<https://www.rfc-editor.org/info/rfc5029>>.
- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.
- [RFC5130] Previdi, S., Shand, M., Ed., and C. Martin, "A Policy Control Mechanism in IS-IS Using Administrative Tags", RFC 5130, DOI 10.17487/RFC5130, February 2008, <<https://www.rfc-editor.org/info/rfc5130>>.

- [RFC5286] Atlas, A., Ed. and A. Zinin, Ed., "Basic Specification for IP Fast Reroute: Loop-Free Alternates", RFC 5286, DOI 10.17487/RFC5286, September 2008, <<https://www.rfc-editor.org/info/rfc5286>>.
- [RFC5301] McPherson, D. and N. Shen, "Dynamic Hostname Exchange Mechanism for IS-IS", RFC 5301, DOI 10.17487/RFC5301, October 2008, <<https://www.rfc-editor.org/info/rfc5301>>.
- [RFC5302] Li, T., Smit, H., and T. Przygienda, "Domain-Wide Prefix Distribution with Two-Level IS-IS", RFC 5302, DOI 10.17487/RFC5302, October 2008, <<https://www.rfc-editor.org/info/rfc5302>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5306] Shand, M. and L. Ginsberg, "Restart Signaling for IS-IS", RFC 5306, DOI 10.17487/RFC5306, October 2008, <<https://www.rfc-editor.org/info/rfc5306>>.
- [RFC5307] Kompella, K., Ed. and Y. Rekhter, Ed., "IS-IS Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 5307, DOI 10.17487/RFC5307, October 2008, <<https://www.rfc-editor.org/info/rfc5307>>.
- [RFC5308] Hopps, C., "Routing IPv6 with IS-IS", RFC 5308, DOI 10.17487/RFC5308, October 2008, <<https://www.rfc-editor.org/info/rfc5308>>.
- [RFC5443] Jork, M., Atlas, A., and L. Fang, "LDP IGP Synchronization", RFC 5443, DOI 10.17487/RFC5443, March 2009, <<https://www.rfc-editor.org/info/rfc5443>>.
- [RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", RFC 5880, DOI 10.17487/RFC5880, June 2010, <<https://www.rfc-editor.org/info/rfc5880>>.
- [RFC5881] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD) for IPv4 and IPv6 (Single Hop)", RFC 5881, DOI 10.17487/RFC5881, June 2010, <<https://www.rfc-editor.org/info/rfc5881>>.
- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.

- [RFC6119] Harrison, J., Berger, J., and M. Bartlett, "IPv6 Traffic Engineering in IS-IS", RFC 6119, DOI 10.17487/RFC6119, February 2011, <<https://www.rfc-editor.org/info/rfc6119>>.
- [RFC6232] Wei, F., Qin, Y., Li, Z., Li, T., and J. Dong, "Purge Originator Identification TLV for IS-IS", RFC 6232, DOI 10.17487/RFC6232, May 2011, <<https://www.rfc-editor.org/info/rfc6232>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.
- [RFC6991] Schoenwaelder, J., Ed., "Common YANG Data Types", RFC 6991, DOI 10.17487/RFC6991, July 2013, <<https://www.rfc-editor.org/info/rfc6991>>.
- [RFC7490] Bryant, S., Filmsils, C., Previdi, S., Shand, M., and N. So, "Remote Loop-Free Alternate (LFA) Fast Reroute (FRR)", RFC 7490, DOI 10.17487/RFC7490, April 2015, <<https://www.rfc-editor.org/info/rfc7490>>.
- [RFC7794] Ginsberg, L., Ed., Decraene, B., Previdi, S., Xu, X., and U. Chunduri, "IS-IS Prefix Attributes for Extended IPv4 and IPv6 Reachability", RFC 7794, DOI 10.17487/RFC7794, March 2016, <<https://www.rfc-editor.org/info/rfc7794>>.
- [RFC7917] Sarkar, P., Ed., Gredler, H., Hegde, S., Litkowski, S., and B. Decraene, "Advertising Node Administrative Tags in IS-IS", RFC 7917, DOI 10.17487/RFC7917, July 2016, <<https://www.rfc-editor.org/info/rfc7917>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.

- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8177] Lindem, A., Ed., Qu, Y., Yeung, D., Chen, I., and J. Zhang, "YANG Data Model for Key Chains", RFC 8177, DOI 10.17487/RFC8177, June 2017, <<https://www.rfc-editor.org/info/rfc8177>>.
- [RFC8294] Liu, X., Qu, Y., Lindem, A., Hopps, C., and L. Berger, "Common YANG Data Types for the Routing Area", RFC 8294, DOI 10.17487/RFC8294, December 2017, <<https://www.rfc-editor.org/info/rfc8294>>.
- [RFC8341] Bierman, A. and M. Bjorklund, "Network Configuration Access Control Model", STD 91, RFC 8341, DOI 10.17487/RFC8341, March 2018, <<https://www.rfc-editor.org/info/rfc8341>>.
- [RFC8342] Bjorklund, M., Schoenwaelder, J., Shafer, P., Watsen, K., and R. Wilton, "Network Management Datastore Architecture (NMDA)", RFC 8342, DOI 10.17487/RFC8342, March 2018, <<https://www.rfc-editor.org/info/rfc8342>>.
- [RFC8343] Bjorklund, M., "A YANG Data Model for Interface Management", RFC 8343, DOI 10.17487/RFC8343, March 2018, <<https://www.rfc-editor.org/info/rfc8343>>.
- [RFC8349] Lhotka, L., Lindem, A., and Y. Qu, "A YANG Data Model for Routing Management (NMDA Version)", RFC 8349, DOI 10.17487/RFC8349, March 2018, <<https://www.rfc-editor.org/info/rfc8349>>.
- [RFC8405] Decraene, B., Litkowski, S., Gredler, H., Lindem, A., Francois, P., and C. Bowers, "Shortest Path First (SPF) Back-Off Delay Algorithm for Link-State IGP", RFC 8405, DOI 10.17487/RFC8405, June 2018, <<https://www.rfc-editor.org/info/rfc8405>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.

- [RFC8570] Ginsberg, L., Ed., Previdi, S., Ed., Giacalone, S., Ward, D., Drake, J., and Q. Wu, "IS-IS Traffic Engineering (TE) Metric Extensions", RFC 8570, DOI 10.17487/RFC8570, March 2019, <<https://www.rfc-editor.org/info/rfc8570>>.

11.2. Informative References

- [I-D.ietf-rtgwg-segment-routing-ti-lfa]
Litkowski, S., Bashandy, A., Filsfils, C., Decraene, B., Francois, P., daniel.voyer@bell.ca, d., Clad, F., and P. Camarillo, "Topology Independent Fast Reroute using Segment Routing", draft-ietf-rtgwg-segment-routing-ti-lfa-01 (work in progress), March 2019.
- [RFC7812] Atlas, A., Bowers, C., and G. Enyedi, "An Architecture for IP/LDP Fast Reroute Using Maximally Redundant Trees (MRT-FRR)", RFC 7812, DOI 10.17487/RFC7812, June 2016, <<https://www.rfc-editor.org/info/rfc7812>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.

Appendix A. Example of IS-IS configuration in XML

This section gives an example of configuration of an IS-IS instance on a device. The example is written in XML.

```
<?xml version="1.0" encoding="utf-8"?>
<data xmlns="urn:ietf:params:xml:ns:netconf:base:1.0">
  <routing xmlns="urn:ietf:params:xml:ns:yang:ietf-routing">
    <name>SLI</name>
    <router-id>192.0.2.1</router-id>
    <control-plane-protocols>
      <control-plane-protocol>
        <name>ISIS-example</name>
        <description/>
        <type>
          <type xmlns:isis="urn:ietf:params:xml:ns:yang:ietf-isis">
            isis:isis
          </type>
        </type>
        <isis xmlns="urn:ietf:params:xml:ns:yang:ietf-isis">
          <enable>true</enable>
          <level-type>level-2</level-type>
          <system-id>87FC.FCDF.4432</system-id>
          <area-address>49.0001</area-address>
          <mpls>
```

```
    <te-rid>
      <ipv4-router-id>192.0.2.1</ipv4-router-id>
    </te-rid>
  </mpls>
  <lsp-lifetime>65535</lsp-lifetime>
  <lsp-refresh>65000</lsp-refresh>
  <metric-type>
    <value>wide-only</value>
  </metric-type>
  <default-metric>
    <value>111111</value>
  </default-metric>
  <address-families>
    <address-family-list>
      <address-family>ipv4</address-family>
      <enable>true</enable>
    </address-family-list>
    <address-family-list>
      <address-family>ipv6</address-family>
      <enable>true</enable>
    </address-family-list>
  </address-families>
  <interfaces>
    <interface>
      <name>Loopback0</name>
      <tag>200</tag>
      <metric>
        <value>0</value>
      </metric>
      <passive>true</passive>
    </interface>
    <interface>
      <name>Eth1</name>
      <level-type>level-2</level-type>
      <interface-type>point-to-point</interface-type>
      <metric>
        <value>167890</value>
      </metric>
    </interface>
  </interfaces>
</isis>
</control-plane-protocol>
</control-plane-protocols>
</routing>
<interfaces xmlns="urn:ietf:params:xml:ns:yang:ietf-interfaces">
  <interface>
    <name>Loopback0</name>
    <description/>
```

```
<type xmlns:ianaift="urn:ietf:params:xml:ns:yang:iana-if-type">
ianaift:softwareLoopback
</type>
<link-up-down-trap-enable>enabled</link-up-down-trap-enable>
<ipv4 xmlns="urn:ietf:params:xml:ns:yang:ietf-ip">
  <address>
    <ip>192.0.2.1</ip>
    <prefix-length>32</prefix-length>
  </address>
</ipv4>
<ipv6 xmlns="urn:ietf:params:xml:ns:yang:ietf-ip">
  <address>
    <ip>2001:DB8::1</ip>
    <prefix-length>128</prefix-length>
  </address>
</ipv6>
</interface>
<interface>
  <name>Eth1</name>
  <description/>
  <type xmlns:ianaift="urn:ietf:params:xml:ns:yang:iana-if-type">
ianaift:ethernetCsmacd
  </type>
  <link-up-down-trap-enable>enabled</link-up-down-trap-enable>
  <ipv4 xmlns="urn:ietf:params:xml:ns:yang:ietf-ip">
    <address>
      <ip>198.51.100.1</ip>
      <prefix-length>30</prefix-length>
    </address>
  </ipv4>
  <ipv6 xmlns="urn:ietf:params:xml:ns:yang:ietf-ip">
    <address>
      <ip>2001:DB8:0:0:FF::1</ip>
      <prefix-length>64</prefix-length>
    </address>
  </ipv6>
</interface>
</interfaces>
</data>
```

Authors' Addresses

Stephane Litkowski
Cisco Systems

Email: slitkows.ietf@gmail.com

Derek Yeung
Arrcus, Inc

Email: derek@arrcus.com

Acee Lindem
Cisco Systems

Email: acee@cisco.com

Jeffrey Zhang
Juniper Networks

Email: zzhang@juniper.net

Ladislav Lhotka
CZ.NIC

Email: lhotka@nic.cz

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: October 9, 2022

P. Psenak, Ed.
Cisco Systems
S. Hegde
Juniper Networks, Inc.
C. Filsfils
Cisco Systems, Inc.
K. Talaulikar
Arrcus, Inc
A. Gulko
Edward Jones
April 7, 2022

IGP Flexible Algorithm
draft-ietf-lsr-flex-algo-19

Abstract

IGP protocols traditionally compute best paths over the network based on the IGP metric assigned to the links. Many network deployments use RSVP-TE based or Segment Routing based Traffic Engineering to steer traffic over a path that is computed using different metrics or constraints than the shortest IGP path. This document proposes a solution that allows IGPs themselves to compute constraint-based paths over the network. This document also specifies a way of using Segment Routing (SR) Prefix-SIDs and SRv6 locators to steer packets along the constraint-based paths.

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 <https://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 October 9, 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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
2. Requirements Language	4
3. Terminology	4
4. Flexible Algorithm	5
5. Flexible Algorithm Definition Advertisement	6
5.1. IS-IS Flexible Algorithm Definition Sub-TLV	6
5.2. OSPF Flexible Algorithm Definition TLV	8
5.3. Common Handling of Flexible Algorithm Definition TLV	9
6. Sub-TLVs of IS-IS FAD Sub-TLV	10
6.1. IS-IS Flexible Algorithm Exclude Admin Group Sub-TLV	11
6.2. IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV	12
6.3. IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV	12
6.4. IS-IS Flexible Algorithm Definition Flags Sub-TLV	13
6.5. IS-IS Flexible Algorithm Exclude SRLG Sub-TLV	14
7. Sub-TLVs of OSPF FAD TLV	15
7.1. OSPF Flexible Algorithm Exclude Admin Group Sub-TLV	15
7.2. OSPF Flexible Algorithm Include-Any Admin Group Sub-TLV	16
7.3. OSPF Flexible Algorithm Include-All Admin Group Sub-TLV	16
7.4. OSPF Flexible Algorithm Definition Flags Sub-TLV	16
7.5. OSPF Flexible Algorithm Exclude SRLG Sub-TLV	17
8. IS-IS Flexible Algorithm Prefix Metric Sub-TLV	18
9. OSPF Flexible Algorithm Prefix Metric Sub-TLV	19
10. OSPF Flexible Algorithm ASBR Reachability Advertisement	21
10.1. OSPFv2 Extended Inter-Area ASBR LSA	21
10.1.1. OSPFv2 Extended Inter-Area ASBR TLV	23
10.2. OSPF Flexible Algorithm ASBR Metric Sub-TLV	23
11. Advertisement of Node Participation in a Flex-Algorithm	25
11.1. Advertisement of Node Participation for Segment Routing	26
11.2. Advertisement of Node Participation for Other Applications	26
12. Advertisement of Link Attributes for Flex-Algorithm	26

13. Calculation of Flexible Algorithm Paths	27
13.1. Multi-area and Multi-domain Considerations	29
14. Flex-Algorithm and Forwarding Plane	31
14.1. Segment Routing MPLS Forwarding for Flex-Algorithm	32
14.2. SRv6 Forwarding for Flex-Algorithm	32
14.3. Other Applications' Forwarding for Flex-Algorithm	33
15. Operational Considerations	33
15.1. Inter-area Considerations	33
15.2. Usage of SRLG Exclude Rule with Flex-Algorithm	34
15.3. Max-metric consideration	35
16. Backward Compatibility	35
17. Security Considerations	35
18. IANA Considerations	36
18.1. IGP IANA Considerations	36
18.1.1. IGP Algorithm Types Registry	36
18.1.2. IGP Metric-Type Registry	36
18.2. Flexible Algorithm Definition Flags Registry	37
18.3. IS-IS IANA Considerations	37
18.3.1. Sub TLVs for Type 242	37
18.3.2. Sub TLVs for for TLVs 135, 235, 236, and 237	37
18.3.3. Sub-Sub-TLVs for Flexible Algorithm Definition Sub-TLV	37
18.4. OSPF IANA Considerations	38
18.4.1. OSPF Router Information (RI) TLVs Registry	38
18.4.2. OSPFv2 Extended Prefix TLV Sub-TLVs	39
18.4.3. OSPFv3 Extended-LSA Sub-TLVs	39
18.4.4. OSPF Flex-Algorithm Prefix Metric Bits	39
18.4.5. OSPF Opaque LSA Option Types	39
18.4.6. OSPFv2 Extended Inter-Area ASBR TLVs	40
18.4.7. OSPFv2 Inter-Area ASBR Sub-TLVs	40
18.4.8. OSPF Flexible Algorithm Definition TLV Sub-TLV Registry	40
18.4.9. Link Attribute Applications Registry	42
19. Acknowledgements	42
20. References	42
20.1. Normative References	42
20.2. Informative References	44
Authors' Addresses	46

1. Introduction

An IGP-computed path based on the shortest IGP metric is often be replaced by a traffic-engineered path due to the traffic requirements which are not reflected by the IGP metric. Some networks engineer the IGP metric assignments in a way that the IGP metric reflects the link bandwidth or delay. If, for example, the IGP metric is reflecting the bandwidth on the link and the application traffic is

delay sensitive, the best IGP path may not reflect the best path from such an application's perspective.

To overcome this limitation, various sorts of traffic engineering have been deployed, including RSVP-TE and SR-TE, in which case the TE component is responsible for computing paths based on additional metrics and/or constraints. Such paths need to be installed in the forwarding tables in addition to, or as a replacement for, the original paths computed by IGPs. Tunnels are often used to represent the engineered paths and mechanisms like one described in [RFC3906] are used to replace the native IGP paths with such tunnel paths.

This document specifies a set of extensions to IS-IS, OSPFv2, and OSPFv3 that enable a router to advertise TLVs that (a) identify calculation-type, (b) specify a metric-type, and (c) describe a set of constraints on the topology, that are to be used to compute the best paths along the constrained topology. A given combination of calculation-type, metric-type, and constraints is known as a "Flexible Algorithm Definition". A router that sends such a set of TLVs also assigns a Flex-Algorithm value to the specified combination of calculation-type, metric-type, and constraints.

This document also specifies a way for a router to use IGPs to associate one or more SR Prefix-SIDs or SRv6 locators with a particular Flex-Algorithm. Each such Prefix-SID or SRv6 locator then represents a path that is computed according to the identified Flex-Algorithm.

2. Requirements Language

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. Terminology

This section defines terms that are often used in this document.

Flexible Algorithm Definition (FAD) - the set consisting of (a) calculation-type, (b) metric-type, and (c) a set of constraints.

Flexible Algorithm - a numeric identifier in the range 128-255 that is associated via configuration with the Flexible-Algorithm Definition.

Local Flexible Algorithm Definition - Flexible Algorithm Definition defined locally on the node.

Remote Flexible Algorithm Definition - Flexible Algorithm Definition received from other nodes via IGP flooding.

Flexible Algorithm Participation - per application configuration state that expresses whether the node is participating in a particular Flexible Algorithm.

IGP Algorithm - value from the the "IGP Algorithm Types" registry defined under "Interior Gateway Protocol (IGP) Parameters" IANA registries. IGP Algorithms represents the triplet (Calculation Type, Metric, Constraints), where the second and third elements of the triple MAY be unspecified.

ABR - Area Border Router. In IS-IS terminology it is also known as L1/L2 router.

ASBR - Autonomous System Border Router.

4. Flexible Algorithm

Many possible constraints may be used to compute a path over a network. Some networks are deployed as multiple planes. A simple form of constraint may be to use a particular plane. A more sophisticated form of constraint can include some extended metric as described in [RFC8570]. Constraints which restrict paths to links with specific affinities or avoid links with specific affinities are also possible. Combinations of these are also possible.

To provide maximum flexibility, we want to provide a mechanism that allows a router to (a) identify a particular calculation-type, (b) metric-type, (c) describe a particular set of constraints, and (d) assign a numeric identifier, referred to as Flex-Algorithm, to the combination of that calculation-type, metric-type, and those constraints. We want the mapping between the Flex-Algorithm and its meaning to be flexible and defined by the user. As long as all routers in the domain have a common understanding as to what a particular Flex-Algorithm represents, the resulting routing computation is consistent and traffic is not subject to any looping.

The set consisting of (a) calculation-type, (b) metric-type, and (c) a set of constraints is referred to as a Flexible-Algorithm Definition.

Flexible-Algorithm is a numeric identifier in the range 128-255 that is associated via configuratin with the Flexible-Algorithm Definition.

IANA "IGP Algorithm Types" registry defines the set of values for IGP Algorithms. We propose to allocate the following values for Flex-Algorithms from this registry:

128-255 - Flex-Algorithms

5. Flexible Algorithm Definition Advertisement

To guarantee the loop-free forwarding for paths computed for a particular Flex-Algorithm, all routers that (a) are configured to participate in a particular Flex-Algorithm, and (b) are in the same Flex-Algorithm definition advertisement scope MUST agree on the definition of the Flex-Algorithm.

5.1. IS-IS Flexible Algorithm Definition Sub-TLV

The IS-IS Flexible Algorithm Definition Sub-TLV (FAD Sub-TLV) is used to advertise the definition of the Flex-Algorithm.

The IS-IS FAD Sub-TLV is advertised as a Sub-TLV of the IS-IS Router Capability TLV-242 that is defined in [RFC7981].

IS-IS FAD Sub-TLV has the following format:

0	1	2	3
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		Length
+	+	+	+
	Calc-Type		Priority
+	+	+	+
	Sub-TLVs		
+	+		
	...		
+	+		
+	+		

where:

Type: 26

Length: variable, dependent on the included Sub-TLVs

Flex-Algorithm: Single octet value between 128 and 255 inclusive.

Metric-Type: Type of metric to be used during the calculation.
Following values are defined:

0: IGP Metric

1: Min Unidirectional Link Delay as defined in [RFC8570], section 4.2, encoded as application specific link attribute as specified in [RFC8919] and Section 12 of this document.

2: Traffic Engineering Default Metric as defined in [RFC5305], section 3.7, encoded as application specific link attribute as specified in [RFC8919] and Section 12 of this document.

Calc-Type: value from 0 to 127 inclusive from the "IGP Algorithm Types" registry defined under "Interior Gateway Protocol (IGP) Parameters" IANA registries. IGP algorithms in the range of 0-127 have a defined triplet (Calculation Type, Metric, Constraints). When used to specify the Calc-Type in the FAD Sub-TLV, only the Calculation Type defined for the specified IGP Algorithm is used. The Metric/Constraints MUST NOT be inherited. If the required calculation type is Shortest Path First, the value 0 SHOULD appear in this field.

Priority: Value between 0 and 255 inclusive that specifies the priority of the advertisement.

Sub-TLVs - optional sub-TLVs.

The IS-IS FAD Sub-TLV MAY be advertised in an LSP of any number. IS-IS router MAY advertise more than one IS-IS FAD Sub-TLV for a given Flexible-Algorithm (see Section 6).

The IS-IS FAD Sub-TLV has an area scope. The Router Capability TLV in which the FAD Sub-TLV is present MUST have the S-bit clear.

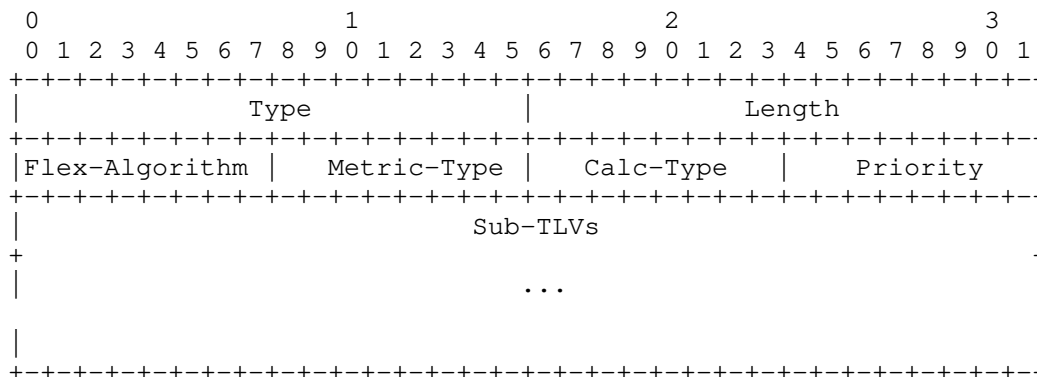
IS-IS L1/L2 router MAY be configured to re-generate the winning FAD from level 2, without any modification to it, to level 1 area. The re-generation of the FAD Sub-TLV from level 2 to level 1 is determined by the L1/L2 router, not by the originator of the FAD advertisement in the level 2. In such case, the re-generated FAD Sub-TLV will be advertised in the level 1 Router Capability TLV originated by the L1/L2 router.

L1/L2 router MUST NOT re-generate any FAD Sub-TLV from level 1 to level 2.

5.2. OSPF Flexible Algorithm Definition TLV

OSPF FAD TLV is advertised as a top-level TLV of the RI LSA that is defined in [RFC7770].

OSPF FAD TLV has the following format:



where:

Type: 16

Length: variable, dependent on the included Sub-TLVs

Flex-Algorithm:: Flex-Algorithm number. Value between 128 and 255 inclusive.

Metric-Type: Type of metric to be used during the calculation.
Following values are defined:

0: IGP Metric

1: Min Unidirectional Link Delay as defined in [RFC7471], section 4.2, encoded as application specific link attribute as specified in [RFC8920] and Section 12 of this document.

2: Traffic Engineering metric as defined in [RFC3630], section 2.5.5, encoded as application specific link attribute as specified in [RFC8920] and Section 12 of this document.

Calc-Type: as described in Section 5.1

Priority: as described in Section 5.1

Sub-TLVs - optional sub-TLVs.

When multiple OSPF FAD TLVs, for the same Flexible-Algorithm, are received from a given router, the receiver MUST use the first occurrence of the TLV in the Router Information LSA. If the OSPF FAD TLV, for the same Flex-Algorithm, appears in multiple Router Information LSAs that have different flooding scopes, the OSPF FAD TLV in the Router Information LSA with the area-scoped flooding scope MUST be used. If the OSPF FAD TLV, for the same algorithm, appears in multiple Router Information LSAs that have the same flooding scope, the OSPF FAD TLV in the Router Information (RI) LSA with the numerically smallest Instance ID MUST be used and subsequent instances of the OSPF FAD TLV MUST be ignored.

The RI LSA can be advertised at any of the defined opaque flooding scopes (link, area, or Autonomous System (AS)). For the purpose of OSPF FAD TLV advertisement, area-scoped flooding is REQUIRED. The Autonomous System flooding scope SHOULD NOT be used by default unless local configuration policy on the originating router indicates domain wide flooding.

5.3. Common Handling of Flexible Algorithm Definition TLV

This section describes the protocol-independent handling of the FAD TLV (OSPF) or FAD Sub-TLV (IS-IS). We will refer to it as FAD TLV in this section, even though in the case of IS-IS it is a Sub-TLV.

The value of the Flex-Algorithm MUST be between 128 and 255 inclusive. If it is not, the FAD TLV MUST be ignored.

Only a subset of the routers participating in the particular Flex-Algorithm need to advertise the definition of the Flex-Algorithm.

Every router, that is configured to participate in a particular Flex-Algorithm, MUST select the Flex-Algorithm definition based on the following ordered rules. This allows for the consistent Flex-Algorithm definition selection in cases where different routers advertise different definitions for a given Flex-Algorithm:

1. From the advertisements of the FAD in the area (including both locally generated advertisements and received advertisements) select the one(s) with the highest priority value.
2. If there are multiple advertisements of the FAD with the same highest priority, select the one that is originated from the router with the highest System-ID, in the case of IS-IS, or Router ID, in the case of OSPFv2 and OSPFv3. For IS-IS, the System-ID is

described in [ISO10589]. For OSPFv2 and OSPFv3, standard Router ID is described in [RFC2328] and [RFC5340] respectively.

A router that is not configured to participate in a particular Flex-Algorithm MUST ignore FAD Sub-TLVs advertisements for such Flex-Algorithm.

A router that is not participating in a particular Flex-Algorithm is allowed to advertise FAD for such Flex-Algorithm. Receiving routers MUST consider FAD advertisement regardless of the Flex-Algorithm participation of the FAD originator.

Any change in the Flex-Algorithm definition may result in temporary disruption of traffic that is forwarded based on such Flex-Algorithm paths. The impact is similar to any other event that requires network-wide convergence.

If a node is configured to participate in a particular Flexible-Algorithm, but there is no valid Flex-Algorithm definition available for it, or the selected Flex-Algorithm definition includes calculation-type, metric-type, constraint, flag, or Sub-TLV that is not supported by the node, it MUST stop participating in such Flexible-Algorithm. That implies that it MUST NOT announce participation for such Flexible-Algorithm as specified in Section 11 and it MUST remove any forwarding state associated with it.

Flex-Algorithm definition is topology independent. It applies to all topologies that a router participates in.

6. Sub-TLVs of IS-IS FAD Sub-TLV

One of the limitations of IS-IS [ISO10589] is that the length of a TLV/sub-TLV is limited to a maximum of 255 octets. For the FAD sub-TLV, there are a number of sub-sub-TLVs (defined below) which are supported. For a given Flex-Algorithm, it is possible that the total number of octets required to completely define a FAD exceeds the maximum length supported by a single FAD sub-TLV. In such cases, the FAD may be split into multiple such sub-TLVs and the content of the multiple FAD sub-TLVs combined to provide a complete FAD for the Flex-Algorithm. In such case, the fixed portion of the FAD (see Section 5.1) MUST be identical in all FAD sub-TLVs for a given Flex-Algorithm from a given IS. In case the fixed portion of such FAD Sub-TLVs differ, the values in the fixed portion in the FAD sub-TLV in the first occurrence in the lowest numbered LSP from a given IS MUST be used.

Any specification that introduces a new ISIS FAD sub-sub-TLV MUST specify whether the FAD sub-TLV may appear multiple times in the set

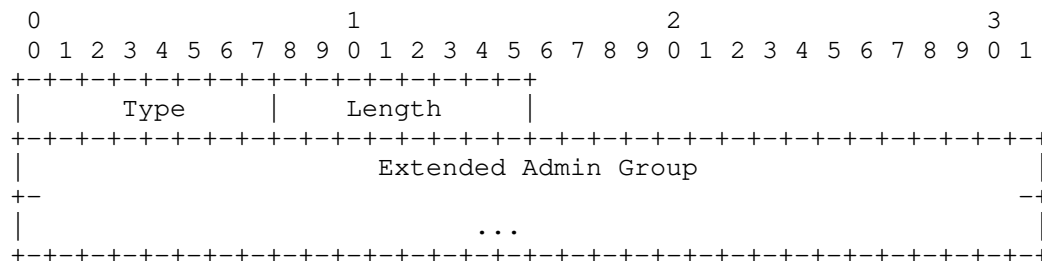
of FAD sub-TLVs for a given Flex-Algorithm from a given IS and how to handle them if multiple are allowed.

6.1. IS-IS Flexible Algorithm Exclude Admin Group Sub-TLV

The Flexible Algorithm definition can specify 'colors' that are used by the operator to exclude links during the Flex-Algorithm path computation.

The IS-IS Flexible Algorithm Exclude Admin Group Sub-TLV is used to advertise the exclude rule that is used during the Flex-Algorithm path calculation as specified in Section 13.

The IS-IS Flexible Algorithm Exclude Admin Group Sub-TLV (FAEAG Sub-TLV) is a Sub-TLV of the IS-IS FAD Sub-TLV. It has the following format:



where:

Type: 1

Length: variable, dependent on the size of the Extended Admin Group. MUST be a multiple of 4 octets.

Extended Administrative Group: Extended Administrative Group as defined in [RFC7308].

The IS-IS FAEAG Sub-TLV MUST NOT appear more than once in a single IS-IS FAD Sub-TLV. If it appears more than once, the IS-IS FAD Sub-TLV MUST be ignored by the receiver.

The IS-IS FAEAG Sub-TLV MUST NOT appear more than once in the set of FAD sub-TLVs for a given Flex-Algorithm from a given IS. If it appears more than once in such set, the IS-IS FAEAG Sub-TLV in the first occurrence in the lowest numbered LSP from a given IS MUST be used and any other occurrences MUST be ignored.

6.2. IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV

The Flexible Algorithm definition can specify 'colors' that are used by the operator to include links during the Flex-Algorithm path computation.

The IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV is used to advertise include-any rule that is used during the Flex-Algorithm path calculation as specified in Section 13.

The format of the IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV is identical to the format of the FAEAG Sub-TLV in Section 6.1.

The IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV Type is 2.

The IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV MUST NOT appear more than once in a single IS-IS FAD Sub-TLV. If it appears more than once, the IS-IS FAD Sub-TLV MUST be ignored by the receiver.

The IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV MUST NOT appear more than once in the set of FAD sub-TLVs for a given Flex-Algorithm from a given IS. If it appears more than once in such set, the IS-IS Flexible Algorithm Include-Any Admin Group Sub-TLV in the first occurrence in the lowest numbered LSP from a given IS MUST be used and any other occurrences MUST be ignored.

6.3. IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV

The Flexible Algorithm definition can specify 'colors' that are used by the operator to include link during the Flex-Algorithm path computation.

The IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV is used to advertise include-all rule that is used during the Flex-Algorithm path calculation as specified in Section 13.

The format of the IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV is identical to the format of the FAEAG Sub-TLV in Section 6.1.

The IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV Type is 3.

The IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV MUST NOT appear more than once in a single IS-IS FAD Sub-TLV. If it appears

more than once, the IS-IS FAD Sub-TLV MUST be ignored by the receiver.

The IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV MUST NOT appear more than once in the set of FAD sub-TLVs for a given Flex-Algorithm from a given IS. If it appears more than once in such set, the IS-IS Flexible Algorithm Include-All Admin Group Sub-TLV in the first occurrence in the lowest numbered LSP from a given IS MUST be used and any other occurrences MUST be ignored.

6.4. IS-IS Flexible Algorithm Definition Flags Sub-TLV

The IS-IS Flexible Algorithm Definition Flags Sub-TLV (FADF Sub-TLV) is a Sub-TLV of the IS-IS FAD Sub-TLV. It has the following format:

```

      0               1               2               3
      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          |      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                                     Flags                                     |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                                     ...                                     |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

where:

Type: 4

Length: variable, non-zero number of octets of the Flags field

Flags:

```

      0 1 2 3 4 5 6 7...
+---+---+---+---+---+---+---+
|M| | | | | | |...
+---+---+---+---+---+---+---+

```

M-flag: when set, the Flex-Algorithm specific prefix metric MUST be used for inter-area and external prefix calculation. This flag is not applicable to prefixes advertised as SRv6 locators.

Bits are defined/sent starting with Bit 0 defined above. Additional bit definitions that may be defined in the future SHOULD be assigned in ascending bit order so as to minimize the number of bits that will need to be transmitted.

Undefined bits MUST be transmitted as 0.

Bits that are NOT transmitted MUST be treated as if they are set to 0 on receipt.

The IS-IS FADF Sub-TLV MUST NOT appear more than once in a single IS-IS FAD Sub-TLV. If it appears more than once, the IS-IS FAD Sub-TLV MUST be ignored by the receiver.

The IS-IS FADF Sub-TLV MUST NOT appear more than once in the set of FAD sub-TLVs for a given Flex-Algorithm from a given IS. If it appears more than once in such set, the IS-IS FADF Sub-TLV in the first occurrence in the lowest numbered LSP from a given IS MUST be used and any other occurrences MUST be ignored.

If the IS-IS FADF Sub-TLV is not present inside the IS-IS FAD Sub-TLV, all the bits are assumed to be set to 0.

If a node is configured to participate in a particular Flexible-Algorithm, but the selected Flex-Algorithm definition includes a bit in the IS-IS FADF Sub-TLV that is not supported by the node, it MUST stop participating in such Flexible-Algorithm.

New flag bits may be defined in the future. Implementations MUST check all advertised flag bits in the received IS-IS FADF Sub-TLV - not just the subset currently defined.

6.5. IS-IS Flexible Algorithm Exclude SRLG Sub-TLV

The Flexible Algorithm definition can specify Shared Risk Link Groups (SRLGs) that the operator wants to exclude during the Flex-Algorithm path computation.

The IS-IS Flexible Algorithm Exclude SRLG Sub-TLV (FAESRLG) is used to advertise the exclude rule that is used during the Flex-Algorithm path calculation as specified in Section 13.

The IS-IS FAESRLG Sub-TLV is a Sub-TLV of the IS-IS FAD Sub-TLV. It has the following format:

```

      0               1               2               3
      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          |      Length      |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Shared Risk Link Group Value                               |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                                                                                               ...                               |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

where:

Type: 5

Length: variable, dependent on number of SRLG values. MUST be a multiple of 4 octets.

Shared Risk Link Group Value: SRLG value as defined in [RFC5307].

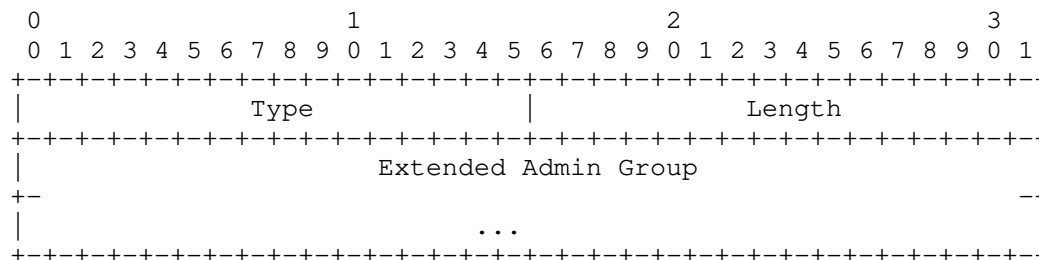
The IS-IS FAESRLG Sub-TLV MUST NOT appear more than once in a single IS-IS FAD Sub-TLV. If it appears more than once, the IS-IS FAD Sub-TLV MUST be ignored by the receiver.

The IS-IS FAESRLG Sub-TLV MAY appear more than once in the set of FAD sub-TLVs for a given Flex-Algorithm from a given IS. This may be necessary in cases where the total number of SRLG values which are specified cause the FAD sub-TLV to exceed the maximum length of a single FAD sub-TLV. In such case the receiver MUST use the union of all values across all IS-IS FAESRLG Sub-TLVs from such set.

7. Sub-TLVs of OSPF FAD TLV

7.1. OSPF Flexible Algorithm Exclude Admin Group Sub-TLV

The Flexible Algorithm Exclude Admin Group Sub-TLV (FAEAG Sub-TLV) is a Sub-TLV of the OSPF FAD TLV. It's usage is described in Section 6.1. It has the following format:



where:

Type: 1

Length: variable, dependent on the size of the Extended Admin Group. MUST be a multiple of 4 octets.

Extended Administrative Group: Extended Administrative Group as defined in [RFC7308].

The OSPF FAEAG Sub-TLV MUST NOT appear more than once in an OSPF FAD TLV. If it appears more than once, the OSPF FAD TLV MUST be ignored by the receiver.

7.2. OSPF Flexible Algorithm Include-Any Admin Group Sub-TLV

The usage of this Sub-TLVs is described in Section 6.2.

The format of the OSPF Flexible Algorithm Include-Any Admin Group Sub-TLV is identical to the format of the OSPF FAEAG Sub-TLV in Section 7.1.

The OSPF Flexible Algorithm Include-Any Admin Group Sub-TLV Type is 2.

The OSPF Flexible Algorithm Include-Any Admin Group Sub-TLV MUST NOT appear more than once in an OSPF FAD TLV. If it appears more than once, the OSPF FAD TLV MUST be ignored by the receiver.

7.3. OSPF Flexible Algorithm Include-All Admin Group Sub-TLV

The usage of this Sub-TLVs is described in Section 6.3.

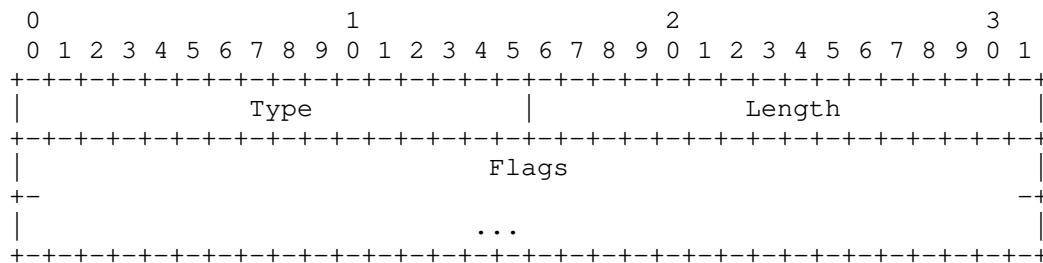
The format of the OSPF Flexible Algorithm Include-All Admin Group Sub-TLV is identical to the format of the OSPF FAEAG Sub-TLV in Section 7.1.

The OSPF Flexible Algorithm Include-All Admin Group Sub-TLV Type is 3.

The OSPF Flexible Algorithm Include-All Admin Group Sub-TLV MUST NOT appear more than once in an OSPF FAD TLV. If it appears more than once, the OSPF FAD TLV MUST be ignored by the receiver.

7.4. OSPF Flexible Algorithm Definition Flags Sub-TLV

The OSPF Flexible Algorithm Definition Flags Sub-TLV (FADF Sub-TLV) is a Sub-TLV of the OSPF FAD TLV. It has the following format:



where:

Type: 4

Length: variable, dependent on the size of the Flags field. MUST be a multiple of 4 octets.

Flags:

```

    0 1 2 3 4 5 6 7...
    +-+-+-+-+-+-+-+...
    |M| | |         ...
    +-+-+-+-+-+-+-+...

```

M-flag: when set, the Flex-Algorithm specific prefix and ASBR metric MUST be used for inter-area and external prefix calculation. This flag is not applicable to prefixes advertised as SRv6 locators.

Bits are defined/sent starting with Bit 0 defined above. Additional bit definitions that may be defined in the future SHOULD be assigned in ascending bit order so as to minimize the number of bits that will need to be transmitted.

Undefined bits MUST be transmitted as 0.

Bits that are NOT transmitted MUST be treated as if they are set to 0 on receipt.

The OSPF FADF Sub-TLV MUST NOT appear more than once in an OSPF FAD TLV. If it appears more than once, the OSPF FAD TLV MUST be ignored by the receiver.

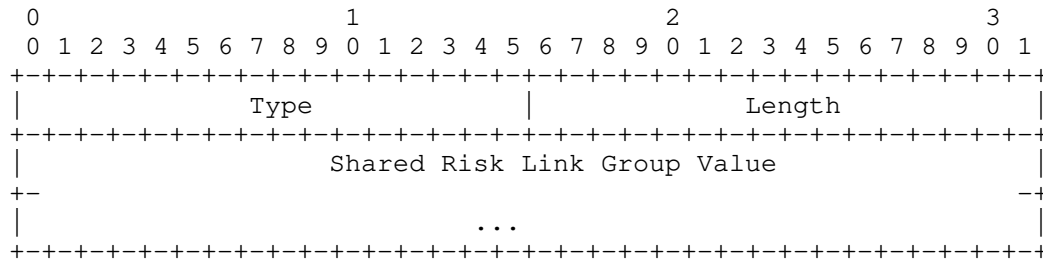
If the OSPF FADF Sub-TLV is not present inside the OSPF FAD TLV, all the bits are assumed to be set to 0.

If a node is configured to participate in a particular Flexible-Algorithm, but the selected Flex-Algorithm definition includes a bit in the OSPF FADF Sub-TLV that is not supported by the node, it MUST stop participating in such Flexible-Algorithm.

New flag bits may be defined in the future. Implementations MUST check all advertised flag bits in the received OSPF FADF Sub-TLV - not just the subset currently defined.

7.5. OSPF Flexible Algorithm Exclude SRLG Sub-TLV

The OSPF Flexible Algorithm Exclude SRLG Sub-TLV (FAESRLG Sub-TLV) is a Sub-TLV of the OSPF FAD TLV. Its usage is described in Section 6.5. It has the following format:



where:

Type: 5

Length: variable, dependent on the number of SRLGs. MUST be a multiple of 4 octets.

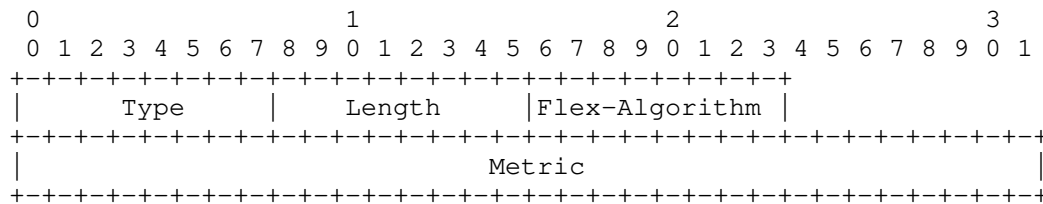
Shared Risk Link Group Value: SRLG value as defined in [RFC4203].

The OSPF FAESRLG Sub-TLV MUST NOT appear more than once in an OSPF FAD TLV. If it appears more than once, the OSPF FAD TLV MUST be ignored by the receiver.

8. IS-IS Flexible Algorithm Prefix Metric Sub-TLV

The IS-IS Flexible Algorithm Prefix Metric (FAPM) Sub-TLV supports the advertisement of a Flex-Algorithm specific prefix metric associated with a given prefix advertisement.

The IS-IS FAPM Sub-TLV is a sub-TLV of TLVs 135, 235, 236, and 237 and has the following format:



where:

Type: 6

Length: 5 octets

Flex-Algorithm: Single octet value between 128 and 255 inclusive.

Metric: 4 octets of metric information

The IS-IS FAPM Sub-TLV MAY appear multiple times in its parent TLV. If it appears more than once with the same Flex-Algorithm value, the first instance MUST be used and any subsequent instances MUST be ignored.

If a prefix is advertised with a Flex-Algorithm prefix metric larger than MAX_PATH_METRIC as defined in [RFC5305] this prefix MUST NOT be considered during the Flexible-Algorithm computation.

The usage of the Flex-Algorithm prefix metric is described in Section 13.

The IS-IS FAPM Sub-TLV MUST NOT be advertised as a sub-TLV of the IS-IS SRv6 Locator TLV [I-D.ietf-lsr-isis-srv6-extensions]. The IS-IS SRv6 Locator TLV includes the Algorithm and Metric fields which MUST be used instead. If the FAPM Sub-TLV is present as a sub-TLV of the IS-IS SRv6 Locator TLV in the received LSP, such FAPM Sub-TLV MUST be ignored.

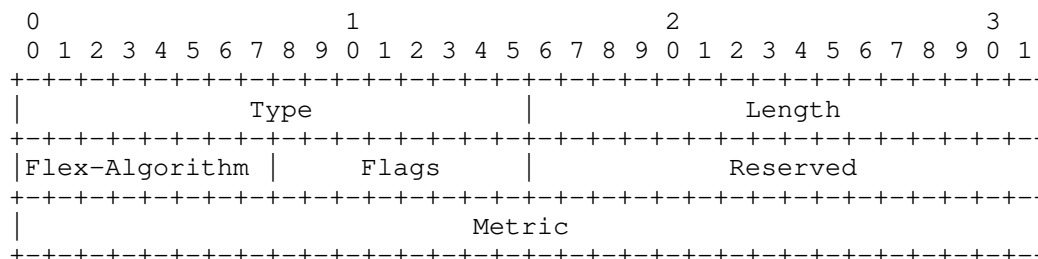
9. OSPF Flexible Algorithm Prefix Metric Sub-TLV

The OSPF Flexible Algorithm Prefix Metric (FAPM) Sub-TLV supports the advertisement of a Flex-Algorithm specific prefix metric associated with a given prefix advertisement.

The OSPF Flex-Algorithm Prefix Metric (FAPM) Sub-TLV is a Sub-TLV of the:

- OSPFv2 Extended Prefix TLV [RFC7684]
- Following OSPFv3 TLVs as defined in [RFC8362]:
 - Inter-Area Prefix TLV
 - External Prefix TLV

OSPF FAPM Sub-TLV has the following format:



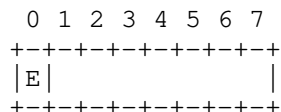
where:

Type: 3 for OSPFv2, 26 for OSPFv3

Length: 8 octets

Flex-Algorithm: Single octet value between 128 and 255 inclusive.

Flags: single octet value



E bit : position 0: The type of external metric. If bit is set, the metric specified is a Type 2 external metric. This bit is applicable only to OSPF External and NSSA external prefixes. This is semantically the same as E bit in section A.4.5 of [RFC2328] and section A.4.7 of [RFC5340] for OSPFv2 and OSPFv3 respectively.

Bits 1 through 7: MUST be cleared by sender and ignored by receiver.

Reserved: Must be set to 0, ignored at reception.

Metric: 4 octets of metric information

The OSPF FAPM Sub-TLV MAY appear multiple times in its parent TLV. If it appears more than once with the same Flex-Algorithm value, the first instance MUST be used and any subsequent instances MUST be ignored.

The usage of the Flex-Algorithm prefix metric is described in Section 13.

10. OSPF Flexible Algorithm ASBR Reachability Advertisement

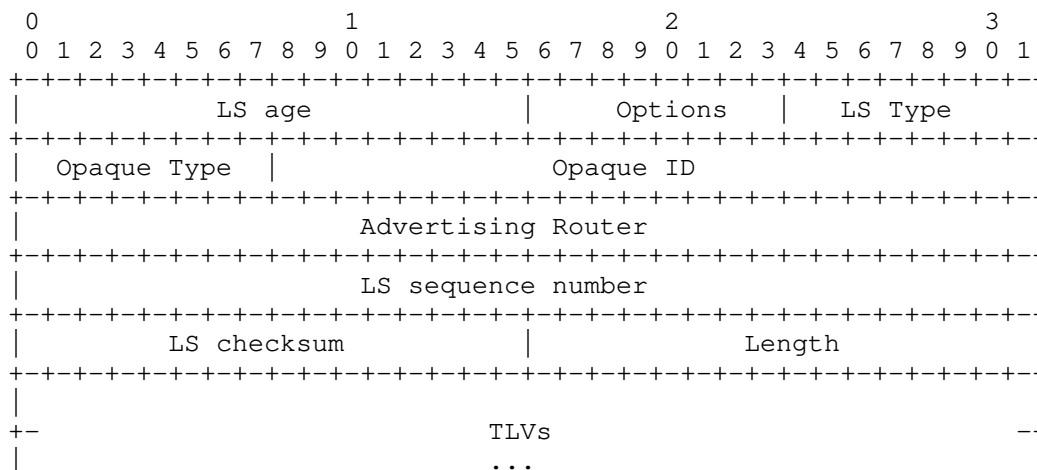
An OSPF ABR advertises the reachability of ASBRs in its attached areas to enable routers within those areas to perform route calculations for external prefixes advertised by the ASBRs. OSPF extensions for advertisement of Flex-Algorithm specific reachability and metric for ASBRs is similarly required for Flex-Algorithm external prefix computations as described further in Section 13.1.

10.1. OSPFv2 Extended Inter-Area ASBR LSA

The OSPFv2 Extended Inter-Area ASBR (EIA-ASBR) LSA is an OSPF Opaque LSA [RFC5250] that is used to advertise additional attributes related to the reachability of the OSPFv2 ASBR that is external to the area yet internal to the OSPF domain. Semantically, the OSPFv2 EIA-ASBR LSA is equivalent to the fixed format Type 4 Summary LSA [RFC2328]. Unlike the Type 4 Summary LSA, the LSID of the EIA-ASBR LSA does not carry the ASBR Router-ID - the ASBR Router-ID is carried in the body of the LSA. OSPFv2 EIA-ASBR LSA is advertised by an OSPFv2 ABR and its flooding is defined to be area-scoped only.

An OSPFv2 ABR generates the EIA-ASBR LSA for an ASBR when it is advertising the Type-4 Summary LSA for it and has the need for advertising additional attributes for that ASBR beyond what is conveyed in the fixed format Type-4 Summary LSA. An OSPFv2 ABR MUST NOT advertise the EIA-ASBR LSA for an ASBR for which it is not advertising the Type 4 Summary LSA. This ensures that the ABR does not generate the EIA-ASBR LSA for an ASBR to which it does not have reachability in the base OSPFv2 topology calculation. The OSPFv2 ABR SHOULD NOT advertise the EIA-ASBR LSA for an ASBR when it does not have additional attributes to advertise for that ASBR.

The OSPFv2 EIA-ASBR LSA has the following format:



The Opaque Type used by the OSPFv2 EIA-ASBR LSA is TBD (suggested value 11). The Opaque Type is used to differentiate the various types of OSPFv2 Opaque LSAs and is described in Section 3 of [RFC5250]. The LS Type MUST be 10, indicating that the Opaque LSA flooding scope is area-local [RFC5250]. The LSA Length field [RFC2328] represents the total length (in octets) of the Opaque LSA, including the LSA header and all TLVs (including padding).

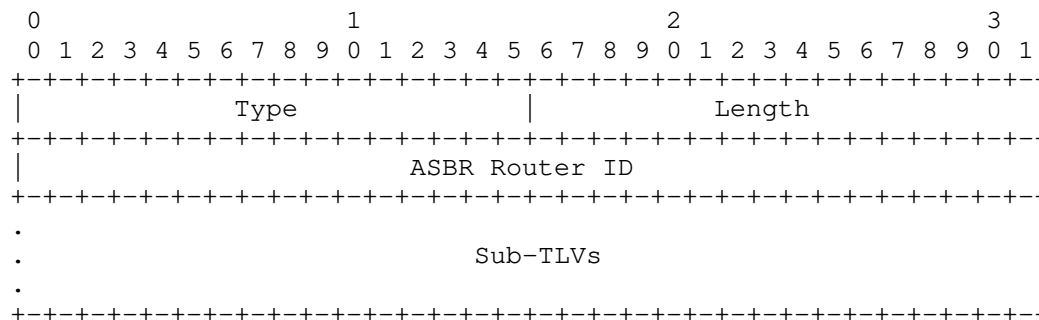
The Opaque ID field is an arbitrary value used to maintain multiple OSPFv2 EIA-ASBR LSAs. For OSPFv2 EIA-ASBR LSAs, the Opaque ID has no semantic significance other than to differentiate OSPFv2 EIA-ASBR LSAs originated by the same OSPFv2 ABR. If multiple OSPFv2 EIA-ASBR LSAs specify the same ASBR, the attributes from the Opaque LSA with the lowest Opaque ID SHOULD be used.

The format of the TLVs within the body of the OSPFv2 EIA-ASBR LSA is the same as the format used by the Traffic Engineering Extensions to OSPFv2 [RFC3630]. The variable TLV section consists of one or more nested TLV tuples. Nested TLVs are also referred to as sub-TLVs. The Length field defines the length of the value portion in octets (thus, a TLV with no value portion would have a length of 0). The TLV is padded to 4-octet alignment; padding is not included in the Length field (so a 3-octet value would have a length of 3, but the total size of the TLV would be 8 octets). Nested TLVs are also 32-bit aligned. For example, a 1-byte value would have the Length field set to 1, and 3 octets of padding would be added to the end of the value portion of the TLV. The padding is composed of zeros.

10.1.1. OSPFv2 Extended Inter-Area ASBR TLV

The OSPFv2 Extended Inter-Area ASBR (EIA-ASBR) TLV is a top-level TLV of the OSPFv2 EIA-ASBR LSA and is used to advertise additional attributes associated with the reachability of an ASBR.

The OSPFv2 EIA-ASBR TLV has the following format:



where:

Type: 1

Length: variable

ASBR Router ID: four octets carrying the OSPF Router ID of the ASBR whose information is being carried.

Sub-TLVs : variable

Only a single OSPFv2 EIA-ASBR TLV MUST be advertised in each OSPFv2 EIA-ASBR LSA and the receiver MUST ignore all instances of this TLV other than the first one in an LSA.

OSPFv2 EIA-ASBR TLV MUST be present inside an OSPFv2 EIA-ASBR LSA with at least a single sub-TLV included, otherwise the OSPFv2 EIA-ASBR LSA MUST be ignored by the receiver.

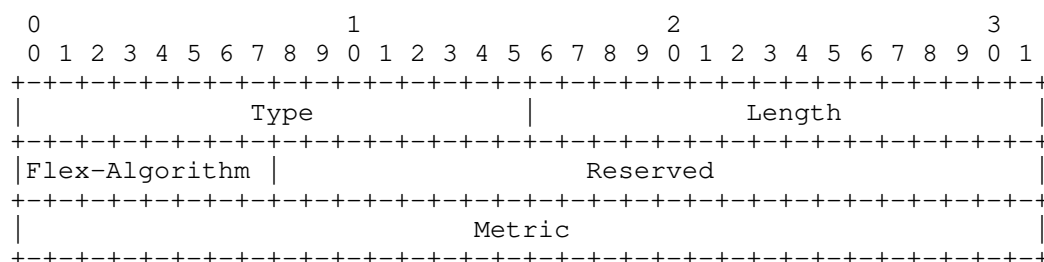
10.2. OSPF Flexible Algorithm ASBR Metric Sub-TLV

The OSPF Flexible Algorithm ASBR Metric (FAAM) Sub-TLV supports the advertisement of a Flex-Algorithm specific metric associated with a given ASBR reachability advertisement by an ABR.

The OSPF Flex-Algorithm ASBR Metric (FAAM) Sub-TLV is a Sub-TLV of the:

- OSPFv2 Extended Inter-Area ASBR TLV as defined in Section 10.1.1
- OSPFv3 Inter-Area-Router TLV defined in [RFC8362]

OSPF FAAM Sub-TLV has the following format:



where:

Type: 1 for OSPFv2, TBD (suggested value 30) for OSPFv3

Length: 8 octets

Flex-Algorithm: Single octet value between 128 and 255 inclusive.

Reserved: Must be set to 0, ignored at reception.

Metric: 4 octets of metric information

The OSPF FAAM Sub-TLV MAY appear multiple times in its parent TLV. If it appears more than once with the same Flex-Algorithm value, the first instance MUST be used and any subsequent instances MUST be ignored.

The advertisement of the ASBR reachability using the OSPF FAAM Sub-TLV inside the OSPFv2 EIA-ASBR LSA follows the section 12.4.3 of [RFC2328] and inside the OSPFv3 E-Inter-Area-Router LSA follows the section 4.8.5 of [RFC5340]. The reachability of the ASBR is evaluated in the context of the specific Flex-Algorithm.

The FAAM computed by the ABR will be equal to the metric to reach the ASBR for a given Flex-Algorithm in a source area or the cumulative metric via other ABR(s) when the ASBR is in a remote area. This is similar in nature to how the metric is set when the ASBR reachability metric is computed in the default algorithm for the metric in the OSPFv2 Type 4 ASBR Summary LSA and the OSPFv3 Inter-Area-Router LSA.

An OSPF ABR MUST NOT include the OSPF FAAM Sub-TLV with a specific Flex-Algorithm in its reachability advertisement for an ASBR between

areas unless that ASBR is reachable for it in the context of that specific Flex-Algorithm.

An OSPF ABR MUST include the OSPF FAAM Sub-TLVs as part of the ASBR reachability advertisement between areas for the Flex-Algorithm for which the winning FAD includes the M-flag and the ASBR is reachable in the context of that specific Flex-Algorithm.

OSPF routers MUST use the OSPF FAAM Sub-TLV to calculate the reachability of the ASBRs if the winning FAD for the specific Flex-Algorithm includes the M-flag. OSPF routers MUST NOT use the OSPF FAAM Sub-TLV to calculate the reachability of the ASBRs for the specific Flex-Algorithm if the winning FAD for such Flex-Algorithm does not include the M-flag. Instead, the OSPFv2 Type 4 Summary LSAs or the OSPFv3 Inter-Area-Router-LSAs MUST be used instead as specified in section 16.2 of [RFC2328] and section 4.8.5 of [RFC5340] for OSPFv2 and OSPFv3 respectively.

The processing of the new or changed OSPF FAAM Sub-TLV triggers the processing of the External routes similar to what is described in section 16.5 of the [RFC2328] for OSPFv2 and section 4.8.5 of [RFC5340] for OSPFv3 for the specific Flex-Algorithm. The External and NSSA External route calculation should be limited to Flex-Algorithm(s) for which the winning FAD(s) includes the M-flag.

Processing of the OSPF FAAM Sub-TLV does not require the existence of the equivalent OSPFv2 Type 4 Summary LSA or the OSPFv3 Inter-Area-Router-LSA that is advertised by the same ABR inside the area. When the OSPFv2 EIA-ASBR LSA or the OSPFv3 E-Inter-Area-Router-LSA are advertised along with the OSPF FAAM Sub-TLV by the ABR for a specific ASBR, it is expected that the same ABR would advertise the reachability of the same ASBR in the equivalent base LSAs - i.e., the OSPFv2 Type 4 Summary LSA or the OSPFv3 Inter-Area-Router-LSA. The presence of the base LSA is not mandatory for the usage of the extended LSA with the OSPF FAAM Sub-TLV. This means that the order in which these LSAs are received is not significant.

11. Advertisement of Node Participation in a Flex-Algorithm

When a router is configured to support a particular Flex-Algorithm, we say it is participating in that Flex-Algorithm.

Paths computed for a specific Flex-Algorithm MAY be used by various applications, each potentially using its own specific data plane for forwarding traffic over such paths. To guarantee the presence of the application specific forwarding state associated with a particular Flex-Algorithm, a router MUST advertise its participation for a particular Flex-Algorithm for each application specifically.

11.1. Advertisement of Node Participation for Segment Routing

[RFC8667], [RFC8665], and [RFC8666] (IGP Segment Routing extensions) describe how the SR-Algorithm is used to compute the IGP best path.

Routers advertise the support for the SR-Algorithm as a node capability as described in the above mentioned IGP Segment Routing extensions. To advertise participation for a particular Flex-Algorithm for Segment Routing, including both SR MPLS and SRv6, the Flex-Algorithm value MUST be advertised in the SR-Algorithm TLV (OSPF) or sub-TLV (IS-IS).

Segment Routing Flex-Algorithm participation advertisement is topology independent. When a router advertises participation in an SR-Algorithm, the participation applies to all topologies in which the advertising node participates.

11.2. Advertisement of Node Participation for Other Applications

This section describes considerations related to how other applications can advertise their participation in a specific Flex-Algorithm.

Application-specific Flex-Algorithm participation advertisements MAY be topology specific or MAY be topology independent, depending on the application itself.

Application-specific advertisement for Flex-Algorithm participation MUST be defined for each application and is outside of the scope of this document.

12. Advertisement of Link Attributes for Flex-Algorithm

Various link attributes may be used during the Flex-Algorithm path calculation. For example, include or exclude rules based on link affinities can be part of the Flex-Algorithm definition as defined in Section 6 and Section 7.

Application-specific link attributes, as specified in [RFC8919] or [RFC8920], that are to be used during Flex-Algorithm calculation MUST use the Application-Specific Link Attribute (ASLA) advertisements defined in [RFC8919] or [RFC8920], unless, in the case of IS-IS, the L-Flag is set in the ASLA advertisement. When the L-Flag is set, then legacy advertisements are to be used, subject to the procedures and constraints defined in [[RFC8919] Section 4.2 and Section 6.

The mandatory use of ASLA advertisements applies to link attributes specifically mentioned in this document (Min Unidirectional Link

Delay, TE Default Metric, Administrative Group, Extended Administrative Group and Shared Risk Link Group) and any other link attributes that may be used in support of Flex-Algorithm in the future.

A new Application Identifier Bit is defined to indicate that the ASLA advertisement is associated with the Flex-Algorithm application. This bit is set in the Standard Application Bit Mask (SABM) defined in [RFC8919] or [RFC8920]:

Bit-3: Flexible Algorithm (X-bit)

ASLA Admin Group Advertisements to be used by the Flexible Algorithm Application MAY use either the Administrative Group or Extended Administrative Group encodings. If the Administrative Group encoding is used, then the first 32 bits of the corresponding FAD sub-TLVs are mapped to the link attribute advertisements as specified in RFC 7308.

A receiver supporting this specification MUST accept both ASLA Administrative Group and Extended Administrative Group TLVs as defined in [RFC8919] or [RFC8920]. In the case of ISIS, if the L-Flag is set in ASLA advertisement, as defined in [RFC8919] Section 4.2, then the receiver MUST be able to accept both Administrative Group TLV as defined in [RFC5305] and Extended Administrative Group TLV as defined in [RFC7308].

13. Calculation of Flexible Algorithm Paths

A router MUST be configured to participate in a given Flex-Algorithm K and MUST select the FAD based on the rules defined in Section 5.3 before it can compute any path for that Flex-Algorithm.

No specific two way connectivity check is performed during the Flex-Algorithm path computation. The result of the existing, Flex-Algorithm agnostic, two way connectivity check is used during the Flex-Algorithm path computation.

As described in Section 11, participation for any particular Flex-Algorithm MUST be advertised on a per-application basis. Calculation of the paths for any particular Flex-Algorithm MUST be application specific.

The way applications handle nodes that do not participate in Flexible-Algorithm is application specific. If the application only wants to consider participating nodes during the Flex-Algorithm calculation, then when computing paths for a given Flex-Algorithm, all nodes that do not advertise participation for that Flex-Algorithm in their application-specific advertisements MUST be pruned from the

topology. Segment Routing, including both SR MPLS and SRv6, are applications that MUST use such pruning when computing Flex-Algorithm paths.

When computing the path for a given Flex-Algorithm, the metric-type that is part of the Flex-Algorithm definition (Section 5) MUST be used.

When computing the path for a given Flex-Algorithm, the calculation-type that is part of the Flex-Algorithm definition (Section 5) MUST be used.

Various link include or exclude rules can be part of the Flex-Algorithm definition. To refer to a particular bit within an AG or EAG we use the term 'color'.

Rules, in the order as specified below, MUST be used to prune links from the topology during the Flex-Algorithm computation.

For all links in the topology:

1. Check if any exclude AG rule is part of the Flex-Algorithm definition. If such exclude rule exists, check if any color that is part of the exclude rule is also set on the link. If such a color is set, the link MUST be pruned from the computation.
2. Check if any exclude SRLG rule is part of the Flex-Algorithm definition. If such exclude rule exists, check if the link is part of any SRLG that is also part of the SRLG exclude rule. If the link is part of such SRLG, the link MUST be pruned from the computation.
3. Check if any include-any AG rule is part of the Flex-Algorithm definition. If such include-any rule exists, check if any color that is part of the include-any rule is also set on the link. If no such color is set, the link MUST be pruned from the computation.
4. Check if any include-all AG rule is part of the Flex-Algorithm definition. If such include-all rule exists, check if all colors that are part of the include-all rule are also set on the link. If all such colors are not set on the link, the link MUST be pruned from the computation.
5. If the Flex-Algorithm definition uses other than IGP metric (Section 5), and such metric is not advertised for the particular link in a topology for which the computation is done, such link

MUST be pruned from the computation. A metric of value 0 MUST NOT be assumed in such case.

13.1. Multi-area and Multi-domain Considerations

Any IGP Shortest Path Tree calculation is limited to a single area. This applies to Flex-Algorithm calculations as well. Given that the computing router does not have visibility of the topology of the next areas or domain, the Flex-Algorithm specific path to an inter-area or inter-domain prefix will be computed for the local area only. The egress L1/L2 router (ABR in OSPF), or ASBR for inter-domain case, will be selected based on the best path for the given Flex-Algorithm in the local area and such egress ABR or ASBR router will be responsible to compute the best Flex-Algorithm specific path over the next area or domain. This may produce an end-to-end path, which is sub-optimal based on Flex-Algorithm constraints. In cases where the ABR or ASBR has no reachability to a prefix for a given Flex-Algorithm in the next area or domain, the traffic may be dropped by the ABR/ASBR.

To allow the optimal end-to-end path for an inter-area or inter-domain prefix for any Flex-Algorithm to be computed, the FAPM has been defined in Section 8 and Section 9. For external route calculation for prefixes originated by ASBRs in remote areas in OSPF, the FAAM has been defined in Section 10.2 for the ABR to indicate its ASBR reachability along with the metric for the specific Flex-Algorithm.

If the FAD selected based on the rules defined in Section 5.3 includes the M-flag, an ABR or ASBR MUST include the FAPM (Section 8, Section 9) when advertising the prefix, that is reachable in a given Flex-Algorithm, between areas or domains. Such metric will be equal to the metric to reach the prefix for that Flex-Algorithm in its source area or domain. This is similar in nature to how the metric is set when prefixes are advertised between areas or domains for the default algorithm. When a prefix is unreachable in its source area or domain in a specific Flex-Algorithm, then an ABR or ASBR MUST NOT include the FAPM for that Flex-Algorithm when advertising the prefix between areas or domains.

If the FAD selected based on the rules defined in Section 5.3 includes the M-flag, the FAPM MUST be used during the calculation of prefix reachability for the inter-area and external prefixes. If the FAPM for the Flex-Algorithm is not advertised with the inter-area or external prefix reachability advertisement, the prefix MUST be considered as unreachable for that Flex-Algorithm. Similarly in the case of OSPF, for ASBRs in remote areas, if the FAAM is not advertised by the local ABR(s), the ASBR MUST be considered as

unreachable for that Flex-Algorithm and the external prefix advertisements from such an ASBR are not considered for that Flex-Algorithm.

Flex-Algorithm prefix metrics and the OSPF Flex-Algorithm ASBR metrics MUST NOT be used during the Flex-Algorithm computation unless the FAD selected based on the rules defined in Section 5.3 includes the M-Flag, as described in (Section 6.4 or Section 7.4).

In the case of OSPF, when calculating external routes in a Flex-Algorithm (with FAD selected includes the M-Flag) where the advertising ASBR is in a remote area, the metric will be the sum of the following:

- o the FAPM for that Flex-Algorithm advertised with the external route by the ASBR
- o the metric to reach the ASBR for that Flex-Algorithm from the local ABR i.e., the FAAM for that Flex-Algorithm advertised by the ABR in the local area for that ASBR
- o the Flex-Algorithm specific metric to reach the local ABR

This is similar in nature to how the metric is calculated for routes learned from remote ASBRs in the default algorithm using the OSPFv2 Type 4 ASBR Summary LSA and the OSPFv3 Inter-Area-Router LSA.

If the FAD selected based on the rules defined in Section 5.3 does not include the M-flag, then the IGP metrics associated with the prefix reachability advertisements used by the base IS-IS and OSPF protocol MUST be used for the Flex-Algorithm route computation. Similarly, in the case of external route calculations in OSPF, the ASBR reachability is determined based on the base OSPFv2 Type 4 Summary LSA and the OSPFv3 Inter-Area-Router LSA.

It is NOT RECOMMENDED to use the Flex-Algorithm for inter-area or inter-domain prefix reachability without the M-flag set. The reason is that without the explicit Flex-Algorithm Prefix Metric advertisement (and the Flex-Algorithm ASBR metric advertisement in the case of OSPF external route calculation), it is not possible to conclude whether the ABR or ASBR has reachability to the inter-area or inter-domain prefix for a given Flex-Algorithm in the next area or domain. Sending the Flex-Algorithm traffic for such prefix towards the ABR or ASBR may result in traffic looping or black-holing.

During the route computation, it is possible for the Flex-Algorithm specific metric to exceed the maximum value that can be stored in an unsigned 32-bit variable. In such scenarios, the value MUST be

considered to be of value 4,294,967,295 during the computation and advertised as such.

The FAPM MUST NOT be advertised with IS-IS L1 or L2 intra-area, OSPFv2 intra-area, or OSPFv3 intra-area routes. If the FAPM is advertised for these route-types, it MUST be ignored during the prefix reachability calculation.

The M-flag in FAD is not applicable to prefixes advertised as SRv6 locators. The IS-IS SRv6 Locator TLV [I-D.ietf-lsr-isis-srv6-extensions] includes the Algorithm and Metric fields. When the SRv6 Locator is advertised between areas or domains, the metric field in the Locator TLV of IS-IS MUST be used irrespective of the M-flag in the FAD advertisement.

OSPF external and NSSA external prefix advertisements MAY include a non-zero forwarding address in the prefix advertisements in the base protocol. In such a scenario, the Flex-Algorithm specific reachability of the external prefix is determined by Flex-Algorithm specific reachability of the forwarding address.

In OSPF, the procedures for translation of NSSA external prefix advertisements into external prefix advertisements performed by an NSSA ABR [RFC3101] remain unchanged for Flex-Algorithm. An NSSA translator MUST include the OSPF FAPM Sub-TLVs for all Flex-Algorithms that are in the original NSSA external prefix advertisement from the NSSA ASBR in the translated external prefix advertisement generated by it regardless of its participation in those Flex-Algorithms or its having reachability to the NSSA ASBR in those Flex-Algorithms.

An area could become partitioned from the perspective of the Flex-Algorithm due to the constraints and/or metric being used for it, while maintaining the continuity in the algorithm 0. When that happens, some destinations inside that area could become unreachable in that Flex-Algorithm. These destinations will not be able to use an inter-area path. This is the consequence of the fact that the inter-area prefix reachability advertisement would not be available for these intra-area destinations within the area. It is RECOMMENDED to avoid such partitioning by providing enough redundancy inside the area for each Flex-Algorithm being used.

14. Flex-Algorithm and Forwarding Plane

This section describes how Flex-Algorithm paths are used in forwarding.

14.1. Segment Routing MPLS Forwarding for Flex-Algorithm

This section describes how Flex-Algorithm paths are used with SR MPLS forwarding.

Prefix SID advertisements include an SR-Algorithm value and, as such, are associated with the specified SR-Algorithm. Prefix-SIDs are also associated with a specific topology which is inherited from the associated prefix reachability advertisement. When the algorithm value advertised is a Flex-Algorithm value, the Prefix SID is associated with paths calculated using that Flex-Algorithm in the associated topology.

A Flex-Algorithm path MUST be installed in the MPLS forwarding plane using the MPLS label that corresponds to the Prefix-SID that was advertised for that Flex-algorithm. If the Prefix SID for a given Flex-algorithm is not known, the Flex-Algorithm specific path cannot be installed in the MPLS forwarding plane.

Traffic that is supposed to be routed via Flex-Algorithm specific paths, MUST be dropped when there are no such paths available.

Loop Free Alternate (LFA) paths for a given Flex-Algorithm MUST be computed using the same constraints as the calculation of the primary paths for that Flex-Algorithm. LFA paths MUST only use Prefix-SIDs advertised specifically for the given algorithm. LFA paths MUST NOT use an Adjacency-SID that belongs to a link that has been pruned from the Flex-Algorithm computation.

If LFA protection is being used to protect a given Flex-Algorithm paths, all routers in the area participating in the given Flex-Algorithm SHOULD advertise at least one Flex-Algorithm specific Node-SID. These Node-SIDs are used to steer traffic over the LFA computed backup path.

14.2. SRv6 Forwarding for Flex-Algorithm

This section describes how Flex-Algorithm paths are used with SRv6 forwarding.

In SRv6 a node is provisioned with topology/algorithm specific locators for each of the topology/algorithm pairs supported by that node. Each locator is an aggregate prefix for all SIDs provisioned on that node which have the matching topology/algorithm.

The SRv6 locator advertisement in IS-IS [I-D.ietf-lsr-isis-srv6-extensions] includes the MTID value that associates the locator with a specific topology. SRv6 locator

advertisements also includes an Algorithm value that explicitly associates the locator with a specific algorithm. When the algorithm value advertised with a locator represents a Flex-Algorithm, the paths to the locator prefix MUST be calculated using the specified Flex-Algorithm in the associated topology.

Forwarding entries for the locator prefixes advertised in IS-IS MUST be installed in the forwarding plane of the receiving SRv6 capable routers when the associated topology/algorithm is participating in them. Forwarding entries for locators associated with Flex-Algorithms in which the node is not participating MUST NOT be installed in the forwarding plane.

When the locator is associated with a Flex-Algorithm, LFA paths to the locator prefix MUST be calculated using such Flex-Algorithm in the associated topology, to guarantee that they follow the same constraints as the calculation of the primary paths. LFA paths MUST only use SRv6 SIDs advertised specifically for the given Flex-Algorithm.

If LFA protection is being used to protect locators associated with a given Flex-Algorithm, all routers in the area participating in the given Flex-Algorithm SHOULD advertise at least one Flex-Algorithm specific locator and END SID per node and one END.X SID for every link that has not been pruned from such Flex-Algorithm computation. These locators and SIDs are used to steer traffic over the LFA-computed backup path.

14.3. Other Applications' Forwarding for Flex-Algorithm

Any application that wants to use Flex-Algorithm specific forwarding needs to install some form of Flex-Algorithm specific forwarding entries.

Application-specific forwarding for Flex-Algorithm MUST be defined for each application and is outside of the scope of this document.

15. Operational Considerations

15.1. Inter-area Considerations

The scope of the Flex-Algorithm computation is an area, so is the scope of the FAD. In IS-IS, the Router Capability TLV in which the FAD Sub-TLV is advertised MUST have the S-bit clear, which prevents it to be flooded outside of the level in which it was originated. Even though in OSPF the FAD Sub-TLV can be flooded in an RI LSA that has AS flooding scope, the FAD selection is performed for each individual area in which it is being used.

There is no requirement for the FAD for a particular Flex-Algorithm to be identical in all areas in the network. For example, traffic for the same Flex-Algorithm may be optimized for minimal delay (e.g., using delay metric) in one area or level, while being optimized for available bandwidth (e.g., using IGP metric) in another area or level.

As described in Section 5.1, IS-IS allows the re-generation of the winning FAD from level 2, without any modification to it, into a level 1 area. This allows the operator to configure the FAD in one or multiple routers in the level 2, without the need to repeat the same task in each level 1 area, if the intent is to have the same FAD for the particular Flex-Algorithm across all levels. This can similarly be achieved in OSPF by using the AS flooding scope of the RI LSA in which the FAD Sub-TLV for the particular Flex-Algorithm is advertised.

Re-generation of FAD from a level 1 area to the level 2 area is not supported in IS-IS, so if the intent is to regenerate the FAD between IS-IS levels, the FAD MUST be defined on router(s) that are in level 2. In OSPF, the FAD definition can be done in any area and be propagated to all routers in the OSPF routing domain by using the AS flooding scope of the RI LSA.

15.2. Usage of SRLG Exclude Rule with Flex-Algorithm

There are two different ways in which SRLG information can be used with Flex-Algorithm:

- In a context of a single Flex-Algorithm, it can be used for computation of backup paths, as described in [I-D.ietf-rtgwg-segment-routing-ti-lfa]. This usage does not require association of any specific SRLG constraint with the given Flex-Algorithm definition.

- In the context of multiple Flex-Algorithms, it can be used for creating disjoint sets of paths by pruning the links belonging to a specific SRLG from the topology on which a specific Flex-Algorithm computes its paths. This usage:

 - Facilitates the usage of already deployed SRLG configurations for setup of disjoint paths between two or more Flex-Algorithms.

 - Requires explicit association of a given Flex-Algorithm with a specific set of SRLG constraints as defined in Section 6.5 and Section 7.5.

The two usages mentioned above are orthogonal.

15.3. Max-metric consideration

Both IS-IS and OSPF have a mechanism to set the IGP metric on a link to a value that would make the link either non-reachable or to serve as the link of last resort. Similar functionality would be needed for the Min Unidirectional Link Delay and TE metric, as these can be used to compute Flex-Algorithm paths.

The link can be made un-reachable for all Flex-Algorithms that use Min Unidirectional Link Delay as metric, as described in Section 5.1, by removing the Flex-Algorithm ASLA Min Unidirectional Link Delay advertisement for the link. The link can be made the link of last resort by setting the delay value in the Flex-Algorithm ASLA delay advertisement for the link to the value of 16,777,215 ($2^{24} - 1$).

The link can be made un-reachable for all Flex-Algorithms that use TE metric, as described in Section 5.1, by removing the Flex-Algorithm ASLA TE metric advertisement for the link. The link can be made the link of last resort by setting the TE metric value in the Flex-Algorithm ASLA delay advertisement for the link to the value of ($2^{24} - 1$) in IS-IS and ($2^{32} - 1$) in OSPF.

16. Backward Compatibility

This extension brings no new backward compatibility issues. IS-IS, OSPFv2 and OSPFv3 all have well defined handling of unrecognized TLVs and sub-TLVs that allows the introduction of the new extensions, similar to those defined here, without introducing any interoperability issues.

17. Security Considerations

This draft adds two new ways to disrupt IGP networks:

An attacker can hijack a particular Flex-Algorithm by advertising a FAD with a priority of 255 (or any priority higher than that of the legitimate nodes).

An attacker could make it look like a router supports a particular Flex-Algorithm when it actually doesn't, or vice versa.

Both of these attacks can be addressed by the existing security extensions as described in [RFC5304] and [RFC5310] for IS-IS, in [RFC2328] and [RFC7474] for OSPFv2, and in [RFC5340] and [RFC4552] for OSPFv3.

18. IANA Considerations

18.1. IGP IANA Considerations

18.1.1. IGP Algorithm Types Registry

This document makes the following registrations in the "IGP Algorithm Types" registry:

Type: 128-255.

Description: Flexible Algorithms.

Reference: This document (Section 4).

18.1.2. IGP Metric-Type Registry

IANA is requested to set up a registry called "IGP Metric-Type Registry" under an "Interior Gateway Protocol (IGP) Parameters" IANA registries. The registration policy for this registry is "Standards Action" ([RFC8126] and [RFC7120]).

Values in this registry come from the range 0-255.

This document registers following values in the "IGP Metric-Type Registry":

Type: 0

Description: IGP metric

Reference: This document (Section 5.1)

Type: 1

Description: Min Unidirectional Link Delay as defined in [RFC8570], section 4.2, and [RFC7471], section 4.2.

Reference: This document (Section 5.1)

Type: 2

Description: Traffic Engineering Default Metric as defined in [RFC5305], section 3.7, and Traffic engineering metric as defined in [RFC3630], section 2.5.5

Reference: This document (Section 5.1)

18.2. Flexible Algorithm Definition Flags Registry

IANA is requested to set up a registry called "IS-IS Flexible Algorithm Definition Flags Registry" under an "Interior Gateway Protocol (IGP) Parameters" IANA registries. The registration policy for this registry is "Standards Action" ([RFC8126] and [RFC7120]).

This document defines the following single bit in Flexible Algorithm Definition Flags registry:

Bit #	Name
-----	-----
0	Prefix Metric Flag (M-flag)

Reference: This document (Section 6.4, Section 7.4).

18.3. IS-IS IANA Considerations

18.3.1. Sub TLVs for Type 242

This document makes the following registrations in the "sub-TLVs for TLV 242" registry.

Type: 26.

Description: Flexible Algorithm Definition.

Reference: This document (Section 5.1).

18.3.2. Sub TLVs for for TLVs 135, 235, 236, and 237

This document makes the following registrations in the "Sub-TLVs for for TLVs 135, 235, 236, and 237" registry.

Type: 6

Description: Flexible Algorithm Prefix Metric.

Reference: This document (Section 8).

18.3.3. Sub-Sub-TLVs for Flexible Algorithm Definition Sub-TLV

This document creates the following Sub-Sub-TLV Registry:

Registry: Sub-Sub-TLVs for Flexible Algorithm Definition Sub-TLV

Registration Procedure: Expert review

Reference: This document (Section 5.1)

This document defines the following Sub-Sub-TLVs in the "Sub-Sub-TLVs for Flexible Algorithm Definition Sub-TLV" registry:

Type: 1

Description: Flexible Algorithm Exclude Admin Group

Reference: This document (Section 6.1).

Type: 2

Description: Flexible Algorithm Include-Any Admin Group

Reference: This document (Section 6.2).

Type: 3

Description: Flexible Algorithm Include-All Admin Group

Reference: This document (Section 6.3).

Type: 4

Description: Flexible Algorithm Definition Flags

Reference: This document (Section 6.4).

Type: 5

Description: Flexible Algorithm Exclude SRLG

Reference: This document (Section 6.5).

18.4. OSPF IANA Considerations

18.4.1. OSPF Router Information (RI) TLVs Registry

This specification updates the OSPF Router Information (RI) TLVs Registry.

Type: 16

Description: Flexible Algorithm Definition TLV.

Reference: This document (Section 5.2).

18.4.2. OSPFv2 Extended Prefix TLV Sub-TLVs

This document makes the following registrations in the "OSPFv2 Extended Prefix TLV Sub-TLVs" registry.

Type: 3

Description: Flexible Algorithm Prefix Metric.

Reference: This document (Section 9).

18.4.3. OSPFv3 Extended-LSA Sub-TLVs

This document makes the following registrations in the "OSPFv3 Extended-LSA Sub-TLVs" registry.

Type: 26

Description: Flexible Algorithm Prefix Metric.

Reference: This document (Section 9).

Type: TBD (suggested value 30)

Description: OSPF Flexible Algorithm ASBR Metric Sub-TLV

Reference: This document (Section 10.2).

18.4.4. OSPF Flex-Algorithm Prefix Metric Bits

This specification requests creation of "OSPF Flex-Algorithm Prefix Metric Bits" registry under the OSPF Parameters Registry with the following initial values.

Bit Number: 0

Description: E bit - External Type

Reference: this document.

The bits 1-7 are unassigned and the registration procedure to be followed for this registry is IETF Review.

18.4.5. OSPF Opaque LSA Option Types

This document makes the following registrations in the "OSPF Opaque LSA Option Types" registry.

Value: TBD (suggested value 11)

Description: OSPFv2 Extended Inter-Area ASBR LSA

Reference: This document (Section 10.1).

18.4.6. OSPFv2 Extended Inter-Area ASBR TLVs

This specification requests creation of "OSPFv2 Extended Inter-Area ASBR TLVs" registry under the OSPFv2 Parameters Registry with the following initial values.

Value: 1

Description : Extended Inter-Area ASBR TLV

Reference: this document

The values 2 to 32767 are unassigned, values 32768 to 33023 are reserved for experimental use while the values 0 and 33024 to 65535 are reserved. The registration procedure to be followed for this registry is IETF Review or IESG Approval.

18.4.7. OSPFv2 Inter-Area ASBR Sub-TLVs

This specification requests creation of "OSPFv2 Extended Inter-Area ASBR Sub-TLVs" registry under the OSPFv2 Parameters Registry with the following initial values.

Value: 1

Description : OSPF Flexible Algorithm ASBR Metric Sub-TLV

Reference: this document

The values 2 to 32767 are unassigned, values 32768 to 33023 are reserved for experimental use while the values 0 and 33024 to 65535 are reserved. The registration procedure to be followed for this registry is IETF Review or IESG Approval.

18.4.8. OSPF Flexible Algorithm Definition TLV Sub-TLV Registry

This document creates the following registry:

Registry: OSPF Flexible Algorithm Definition TLV sub-TLV

Registration Procedure: Expert review

Reference: This document (Section 5.2)

The "OSPF Flexible Algorithm Definition TLV sub-TLV" registry will define sub-TLVs at any level of nesting for the Flexible Algorithm TLV and should be added to the "Open Shortest Path First (OSPF) Parameters" registries group. New values can be allocated via IETF Review or IESG Approval.

This document registers following Sub-TLVs in the "TLVs for Flexible Algorithm Definition TLV" registry:

Type: 1

Description: Flexible Algorithm Exclude Admin Group

Reference: This document (Section 7.1).

Type: 2

Description: Flexible Algorithm Include-Any Admin Group

Reference: This document (Section 7.2).

Type: 3

Description: Flexible Algorithm Include-All Admin Group

Reference: This document (Section 7.3).

Type: 4

Description: Flexible Algorithm Definition Flags

Reference: This document (Section 7.4).

Type: 5

Description: Flexible Algorithm Exclude SRLG

Reference: This document (Section 7.5).

Types in the range 32768-33023 are for experimental use; these will not be registered with IANA, and MUST NOT be mentioned by RFCs.

Types in the range 33024-65535 are not to be assigned at this time. Before any assignments can be made in the 33024-65535 range, there MUST be an IETF specification that specifies IANA Considerations that covers the range being assigned.

18.4.9. Link Attribute Applications Registry

This document registers following bit in the Link Attribute Applications Registry:

Bit-3

Description: Flexible Algorithm (X-bit)

Reference: This document (Section 12).

19. Acknowledgements

This draft, among other things, is also addressing the problem that the [I-D.gulkohegde-routing-planes-using-sr] was trying to solve. All authors of that draft agreed to join this draft.

Thanks to Eric Rosen, Tony Przygienda, William Britto A J, Gunter Van De Velde, Dirk Goethals, Manju Sivaji and, Baalajee S for their detailed review and excellent comments.

Thanks to Cengiz Halit for his review and feedback during initial phase of the solution definition.

Thanks to Kenji Kumaki for his comments.

Thanks to Acee Lindem for editorial comments.

20. References

20.1. Normative References

[I-D.ietf-lsr-isis-srv6-extensions]

Psenak, P., Filshie, C., Bashandy, A., Decraene, B., and Z. Hu, "IS-IS Extensions to Support Segment Routing over IPv6 Dataplane", draft-ietf-lsr-isis-srv6-extensions-18 (work in progress), October 2021.

[ISO10589]

International Organization for Standardization, "Intermediate system to Intermediate system intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)", ISO/IEC 10589:2002, Second Edition, Nov 2002.

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC4203] Kompella, K., Ed. and Y. Rekhter, Ed., "OSPF Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 4203, DOI 10.17487/RFC4203, October 2005, <<https://www.rfc-editor.org/info/rfc4203>>.
- [RFC5250] Berger, L., Bryskin, I., Zinin, A., and R. Coltun, "The OSPF Opaque LSA Option", RFC 5250, DOI 10.17487/RFC5250, July 2008, <<https://www.rfc-editor.org/info/rfc5250>>.
- [RFC5307] Kompella, K., Ed. and Y. Rekhter, Ed., "IS-IS Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 5307, DOI 10.17487/RFC5307, October 2008, <<https://www.rfc-editor.org/info/rfc5307>>.
- [RFC7308] Osborne, E., "Extended Administrative Groups in MPLS Traffic Engineering (MPLS-TE)", RFC 7308, DOI 10.17487/RFC7308, July 2014, <<https://www.rfc-editor.org/info/rfc7308>>.
- [RFC7684] Psenak, P., Gredler, H., Shakir, R., Henderickx, W., Tantsura, J., and A. Lindem, "OSPFv2 Prefix/Link Attribute Advertisement", RFC 7684, DOI 10.17487/RFC7684, November 2015, <<https://www.rfc-editor.org/info/rfc7684>>.
- [RFC7770] Lindem, A., Ed., Shen, N., Vasseur, JP., Aggarwal, R., and S. Shaffer, "Extensions to OSPF for Advertising Optional Router Capabilities", RFC 7770, DOI 10.17487/RFC7770, February 2016, <<https://www.rfc-editor.org/info/rfc7770>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8362] Lindem, A., Roy, A., Goethals, D., Reddy Vallem, V., and F. Baker, "OSPFv3 Link State Advertisement (LSA) Extensibility", RFC 8362, DOI 10.17487/RFC8362, April 2018, <<https://www.rfc-editor.org/info/rfc8362>>.

- [RFC8665] Psenak, P., Ed., Previdi, S., Ed., Filsfils, C., Gredler, H., Shakir, R., Henderickx, W., and J. Tantsura, "OSPF Extensions for Segment Routing", RFC 8665, DOI 10.17487/RFC8665, December 2019, <<https://www.rfc-editor.org/info/rfc8665>>.
- [RFC8666] Psenak, P., Ed. and S. Previdi, Ed., "OSPFv3 Extensions for Segment Routing", RFC 8666, DOI 10.17487/RFC8666, December 2019, <<https://www.rfc-editor.org/info/rfc8666>>.
- [RFC8667] Previdi, S., Ed., Ginsberg, L., Ed., Filsfils, C., Bashandy, A., Gredler, H., and B. Decraene, "IS-IS Extensions for Segment Routing", RFC 8667, DOI 10.17487/RFC8667, December 2019, <<https://www.rfc-editor.org/info/rfc8667>>.
- [RFC8919] Ginsberg, L., Psenak, P., Previdi, S., Henderickx, W., and J. Drake, "IS-IS Application-Specific Link Attributes", RFC 8919, DOI 10.17487/RFC8919, October 2020, <<https://www.rfc-editor.org/info/rfc8919>>.
- [RFC8920] Psenak, P., Ed., Ginsberg, L., Henderickx, W., Tantsura, J., and J. Drake, "OSPF Application-Specific Link Attributes", RFC 8920, DOI 10.17487/RFC8920, October 2020, <<https://www.rfc-editor.org/info/rfc8920>>.

20.2. Informative References

- [I-D.gulkohegde-routing-planes-using-sr] Hegde, S. and A. Gulko, "Separating Routing Planes using Segment Routing", draft-gulkohegde-routing-planes-using-sr-00 (work in progress), March 2017.
- [I-D.ietf-rtgwg-segment-routing-ti-lfa] Litkowski, S., Bashandy, A., Filsfils, C., Francois, P., Decraene, B., and D. Voyer, "Topology Independent Fast Reroute using Segment Routing", draft-ietf-rtgwg-segment-routing-ti-lfa-08 (work in progress), January 2022.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC3101] Murphy, P., "The OSPF Not-So-Stubby Area (NSSA) Option", RFC 3101, DOI 10.17487/RFC3101, January 2003, <<https://www.rfc-editor.org/info/rfc3101>>.

- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", RFC 3630, DOI 10.17487/RFC3630, September 2003, <<https://www.rfc-editor.org/info/rfc3630>>.
- [RFC3906] Shen, N. and H. Smit, "Calculating Interior Gateway Protocol (IGP) Routes Over Traffic Engineering Tunnels", RFC 3906, DOI 10.17487/RFC3906, October 2004, <<https://www.rfc-editor.org/info/rfc3906>>.
- [RFC4552] Gupta, M. and N. Melam, "Authentication/Confidentiality for OSPFv3", RFC 4552, DOI 10.17487/RFC4552, June 2006, <<https://www.rfc-editor.org/info/rfc4552>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC7120] Cotton, M., "Early IANA Allocation of Standards Track Code Points", BCP 100, RFC 7120, DOI 10.17487/RFC7120, January 2014, <<https://www.rfc-editor.org/info/rfc7120>>.
- [RFC7471] Giacalone, S., Ward, D., Drake, J., Atlas, A., and S. Previdi, "OSPF Traffic Engineering (TE) Metric Extensions", RFC 7471, DOI 10.17487/RFC7471, March 2015, <<https://www.rfc-editor.org/info/rfc7471>>.
- [RFC7474] Bhatia, M., Hartman, S., Zhang, D., and A. Lindem, Ed., "Security Extension for OSPFv2 When Using Manual Key Management", RFC 7474, DOI 10.17487/RFC7474, April 2015, <<https://www.rfc-editor.org/info/rfc7474>>.

- [RFC8126] Cotton, M., Leiba, B., and T. Narten, "Guidelines for Writing an IANA Considerations Section in RFCs", BCP 26, RFC 8126, DOI 10.17487/RFC8126, June 2017, <<https://www.rfc-editor.org/info/rfc8126>>.
- [RFC8570] Ginsberg, L., Ed., Previdi, S., Ed., Giacalone, S., Ward, D., Drake, J., and Q. Wu, "IS-IS Traffic Engineering (TE) Metric Extensions", RFC 8570, DOI 10.17487/RFC8570, March 2019, <<https://www.rfc-editor.org/info/rfc8570>>.

Authors' Addresses

Peter Psenak (editor)
Cisco Systems
Apollo Business Center
Mlynske nivy 43
Bratislava, 82109
Slovakia

Email: ppsenak@cisco.com

Shraddha Hegde
Juniper Networks, Inc.
Embassy Business Park
Bangalore, KA, 560093
India

Email: shraddha@juniper.net

Clarence Filsfils
Cisco Systems, Inc.
Brussels
Belgium

Email: cfilsfil@cisco.com

Ketan Talaulikar
Arrcus, Inc
India

Email: ketant.ietf@gmail.com

Arkadiy Gulko
Edward Jones

Email: arkadiy.gulko@edwardjones.com

Internet
Internet-Draft
Intended status: Standards Track
Expires: 6 July 2022

D. Yeung
Arrcus
Y. Qu
Futurewei
J. Zhang
Juniper Networks
I. Chen
The MITRE Corporation
A. Lindem
Cisco Systems
2 January 2022

YANG Data Model for OSPF Segment Routing
draft-ietf-ospf-sr-yang-17

Abstract

This document defines a YANG data module that can be used to configure and manage OSPF Extensions for Segment Routing. It also defines a module for management of Signaling Maximum SID Depth (MSD) Using OSPF.

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 <https://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 6 July 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Overview	2
1.1. Requirements Language	3
1.2. Tree Diagrams	3
2. OSPF MSD	3
2.1. OSPF MSD YANG Module	4
3. OSPF Segment Routing	11
3.1. OSPF Segment Routing YANG Module	16
4. Security Considerations	30
5. Acknowledgements	31
6. IANA Considerations	31
7. References	32
7.1. Normative References	32
7.2. Informative References	34
Appendix A. Contributors' Addresses	34
Authors' Addresses	34

1. Overview

YANG [RFC7950] is a data definition language used to define the contents of a conceptual data store that allows networked devices to be managed using NETCONF [RFC6241]. YANG is proving relevant beyond its initial confines, as bindings to other interfaces (e.g., ReST) and encodings other than XML (e.g., JSON) are being defined. Furthermore, YANG data models can be used as the basis for implementation of other interfaces, such as CLI and programmatic APIs.

This document defines a YANG data model that can be used to configure and manage OSPFv2 extensions for Segment Routing [RFC8665] and it is an augmentation to the OSPF YANG data model.

This document also defines a YANG data model for Signaling Maximum SID Depth (MSD) Using OSPF [RFC8476], which augments the base OSPF YANG data model.

The YANG module in this document conforms to the Network Management Datastore Architecture (NMDA) [RFC8342].

1.1. Requirements Language

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 [RFC2119].

1.2. Tree Diagrams

This document uses the graphical representation of data models defined in [RFC8340].

2. OSPF MSD

This document defines a model for Signaling Maximum SID Depth (MSD) Using OSPF [RFC8476]. It is an augmentation of the OSPF base model.

```

module: ietf-ospf-msd
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
    /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
    /ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
    /ospf:body/ospf:opaque/ospf:ri-opaque:
  +--ro node-msd-tlv
    +--ro node-msds* [msd-type]
      +--ro msd-type      identityref
      +--ro msd-value?    uint8
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf/ospf:database
    /ospf:as-scope-lsa-type/ospf:as-scope-lsas/ospf:as-scope-lsa
    /ospf:version/ospf:ospfv2/ospf:ospfv2/ospf:body/ospf:opaque
    /ospf:ri-opaque:
  +--ro node-msd-tlv
    +--ro node-msds* [msd-type]
      +--ro msd-type      identityref
      +--ro msd-value?    uint8
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
    /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
    /ospf:area-scope-lsa/ospf:version/ospf:ospfv3/ospf:ospfv3
    /ospf:body/ospf:router-information:
  +--ro node-msd-tlv
    +--ro node-msds* [msd-type]
      +--ro msd-type      identityref
      +--ro msd-value?    uint8
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf/ospf:database
    /ospf:as-scope-lsa-type/ospf:as-scope-lsas/ospf:as-scope-lsa
    /ospf:version/ospf:ospfv3/ospf:ospfv3/ospf:body

```

```

        /ospf:router-information:
+--ro node-msd-tlv
  +--ro node-msds* [msd-type]
    +--ro msd-type      identityref
    +--ro msd-value?    uint8
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
  /ospf:interfaces/ospf:interface/ospf:database
  /ospf:link-scope-lsa-type/ospf:link-scope-lsas
  /ospf:link-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
  /ospf:body/ospf:opaque/ospf:extended-link-opaque
  /ospf:extended-link-tlv:
+--ro link-msd-sub-tlv
  +--ro link-msds* [msd-type]
    +--ro msd-type      identityref
    +--ro msd-value?    uint8
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
  /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
  /ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
  /ospf:body/ospf:opaque/ospf:extended-link-opaque
  /ospf:extended-link-tlv:
+--ro link-msd-sub-tlv
  +--ro link-msds* [msd-type]
    +--ro msd-type      identityref
    +--ro msd-value?    uint8
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:database
  /ospf:as-scope-lsa-type/ospf:as-scope-lsas/ospf:as-scope-lsa
  /ospf:version/ospf:ospfv2/ospf:ospfv2/ospf:body/ospf:opaque
  /ospf:extended-link-opaque/ospf:extended-link-tlv:
+--ro link-msd-sub-tlv
  +--ro link-msds* [msd-type]
    +--ro msd-type      identityref
    +--ro msd-value?    uint8
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
  /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
  /ospf:area-scope-lsa/ospf:version/ospf:ospfv3/ospf:ospfv3
  /ospf:body/ospfv3-e-lsa:e-router/ospfv3-e-lsa:e-router-tlvs:
+--ro link-msd-sub-tlv
  +--ro link-msds* [msd-type]
    +--ro msd-type      identityref
    +--ro msd-value?    uint8

```

2.1. OSPF MSD YANG Module

```
<CODE BEGINS> file "ietf-ospf-msd@2022-01-02.yang"
module ietf-ospf-msd {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:yang:ietf-ospf-msd";
  prefix ospf-msd;

  import ietf-routing {
    prefix rt;
    reference "RFC 8349: A YANG Data Model for Routing
              Management (NMDA Version)";
  }

  import ietf-ospf {
    prefix ospf;
  }

  import ietf-ospfv3-extended-lsa {
    prefix ospfv3-e-lsa;
  }

  organization
    "IETF LSR - LSR Working Group";
  contact
    "WG Web:  <https://tools.ietf.org/wg/mppls/>
    WG List:  <mailto:mppls@ietf.org>

    Author:   Yingzhen Qu
              <mailto:yingzhen.qu@futurewei.com>
    Author:   Acee Lindem
              <mailto:acee@cisco.com>
    Author:   Stephane Litkowski
              <mailto:slitkows.ietf@gmail.com>
    Author:   Jeff Tantsura
              <jefftant.ietf@gmail.com>

";
  description
    "The YANG module augments the base OSPF model to
    manage different types of MSDs.

    This YANG model conforms to the Network Management
    Datastore Architecture (NMDA) as described in RFC 8342.

    Copyright (c) 2022 IETF Trust and the persons identified as
    authors of the code. All rights reserved.

    Redistribution and use in source and binary forms, with or
    without modification, is permitted pursuant to, and subject to
```

the license terms contained in, the Revised BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.

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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.";

```
reference "RFC XXXX: YANG Data Model for OSPF MSD.";

revision 2022-01-02 {
  description
    "Initial Version";
  reference "RFC XXXX: YANG Data Model for OSPF MSD.";
}

identity msd-base-type {
  description
    "Base identity for MSD Type";
}

identity base-mpls-msd {
  base msd-base-type;
  description
    "Base MPLS Imposition MSD.";
  reference
    "RFC 8491: Singling MSD using IS-IS.";
}

identity erld-msd {
  base msd-base-type;
  description
    "ERLD-MSD is defined to advertise the ERLD.";
  reference
    "RFC 8662: Entropy Label for Source Packet Routing in
      Networking (SPRING) Tunnels";
}

grouping node-msd-tlv {
  description
    "Grouping for node MSD.";
```

```
    container node-msd-tlv {
      list node-msds {
        key "msd-type";
        leaf msd-type {
          type identityref {
            base msd-base-type;
          }
          description
            "MSD-Types";
        }
        leaf msd-value {
          type uint8;
          description
            "MSD value, in the range of 0-255.";
        }
        description
          "Node MSD is the smallest link MSD supported by
           the node.";
      }
      description
        "Node MSD is the number of SIDs supported by a node.";
      reference
        "RFC 8476: Signaling Maximum SID Depth (MSD) Using OSPF";
    }
  }

  grouping link-msd-sub-tlv {
    description
      "Link Maximum SID Depth (MSD) grouping for an interface.";
    container link-msd-sub-tlv {
      list link-msds {
        key "msd-type";
        leaf msd-type {
          type identityref {
            base msd-base-type;
          }
          description
            "MSD-Types";
        }
        leaf msd-value {
          type uint8;
          description
            "MSD value, in the range of 0-255.";
        }
        description
          "List of link MSDs";
      }
    }
    description
```

```

        "Link MSD sub-tlvs.";
    }
}

/* Node MSD TLV */
augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/"
+ "ospf:area/ospf:database/"
+ "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
+ "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:ri-opaque" {
when "../../../../../../../../../../../../../../../"
+ "rt:type = 'ospf:ospfv2'" {
description
    "This augmentation is only valid for OSPFv2.";
}
description
    "Node MSD TLV is an optional TLV of OSPFv2 RI Opaque
    LSA (RFC7770) and has a type of 12.";

uses node-msd-tlv;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:database/"
+ "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
+ "ospf:as-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:ri-opaque" {
when "../../../../../../../../../../../../../../../"
+ "rt:type = 'ospf:ospfv2'" {
description
    "This augmentation is only valid for OSPFv2.";
}
description
    "Node MSD TLV is an optional TLV of OSPFv2 RI Opaque
    LSA (RFC7770) and has a type of 12.";

uses node-msd-tlv;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/"
+ "ospf:area/ospf:database/"

```

```

        + "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
        + "ospf:area-scope-lsa/ospf:version/ospf:ospfv3/"
        + "ospf:ospfv3/ospf:body/ospf:router-information" {
when "../.../.../.../.../.../.../.../.../..."
        + "rt:type = 'ospf:ospfv3'" {
        description
            "This augmentation is only valid for OSPFv3.";
        }
        description
            "Node MSD TLV is an optional TLV of OSPFv3 RI Opaque
            LSA (RFC7770) and has a type of 12.";

        uses node-msd-tlv;
    }

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:database/"
    + "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
    + "ospf:as-scope-lsa/ospf:version/ospf:ospfv3/"
    + "ospf:ospfv3/ospf:body/ospf:router-information" {
when "../.../.../.../.../.../.../.../..."
        + "rt:type = 'ospf:ospfv3'" {
        description
            "This augmentation is only valid for OSPFv3.";
        }
        description
            "Node MSD TLV is an optional TLV of OSPFv3 RI Opaque
            LSA (RFC7770) and has a type of 12.";

        uses node-msd-tlv;
    }

/* link MSD sub-tlv */
augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:areas/ospf:area/"
    + "ospf:interfaces/ospf:interface/ospf:database/"
    + "ospf:link-scope-lsa-type/ospf:link-scope-lsas/"
    + "ospf:link-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/"
    + "ospf:extended-link-opaque/ospf:extended-link-tlv" {
when "../.../.../.../.../.../.../.../.../.../.../..."
        + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
        }
        description

```

```

    "Link MSD sub-TLV is an optional sub-TLV of OSPFv2 extended
    link TLV as defined in RFC 7684 and has a type of 6.";

    uses link-msd-sub-tlv;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/"
+ "ospf:area/ospf:database/"
+ "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
+ "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:extended-link-opaque/ospf:extended-link-tlv" {
when "../.../.../.../.../.../.../.../.../.../.../..."
+ "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
}
description
    "Link MSD sub-TLV is an optional sub-TLV of OSPFv2 extended
    link TLV as defined in RFC 7684 and has a type of 6.";

    uses link-msd-sub-tlv;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:database/"
+ "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
+ "ospf:as-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:extended-link-opaque/ospf:extended-link-tlv" {
when "../.../.../.../.../.../.../.../.../.../.../..."
+ "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
}
description
    "Link MSD sub-TLV is an optional sub-TLV of OSPFv2 extended
    link TLV as defined in RFC 7684 and has a type of 6.";

    uses link-msd-sub-tlv;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/ospf:area/ospf:database/"

```



```

+ "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
+ "ospf:area-scope-lsa/ospf:version/ospf:ospfv3/"
+ "ospf:ospfv3/ospf:body/ospfv3-e-lsa:e-router"
+ "/ospfv3-e-lsa:e-router-tlvs" {
when "ospf:.../.../.../.../.../.../.../..."
+ "rt:type" = 'ospf:ospfv3' {
description
    "This augmentation is only valid for OSPFv3
    E-Router LSAs";
}
description
    "Augment OSPFv3 Area scope router-link TLV.";

uses link-msd-sub-tlv;
}
}
<CODE ENDS>

```

3. OSPF Segment Routing

This document defines a model for OSPF Segment Routing feature [RFC8665]. It is an augmentation of the OSPF base model.

The OSPF SR YANG module requires support for the base segment routing module [RFC9020], which defines the global segment routing configuration independent of any specific routing protocol configuration, and support of OSPF base model[I-D.ietf-ospf-yang] which defines basic OSPF configuration and state.

```

module: ietf-ospf-sr
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf:
      +--rw segment-routing
      |   +--rw enabled?      boolean
      |   +--rw bindings {mapping-server}?
      |   |   +--rw advertise
      |   |   |   +--rw policies* -> /rt:routing/sr:segment-routing
      |   |   |   |   /sr-mpls:sr-mpls/bindings
      |   |   |   |   /mapping-server/policy/name
      |   |   +--rw receive?    boolean
      +--rw protocol-srgb {sr-mpls:protocol-srgb}?
      |   +--rw srgb* [lower-bound upper-bound]
      |   +--rw lower-bound    uint32
      |   +--rw upper-bound    uint32
  augment /rt:routing/rt:control-plane-protocols
    /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
      /ospf:interfaces/ospf:interface:

```

```

+--rw segment-routing
  +--rw adjacency-sid
    +--rw adj-sids* [value]
      |   +--rw value-type?  enumeration
      |   +--rw value        uint32
      |   +--rw protected?   boolean
      |   +--rw weight?      uint8
      +--rw advertise-adj-group-sid* [group-id]
        |   +--rw group-id    uint32
        +--rw advertise-protection? enumeration
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:interfaces/ospf:interface/ospf:fast-reroute:
+--rw ti-lfa {ti-lfa}?
  +--rw enable?    boolean
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:interfaces/ospf:interface/ospf:database
/ospf:link-scope-lsa-type/ospf:link-scope-lsas
/ospf:link-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
/ospf:body/ospf:opaque/ospf:extended-prefix-opaque:
+--ro extended-prefix-range-tlvs
  +--ro extended-prefix-range-tlv* []
    +--ro prefix-length?            uint8
    +--ro af?                      uint8
    +--ro range-size?               uint16
    +--ro extended-prefix-range-flags
      |   +--ro bits*    identityref
    +--ro prefix?                  inet:ip-prefix
    +--ro prefix-sid-sub-tlvs
      |   +--ro prefix-sid-sub-tlv* []
      |   |   +--ro prefix-sid-flags
      |   |   |   +--ro bits*    identityref
      |   |   +--ro mt-id?        uint8
      |   |   +--ro algorithm?    uint8
      |   |   +--ro sid?          uint32
    +--ro unknown-tlvs
      +--ro unknown-tlv* []
        +--ro type?    uint16
        +--ro length?  uint16
        +--ro value?   yang:hex-string
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
/ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
/ospf:body/ospf:opaque/ospf:extended-prefix-opaque:
+--ro extended-prefix-range-tlvs
  +--ro extended-prefix-range-tlv* []

```

```

+--ro prefix-length?                uint8
+--ro af?                           uint8
+--ro range-size?                   uint16
+--ro extended-prefix-range-flags
|   +--ro bits*    identityref
+--ro prefix?                       inet:ip-prefix
+--ro prefix-sid-sub-tlvs
|   +--ro prefix-sid-sub-tlv* []
|   |   +--ro prefix-sid-flags
|   |   |   +--ro bits*    identityref
|   |   +--ro mt-id?        uint8
|   |   +--ro algorithm?    uint8
|   |   +--ro sid?          uint32
+--ro unknown-tlvs
|   +--ro unknown-tlv* []
|   |   +--ro type?        uint16
|   |   +--ro length?      uint16
|   |   +--ro value?       yang:hex-string
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:database
/ospf:as-scope-lsa-type/ospf:as-scope-lsas
/ospf:as-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
/ospf:body/ospf:opaque/ospf:extended-prefix-opaque:
+--ro extended-prefix-range-tlvs
|   +--ro extended-prefix-range-tlv* []
|   |   +--ro prefix-length?                uint8
|   |   +--ro af?                           uint8
|   |   +--ro range-size?                   uint16
|   |   +--ro extended-prefix-range-flags
|   |   |   +--ro bits*    identityref
|   |   +--ro prefix?                       inet:ip-prefix
|   |   +--ro prefix-sid-sub-tlvs
|   |   |   +--ro prefix-sid-sub-tlv* []
|   |   |   |   +--ro prefix-sid-flags
|   |   |   |   |   +--ro bits*    identityref
|   |   |   |   +--ro mt-id?        uint8
|   |   |   |   +--ro algorithm?    uint8
|   |   |   |   +--ro sid?          uint32
|   |   +--ro unknown-tlvs
|   |   |   +--ro unknown-tlv* []
|   |   |   |   +--ro type?        uint16
|   |   |   |   +--ro length?      uint16
|   |   |   |   +--ro value?       yang:hex-string
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:interfaces/ospf:interface/ospf:database
/ospf:link-scope-lsa-type/ospf:link-scope-lsas
/ospf:link-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2

```

```

        /ospf:body/ospf:opaque/ospf:extended-prefix-opaque
        /ospf:extended-prefix-tlv:
+---ro prefix-sid-sub-tlvs
  +---ro prefix-sid-sub-tlv* []
    +---ro prefix-sid-flags
      | +---ro bits* identityref
    +---ro mt-id? uint8
    +---ro algorithm? uint8
    +---ro sid? uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
  /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
  /ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
  /ospf:body/ospf:opaque/ospf:extended-prefix-opaque
  /ospf:extended-prefix-tlv:
+---ro prefix-sid-sub-tlvs
  +---ro prefix-sid-sub-tlv* []
    +---ro prefix-sid-flags
      | +---ro bits* identityref
    +---ro mt-id? uint8
    +---ro algorithm? uint8
    +---ro sid? uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:database
  /ospf:as-scope-lsa-type/ospf:as-scope-lsas/ospf:as-scope-lsa
  /ospf:version/ospf:ospfv2/ospf:ospfv2/ospf:body/ospf:opaque
  /ospf:extended-prefix-opaque/ospf:extended-prefix-tlv:
+---ro prefix-sid-sub-tlvs
  +---ro prefix-sid-sub-tlv* []
    +---ro prefix-sid-flags
      | +---ro bits* identityref
    +---ro mt-id? uint8
    +---ro algorithm? uint8
    +---ro sid? uint32
augment /rt:routing/rt:control-plane-protocols
  /rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
  /ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
  /ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
  /ospf:body/ospf:opaque/ospf:extended-link-opaque
  /ospf:extended-link-tlv:
+---ro adj-sid-sub-tlvs
  | +---ro adj-sid-sub-tlv* []
  | +---ro adj-sid-flags
  | | +---ro bits* identityref
  | +---ro mt-id? uint8
  | +---ro weight? uint8
  | +---ro sid? uint32
+---ro lan-adj-sid-sub-tlvs

```

```

    +---ro lan-adj-sid-sub-tlv* []
    +---ro lan-adj-sid-flags
    |   +---ro bits*      identityref
    +---ro mt-id?          uint8
    +---ro weight?         uint8
    +---ro neighbor-router-id?  yang:dotted-quad
    +---ro sid?            uint32
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:interfaces/ospf:interface/ospf:database
/ospf:link-scope-lsa-type/ospf:link-scope-lsas
/ospf:link-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
/ospf:body/ospf:opaque/ospf:ri-opaque:
+---ro sr-algorithm-tlv
|   +---ro sr-algorithm*   uint8
+---ro sid-range-tlvs
|   +---ro sid-range-tlv* []
|   +---ro range-size?    uint24
|   +---ro sid-sub-tlv
|   |   +---ro sid?      uint32
+---ro local-block-tlvs
|   +---ro local-block-tlv* []
|   +---ro range-size?    uint24
|   +---ro sid-sub-tlv
|   |   +---ro sid?      uint32
+---ro srms-preference-tlv
    +---ro preference?    uint8
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:areas/ospf:area
/ospf:database/ospf:area-scope-lsa-type/ospf:area-scope-lsas
/ospf:area-scope-lsa/ospf:version/ospf:ospfv2/ospf:ospfv2
/ospf:body/ospf:opaque/ospf:ri-opaque:
+---ro sr-algorithm-tlv
|   +---ro sr-algorithm*   uint8
+---ro sid-range-tlvs
|   +---ro sid-range-tlv* []
|   +---ro range-size?    uint24
|   +---ro sid-sub-tlv
|   |   +---ro sid?      uint32
+---ro local-block-tlvs
|   +---ro local-block-tlv* []
|   +---ro range-size?    uint24
|   +---ro sid-sub-tlv
|   |   +---ro sid?      uint32
+---ro srms-preference-tlv
    +---ro preference?    uint8
augment /rt:routing/rt:control-plane-protocols
/rt:control-plane-protocol/ospf:ospf/ospf:database

```

```

        /ospf:as-scope-lsa-type/ospf:as-scope-lsas/ospf:as-scope-lsa
        /ospf:version/ospf:ospfv2/ospf:ospfv2/ospf:body/ospf:opaque
        /ospf:ri-opaque:
+--ro sr-algorithm-tlv
|   +--ro sr-algorithm*      uint8
+--ro sid-range-tlvs
|   +--ro sid-range-tlv* []
|       +--ro range-size?    uint24
|       +--ro sid-sub-tlv
|           +--ro sid?      uint32
+--ro local-block-tlvs
|   +--ro local-block-tlv* []
|       +--ro range-size?    uint24
|       +--ro sid-sub-tlv
|           +--ro sid?      uint32
+--ro srms-preference-tlv
    +--ro preference?      uint8

```

3.1. OSPF Segment Routing YANG Module

```

<CODE BEGINS> file "ietf-ospf-sr@2022-01-02.yang"
module ietf-ospf-sr {
  yang-version 1.1;
  namespace "urn:ietf:params:xml:ns:yang:ietf-ospf-sr";

  prefix ospf-sr;

  import ietf-inet-types {
    prefix "inet";
    reference "RFC 6991 - Common YANG Data Types";
  }

  import ietf-yang-types {
    prefix "yang";
    reference "RFC 6991 - Common YANG Data Types";
  }

  import ietf-routing {
    prefix "rt";
    reference "RFC 8349 - A YANG Data Model for Routing
              Management (NMDA Version)";
  }

  import ietf-segment-routing-common {
    prefix "sr-cmn";
    reference "RFC 9020 - YANG Data Model for Segment
              Routing";
  }

```

```
}
import ietf-segment-routing-mpls {
  prefix "sr-mpls";
  reference "RFC 9020 - YANG Data Model for Segment
    Routing";
}
import ietf-ospf {
  prefix "ospf";
}

organization
  "IETF LSR - Link State Routing Working Group";

contact
  "WG Web:    <http://tools.ietf.org/wg/lsr/>
  WG List:    <mailto:lsr@ietf.org>

  Editor:     Derek Yeung
               <mailto:derek@arccus.com>
  Author:     Derek Yeung
               <mailto:derek@arccus.com>
  Author:     Yingzhen Qu
               <mailto:yingzhen.qu@futurewei.com>
  Author:     Acee Lindem
               <mailto:acee@cisco.com>
  Author:     Jeffrey Zhang
               <mailto:zzhang@juniper.net>
  Author:     Ing-Wher Chen
               <mailto:ingwherchen@mitre.org>
  Author:     Greg Hankins
               <mailto:greg.hankins@alcatel-lucent.com>";
```

description

"This YANG module defines the generic configuration and operational state for OSPF Segment Routing, which is common across all of the vendor implementations. It is intended that the module will be extended by vendors to define vendor-specific OSPF Segment Routing configuration and operational parameters and policies.

This YANG model conforms to the Network Management Datastore Architecture (NMDA) as described in RFC 8342.

Copyright (c) 2022 IETF Trust and the persons identified as authors of the code. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, is permitted pursuant to, and subject to

the license terms contained in, the Revised BSD License set forth in Section 4.c of the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>).

This version of this YANG module is part of RFC XXXX (<https://www.rfc-editor.org/info/rfcXXXX>); see the RFC itself for full legal notices.

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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.";

reference "RFC XXXX";

```
revision 2022-01-02 {  
  description  
    "Initial revision.";  
  reference  
    "RFC XXXX: A YANG Data Model for OSPF Segment Routing.";  
}
```

```
feature ti-lfa {  
  description  
    "Topology-Independent Loop-Free Alternate (TI-LFA)  
    computation using segment routing.";  
}
```

```
identity prefix-sid-bit {  
  description  
    "Base identity for prefix sid sub-tlv bits.";  
}
```

```
identity np-bit {  
  base prefix-sid-bit;  
  description  
    "No-PHP flag.";  
}
```

```
identity m-bit {  
  base prefix-sid-bit;  
  description
```



```
        "Mapping server flag.";
    }

    identity e-bit {
        base prefix-sid-bit;
        description
            "Explicit-NULL flag.";
    }

    identity v-bit {
        base prefix-sid-bit;
        description
            "Value/Index flag.";
    }

    identity l-bit {
        base prefix-sid-bit;
        description
            "Local flag.";
    }

    identity extended-prefix-range-bit {
        description
            "Base identity for extended prefix range TLV bits.";
    }

    identity ia-bit {
        base extended-prefix-range-bit;
        description
            "Inter-Area flag. If set, advertisement is of inter-area type.";
    }

    identity adj-sid-bit {
        description
            "Base identity for adj sid sub-tlv bits.";
    }

    identity b-bit {
        base adj-sid-bit;
        description
            "Backup flag.";
    }

    identity vi-bit {
        base adj-sid-bit;
        description
            "Value/Index flag.";
    }
}
```

```
identity lo-bit {
    base adj-sid-bit;
    description
        "Local/Global flag.";
}

identity g-bit {
    base adj-sid-bit;
    description
        "Group flag.";
}

identity p-bit {
    base adj-sid-bit;
    description
        "Persistent flag.";
}

typedef uint24 {
    type uint32 {
        range "0 .. 16777215";
    }
    description
        "24-bit unsigned integer.";
}

/* Groupings */
grouping sid-sub-tlv {
    description "SID/Label sub-TLV grouping.";
    container sid-sub-tlv {
        description
            "Used to advertise the SID/Label associated with a
            prefix or adjacency.";
        leaf sid {
            type uint32;
            description
                "Segment Identifier (SID) - A 20 bit label or
                32 bit SID.";
        }
    }
}

grouping prefix-sid-sub-tlvs {
    description "Prefix Segment ID (SID) sub-TLVs.";
    container prefix-sid-sub-tlvs {
        description "Prefix SID sub-TLV.";
        list prefix-sid-sub-tlv {
            description "Prefix SID sub-TLV.";
        }
    }
}
```

```
    container prefix-sid-flags {
      leaf-list bits {
        type identityref {
          base prefix-sid-bit;
        }
        description
          "Prefix SID Sub-TLV flag bits list.";
      }
      description "Segment Identifier (SID) Flags.";
    }
    leaf mt-id {
      type uint8;
      description "Multi-topology ID.";
    }
    leaf algorithm {
      type uint8;
      description
        "The algorithm associated with the prefix-SID.";
    }
    leaf sid {
      type uint32;
      description "An index or label.";
    }
  }
}

grouping extended-prefix-range-tlvs {
  description "Extended prefix range TLV grouping.";

  container extended-prefix-range-tlvs {
    description "The list of range of prefixes.";
    list extended-prefix-range-tlv {
      description "The range of prefixes.";
      leaf prefix-length {
        type uint8;
        description "Length of prefix in bits.";
      }
    }
    leaf af {
      type uint8;
      description "Address family for the prefix.";
    }
    leaf range-size {
      type uint16;
      description "The number of prefixes covered by the
        advertisement.";
    }
    container extended-prefix-range-flags {
```

```

        leaf-list bits {
            type identityref {
                base extended-prefix-range-bit;
            }
            description "Extended prefix range TLV flags list.";
        }
        description "Extended Prefix Range TLV flags.";
    }
    leaf prefix {
        type inet:ip-prefix;
        description "Address prefix.";
    }
    uses prefix-sid-sub-tlvs;
    uses ospf:unknown-tlvs;
}
}

grouping sr-algorithm-tlv {
    description "SR algorithm TLV grouping.";
    container sr-algorithm-tlv {
        description "All SR algorithm TLVs.";
        leaf-list sr-algorithm {
            type uint8;
            description
                "The Segment Routing (SR) algorithms that the router is
                currently using.";
        }
    }
}

grouping sid-range-tlvs {
    description "SID Range TLV grouping.";
    container sid-range-tlvs {
        description "List of SID range TLVs.";
        list sid-range-tlv {
            description "SID range TLV.";
            leaf range-size {
                type uint24;
                description "The SID range.";
            }
            uses sid-sub-tlv;
        }
    }
}

grouping local-block-tlvs {
    description "The SR local block TLV contains the

```

```

        range of labels reserved for local SIDs.";
    container local-block-tlvs {
        description "List of SRLB TLVs.";
        list local-block-tlv {
            description "SRLB TLV.";
            leaf range-size {
                type uint24;
                description "The SID range.";
            }
            uses sid-sub-tlv;
        }
    }
}

grouping srms-preference-tlv {
    description "The SRMS preference TLV is used to advertise
        a preference associated with the node that acts
        as an SR Mapping Server.";
    container srms-preference-tlv {
        description "SRMS Preference TLV.";
        leaf preference {
            type uint8 {
                range "0 .. 255";
            }
            description "SRMS preference TLV, value from 0 to 255.";
        }
    }
}

/* Configuration */
augment "/rt:routing/rt:control-plane-protocols"
    + "/rt:control-plane-protocol/ospf:ospf" {
    when "../rt:type = 'ospf:ospfv2' or "
    + "../rt:type = 'ospf:ospfv3'" {
        description
            "This augments the OSPF routing protocol when used.";
    }
    description
        "This augments the OSPF protocol configuration
        with segment routing.";
    uses sr-mpls:sr-control-plane;
    container protocol-srgb {
        if-feature sr-mpls:protocol-srgb;
        uses sr-cmn:srgb;
        description
            "Per-protocol SRGB.";
    }
}

```

```

augment "/rt:routing/rt:control-plane-protocols/"
  + "rt:control-plane-protocol/ospf:ospf/"
  + "ospf:areas/ospf:area/ospf:interfaces/ospf:interface" {
when "../.../.../rt:type = 'ospf:ospfv2' or "
  + "../.../.../rt:type = 'ospf:ospfv3'" {
  description
    "This augments the OSPF interface configuration
    when used.";
}
description
  "This augments the OSPF protocol interface
  configuration with segment routing.";

  uses sr-mpls:igp-interface;
}

augment "/rt:routing/rt:control-plane-protocols/"
  + "rt:control-plane-protocol/ospf:ospf/"
  + "ospf:areas/ospf:area/ospf:interfaces/ospf:interface/"
  + "ospf:fast-reroute" {
when "../.../.../rt:type = 'ospf:ospfv2' or "
  + "../.../.../rt:type = 'ospf:ospfv3'" {
  description
    "This augments the OSPF routing protocol when used.";
}
description
  "This augments the OSPF protocol IP-FRR with TI-LFA.";

  container ti-lfa {
    if-feature ti-lfa;
    leaf enable {
      type boolean;
      description
        "Enables TI-LFA computation.";
    }
    description
      "Topology Independent Loop Free Alternate
      (TI-LFA) support.";
  }
}

/* Database */
augment "/rt:routing/"
  + "rt:control-plane-protocols/rt:control-plane-protocol/"
  + "ospf:ospf/ospf:areas/ospf:area/"
  + "ospf:interfaces/ospf:interface/ospf:database/"
  + "ospf:link-scope-lsa-type/ospf:link-scope-lsas/"
  + "ospf:link-scope-lsa/ospf:version/ospf:ospfv2/"

```

```

        + "ospf:ospfv2/ospf:body/ospf:opaque/"
        + "ospf:extended-prefix-opaque" {
when "../.../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 extended prefix LSA
        in type 9 opaque LSA.";

    uses extended-prefix-range-tlvs;
}

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:areas/"
    + "ospf:area/ospf:database/"
    + "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
    + "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/"
    + "ospf:extended-prefix-opaque" {
when "../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 extended prefix LSA
        in type 10 opaque LSA.";

    uses extended-prefix-range-tlvs;
}

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:database/"
    + "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
    + "ospf:as-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/"
    + "ospf:extended-prefix-opaque" {
when "../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 extended prefix LSA

```

```

        in type 11 opaque LSA.";

    uses extended-prefix-range-tlvs;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/ospf:area/"
+ "ospf:interfaces/ospf:interface/ospf:database/"
+ "ospf:link-scope-lsa-type/ospf:link-scope-lsas/"
+ "ospf:link-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:extended-prefix-opaque/ospf:extended-prefix-tlv" {
when "../../../../../../../../../../../../../../../"
+ "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
}
description
    "SR specific TLVs for OSPFv2 extended prefix TLV
    in type 9 opaque LSA.";
    uses prefix-sid-sub-tlvs;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/"
+ "ospf:area/ospf:database/"
+ "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
+ "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
+ "ospf:ospfv2/ospf:body/ospf:opaque/"
+ "ospf:extended-prefix-opaque/ospf:extended-prefix-tlv" {
when "../../../../../../../../../../../"
+ "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
}
description
    "SR specific TLVs for OSPFv2 extended prefix TLV
    in type 10 opaque LSA.";
    uses prefix-sid-sub-tlvs;
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:database/"
+ "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
+ "ospf:as-scope-lsa/ospf:version/ospf:ospfv2/"

```



```

    + "ospf:ospfv2/ospf:body/ospf:opaque/"
    + "ospf:extended-prefix-opaque/ospf:extended-prefix-tlv" {
when "../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 extended prefix TLV
        in type 11 opaque LSA.";
    uses prefix-sid-sub-tlvs;
}

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:areas/"
    + "ospf:area/ospf:database/"
    + "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
    + "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/"
    + "ospf:extended-link-opaque/ospf:extended-link-tlv" {
when "../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
    description
        "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 extended link TLV
        in type 10 opaque LSA.";

    container adj-sid-sub-tlvs {
        description "Adjacency SID optional sub-TLVs.";
        list adj-sid-sub-tlv {
            description "List of Adjacency SID sub-TLVs.";
            container adj-sid-flags {
                leaf-list bits {
                    type identityref {
                        base adj-sid-bit;
                    }
                    description "Adj sid sub-tlv flags list.";
                }
                description "Adj-sid sub-tlv flags.";
            }
            leaf mt-id {
                type uint8;
                description "Multi-topology ID.";
            }
            leaf weight {

```

```

        type uint8;
        description "Weight used for load-balancing.";
    }
    leaf sid {
        type uint32;
        description "Segment Identifier (SID) index/label.";
    }
}

container lan-adj-sid-sub-tlvs {
    description "LAN Adjacency SID optional sub-TLVs.";
    list lan-adj-sid-sub-tlv {
        description "List of LAN adjacency SID sub-TLVs.";
        container lan-adj-sid-flags {
            leaf-list bits {
                type identityref {
                    base adj-sid-bit;
                }
                description "LAN adj sid sub-tlv flags list.";
            }
            description "LAN adj-sid sub-tlv flags.";
        }
        leaf mt-id {
            type uint8;
            description "Multi-topology ID.";
        }
        leaf weight {
            type uint8;
            description "Weight used for load-balancing.";
        }
        leaf neighbor-router-id {
            type yang:dotted-quad;
            description "Neighbor router ID.";
        }
        leaf sid {
            type uint32;
            description "Segment Identifier (SID) index/label.";
        }
    }
}

augment "/rt:routing/"
+ "rt:control-plane-protocols/rt:control-plane-protocol/"
+ "ospf:ospf/ospf:areas/ospf:area/"
+ "ospf:interfaces/ospf:interface/ospf:database/"
+ "ospf:link-scope-lsa-type/ospf:link-scope-lsas/"

```

```

        + "ospf:link-scope-lsa/ospf:version/ospf:ospfv2/"
        + "ospf:ospfv2/ospf:body/ospf:opaque/ospf:ri-opaque" {
when "../.../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }

description
    "SR specific TLVs for OSPFv2 type 9 opaque LSA.";

uses sr-algorithm-tlv;
uses sid-range-tlvs;
uses local-block-tlvs;
uses srms-preference-tlv;
}

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:areas/"
    + "ospf:area/ospf:database/"
    + "ospf:area-scope-lsa-type/ospf:area-scope-lsas/"
    + "ospf:area-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/ospf:ri-opaque" {
when "../.../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }

description
    "SR specific TLVs for OSPFv2 type 10 opaque LSA.";

uses sr-algorithm-tlv;
uses sid-range-tlvs;
uses local-block-tlvs;
uses srms-preference-tlv;
}

augment "/rt:routing/"
    + "rt:control-plane-protocols/rt:control-plane-protocol/"
    + "ospf:ospf/ospf:database/"
    + "ospf:as-scope-lsa-type/ospf:as-scope-lsas/"
    + "ospf:as-scope-lsa/ospf:version/ospf:ospfv2/"
    + "ospf:ospfv2/ospf:body/ospf:opaque/ospf:ri-opaque" {
when "../.../.../.../.../.../.../.../.../.../.../..."
    + "rt:type = 'ospf:ospfv2'" {
        description

```

```

        "This augmentation is only valid for OSPFv2.";
    }
    description
        "SR specific TLVs for OSPFv2 type 11 opaque LSA.";

    uses sr-algorithm-tlv;
    uses sid-range-tlvs;
    uses local-block-tlvs;
    uses srms-preference-tlv;
}
}
<CODE ENDS>

```

4. Security Considerations

The YANG modules specified in this document define a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC8446].

The NETCONF Configuration Access Control model (NACM) [RFC8341] provides the means to restrict access for particular NETCONF or RESTCONF users to a pre-configured subset of all available NETCONF or RESTCONF protocol operations and content.

There are a number of data nodes defined in the modules that are writable/creatable/deletable (i.e., config true, which is the default). These data nodes may be considered sensitive or vulnerable in some network environments. Write operations (e.g., edit-config) to these data nodes without proper protection can have a negative effect on network operations. These are the subtrees and data nodes and their sensitivity/vulnerability:

/ospf:ospf/segment-routing/enabled - Modification to the enablement for SR could result in a Denial-of-Service (Dos) attack. If an attacker disables SR, it will cause traffic disruption.

/ospf:ospf/segment-routing/bindings - Modification to the local bindings could result in a Denial-of-Service (Dos) attack.

/ospf:ospf/protocol-srgb - Modification of the protocol SRGB could be used to mount a DoS attack. For example, if the protocol SRBG size is reduced to a very small value, a lot of existing segments could no longer be installed leading to a traffic disruption.

/ospf:interfaces/ospf:interface/segment-routing - Modification of the Adjacency Segment Identifier (Adj-SID) could be used to mount a DoS attack. Change of an Adj-SID could be used to redirect traffic.

/ospf:interfaces/ospf:interface/ospf:fast-reroute/ti-lfa - Modification of the TI-LFA enablement could lead to traffic disruption.

Some of the readable data nodes in the modules may be considered sensitive or vulnerable in some network environments. It is thus important to control read access (e.g., via get, get-config, or notification) to these data nodes.

Both module ietf-ospf-sr and ietf-ospf-msd augment base OSPF module data base with various TLVs. Knowledge of these data nodes can be used to attack other routers in the OSPF domain.

5. Acknowledgements

The authors wish to thank Yi Yang, Alexander Clemm, Gaurav Gupta, Ladislav Lhotka, Stephane Litkowski, Greg Hankins, Manish Gupta and Alan Davey for their thorough reviews and helpful comments.

This document was produced using Marshall Rose's xml2rfc tool.

Author affiliation with The MITRE Corporation is provided for identification purposes only, and is not intended to convey or imply MITRE's concurrence with, or support for, the positions, opinions or viewpoints expressed. MITRE has approved this document for Public Release, Distribution Unlimited, with Public Release Case Number 18-3281.

6. IANA Considerations

This document registers a URI in the IETF XML registry [RFC3688]. Following the format in [RFC3688], the following registration is requested to be made:

URI: urn:ietf:params:xml:ns:yang:ietf-ospf-sr
Registrant Contact: The IESG.
XML: N/A, the requested URI is an XML namespace.

URI: urn:ietf:params:xml:ns:yang:ietf-ospf-msd
Registrant Contact: The IESG.
XML: N/A, the requested URI is an XML namespace.

This document registers a YANG module in the YANG Module Names registry [RFC6020].

```
name: ietf-ospf-sr
namespace: urn:ietf:params:xml:ns:yang:ietf-ospf-sr
prefix: ospf-sr
reference: RFC XXXX

name: ietf-ospf-msd
namespace: urn:ietf:params:xml:ns:yang:ietf-ospf-msd
prefix: ospf-msd
reference: RFC XXXX
```

7. References

7.1. Normative References

- [I-D.ietf-ospf-yang]
Yeung, D., Qu, Y., Zhang, J., Chen, I., and A. Lindem,
"YANG Data Model for OSPF Protocol", Work in Progress,
Internet-Draft, draft-ietf-ospf-yang-29, 17 October 2019,
<<https://www.ietf.org/archive/id/draft-ietf-ospf-yang-29.txt>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate
Requirement Levels", BCP 14, RFC 2119,
DOI 10.17487/RFC2119, March 1997,
<<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328,
DOI 10.17487/RFC2328, April 1998,
<<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688,
DOI 10.17487/RFC3688, January 2004,
<<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC4750] Joyal, D., Ed., Galecki, P., Ed., Giacalone, S., Ed.,
Coltun, R., and F. Baker, "OSPF Version 2 Management
Information Base", RFC 4750, DOI 10.17487/RFC4750,
December 2006, <<https://www.rfc-editor.org/info/rfc4750>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF
for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008,
<<https://www.rfc-editor.org/info/rfc5340>>.

- [RFC5643] Joyal, D., Ed. and V. Manral, Ed., "Management Information Base for OSPFv3", RFC 5643, DOI 10.17487/RFC5643, August 2009, <<https://www.rfc-editor.org/info/rfc5643>>.
- [RFC5838] Lindem, A., Ed., Mirtorabi, S., Roy, A., Barnes, M., and R. Aggarwal, "Support of Address Families in OSPFv3", RFC 5838, DOI 10.17487/RFC5838, April 2010, <<https://www.rfc-editor.org/info/rfc5838>>.
- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.
- [RFC7223] Bjorklund, M., "A YANG Data Model for Interface Management", RFC 7223, DOI 10.17487/RFC7223, May 2014, <<https://www.rfc-editor.org/info/rfc7223>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8341] Bierman, A. and M. Bjorklund, "Network Configuration Access Control Model", STD 91, RFC 8341, DOI 10.17487/RFC8341, March 2018, <<https://www.rfc-editor.org/info/rfc8341>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.
- [RFC8476] Tantsura, J., Chunduri, U., Aldrin, S., and P. Psenak, "Signaling Maximum SID Depth (MSD) Using OSPF", RFC 8476, DOI 10.17487/RFC8476, December 2018, <<https://www.rfc-editor.org/info/rfc8476>>.

- [RFC8665] Psenak, P., Ed., Previdi, S., Ed., Filsfils, C., Gredler, H., Shakir, R., Henderickx, W., and J. Tantsura, "OSPF Extensions for Segment Routing", RFC 8665, DOI 10.17487/RFC8665, December 2019, <<https://www.rfc-editor.org/info/rfc8665>>.
- [RFC9020] Litkowski, S., Qu, Y., Lindem, A., Sarkar, P., and J. Tantsura, "YANG Data Model for Segment Routing", RFC 9020, DOI 10.17487/RFC9020, May 2021, <<https://www.rfc-editor.org/info/rfc9020>>.

7.2. Informative References

- [RFC8022] Lhotka, L. and A. Lindem, "A YANG Data Model for Routing Management", RFC 8022, DOI 10.17487/RFC8022, November 2016, <<https://www.rfc-editor.org/info/rfc8022>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.
- [RFC8342] Bjorklund, M., Schoenwaelder, J., Shafer, P., Watsen, K., and R. Wilton, "Network Management Datastore Architecture (NMDA)", RFC 8342, DOI 10.17487/RFC8342, March 2018, <<https://www.rfc-editor.org/info/rfc8342>>.

Appendix A. Contributors' Addreses

Dean Bogdanovic
Volta Networks, Inc.

EMail: dean@voltanet.io

Kiran Koushik Agrahara Sreenivasa
Cisco Systems
12515 Research Blvd, Bldg 4
Austin, TX 78681
USA

EMail: kkoushik@cisco.com

Authors' Addresses

Derek Yeung
Arrcus

Email: derek@arrcus.com

Yingzhen Qu
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
United States of America

Email: yingzhen.qu@futurewei.com

Jeffrey Zhang
Juniper Networks
10 Technology Park Drive
Westford, MA 01886
United States of America

Email: zzhang@juniper.net

Ing-Wher Chen
The MITRE Corporation

Email: ingwherchen@mitre.org

Acee Lindem
Cisco Systems
301 Midenhall Way
Cary, NC 27513

Email: acee@cisco.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 1, 2021

P. Psenak, Ed.
L. Ginsberg
Cisco Systems
W. Henderickx
Nokia
J. Tantsura
Apstra
J. Drake
Juniper Networks
June 30, 2020

OSPF Application-Specific Link Attributes
draft-ietf-ospf-te-link-attr-reuse-16.txt

Abstract

Existing traffic engineering related link attribute advertisements have been defined and are used in RSVP-TE deployments. Since the original RSVP-TE use case was defined, additional applications (e.g., Segment Routing Policy, Loop Free Alternate) have been defined that also make use of the link attribute advertisements. In cases where multiple applications wish to make use of these link attributes the current advertisements do not support application specific values for a given attribute nor do they support indication of which applications are using the advertised value for a given link. This document introduces new link attribute advertisements in OSPFv2 and OSPFv3 that address both of these shortcomings.

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 <https://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 1, 2021.

Copyright Notice

Copyright (c) 2020 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 (<https://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
2. Requirements Language	4
3. Requirements Discussion	4
4. Existing Advertisement of Link Attributes	5
5. Advertisement of Link Attributes	5
5.1. OSPFv2 Extended Link Opaque LSA and OSPFv3 E-Router-LSA	5
6. Advertisement of Application-Specific Values	6
7. Reused TE link attributes	9
7.1. Shared Risk Link Group (SRLG)	10
7.2. Extended Metrics	10
7.3. Administrative Group	11
7.4. Traffic Engineering Metric	11
8. Maximum Link Bandwidth	11
9. Considerations for Extended TE Metrics	12
10. Local Interface IPv6 Address Sub-TLV	12
11. Remote Interface IPv6 Address Sub-TLV	12
12. Attribute Advertisements and Enablement	13
13. Deployment Considerations	14
13.1. Use of Legacy RSVP-TE LSA Advertisements	14
13.2. Interoperability, Backwards Compatibility and Migration Concerns	15
13.2.1. Multiple Applications: Common Attributes with RSVP-TE	15
13.2.2. Multiple Applications: Some Attributes Not Shared with RSVP-TE	15
13.2.3. Interoperability with Legacy Routers	15
13.2.4. Use of Application-Specific Advertisements for RSVP-TE	16
14. Security Considerations	16
15. IANA Considerations	17
15.1. OSPFv2	17

15.2. OSPFv3	18
16. Contributors	19
17. Acknowledgments	19
18. References	19
18.1. Normative References	19
18.2. Informative References	21
Authors' Addresses	22

1. Introduction

Advertisement of link attributes by the OSPFv2 [RFC2328] and OSPFv3 [RFC5340] protocols in support of traffic engineering (TE) was introduced by [RFC3630] and [RFC5329] respectively. It has been extended by [RFC4203], [RFC7308] and [RFC7471]. Use of these extensions has been associated with deployments supporting Traffic Engineering over Multiprotocol Label Switching (MPLS) in the presence of the Resource Reservation Protocol (RSVP) - more succinctly referred to as RSVP-TE [RFC3209].

For the purposes of this document an application is a technology that makes use of link attribute advertisements, examples of which are listed in Section 6.

In recent years new applications have been introduced that have use cases for many of the link attributes historically used by RSVP-TE. Such applications include Segment Routing (SR) Policy [I-D.ietf-spring-segment-routing-policy] and Loop Free Alternates (LFA) [RFC5286]. This has introduced ambiguity in that if a deployment includes a mix of RSVP-TE support and SR Policy support (for example) it is not possible to unambiguously indicate which advertisements are to be used by RSVP-TE and which advertisements are to be used by SR Policy. If the topologies are fully congruent this may not be an issue, but any incongruence leads to ambiguity.

An example where this ambiguity causes a problem is a network in that RSVP-TE is enabled only on a subset of its links. A link attribute is advertised for the purpose of another application (e.g. SR Policy) for a link that is not enabled for RSVP-TE. As soon as the router that is an RSVP-TE head-end sees the link attribute being advertised for that link, it assumes RSVP-TE is enabled on that link, even though it is not. If such RSVP-TE head-end router tries to setup an RSVP-TE path via that link, it will result in the path setup failure.

An additional issue arises in cases where both applications are supported on a link but the link attribute values associated with each application differ. Current advertisements do not support

advertising application-specific values for the same attribute on a specific link.

This document defines extensions that address these issues. Also, as evolution of use cases for link attributes can be expected to continue in the years to come, this document defines a solution that is easily extensible for the introduction of new applications and new use cases.

2. Requirements Language

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. Requirements Discussion

As stated previously, evolution of use cases for link attributes can be expected to continue. Therefore, any discussion of existing use cases is limited to requirements that are known at the time of this writing. However, in order to determine the functionality required beyond what already exists in OSPF, it is only necessary to discuss use cases that justify the key points identified in the introduction, which are:

1. Support for indicating which applications are using the link attribute advertisements on a link
2. Support for advertising application-specific values for the same attribute on a link

[RFC7855] discusses use cases/requirements for Segment Routing (SR). Included among these use cases is SR Policy which is defined in [I-D.ietf-spring-segment-routing-policy]. If both RSVP-TE and SR Policy are deployed in a network, link attribute advertisements can be used by one or both of these applications. As there is no requirement for the link attributes advertised on a given link used by SR Policy to be identical to the link attributes advertised on that same link used by RSVP-TE, there is a clear requirement to indicate independently which link attribute advertisements are to be used by each application.

As the number of applications that may wish to utilize link attributes may grow in the future, an additional requirement is that the extensions defined allow the association of additional

applications to link attributes without altering the format of the advertisements or introducing new backwards compatibility issues.

Finally, there may still be many cases where a single attribute value can be shared among multiple applications, so the solution must minimize advertising duplicate link/attribute pairs whenever possible.

4. Existing Advertisement of Link Attributes

There are existing advertisements used in support of RSVP-TE. These advertisements are carried in the OSPFv2 TE Opaque LSA [RFC3630] and OSPFv3 Intra-Area-TE-LSA [RFC5329]. Additional RSVP-TE link attributes have been defined by [RFC4203], [RFC7308] and [RFC7471].

Extended Link Opaque LSAs as defined in [RFC7684] for OSPFv2 and Extended Router-LSAs [RFC8362] for OSPFv3 are used to advertise link attributes that are used by applications other than RSVP-TE or GMPLS [RFC4203]. These LSAs were defined as a generic containers for distribution of the extended link attributes.

5. Advertisement of Link Attributes

This section outlines the solution for advertising link attributes originally defined for RSVP-TE or GMPLS when they are used for other applications.

5.1. OSPFv2 Extended Link Opaque LSA and OSPFv3 E-Router-LSA

Advantages of Extended Link Opaque LSAs as defined in [RFC7684] for OSPFv2 and Extended Router-LSAs [RFC8362] for OSPFv3 with respect to advertisement of link attributes originally defined for RSVP-TE when used in packet networks and in GMPLS:

1. Advertisement of the link attributes does not make the link part of the RSVP-TE topology. It avoids any conflicts and is fully compatible with [RFC3630] and [RFC5329].
2. The OSPFv2 TE Opaque LSA and OSPFv3 Intra-Area-TE-LSA remains truly opaque to OSPFv2 and OSPFv3 as originally defined in [RFC3630] and [RFC5329] respectively. Their contents are not inspected by OSPF, which instead acts as a pure transport.
3. There is a clear distinction between link attributes used by RSVP-TE and link attributes used by other OSPFv2 or OSPFv3 applications.

4. All link attributes that are used by other applications are advertised in a single LSA, the Extended Link Opaque LSA in OSPFv2 or the OSPFv3 E-Router-LSA [RFC8362] in OSPFv3.

The disadvantage of this approach is that in rare cases, the same link attribute is advertised in both the TE Opaque and Extended Link Attribute LSAs in OSPFv2 or the Intra-Area-TE-LSA and E-Router-LSA in OSPFv3.

Extended Link Opaque LSA [RFC7684] and E-Router-LSA [RFC8362] are used to advertise any link attributes used for non-RSVP-TE applications in OSPFv2 or OSPFv3 respectively, including those that have been originally defined for RSVP-TE applications (See Section 7).

TE link attributes used for RSVP-TE/GMPLS continue to use OSPFv2 TE Opaque LSA [RFC3630] and OSPFv3 Intra-Area-TE-LSA [RFC5329].

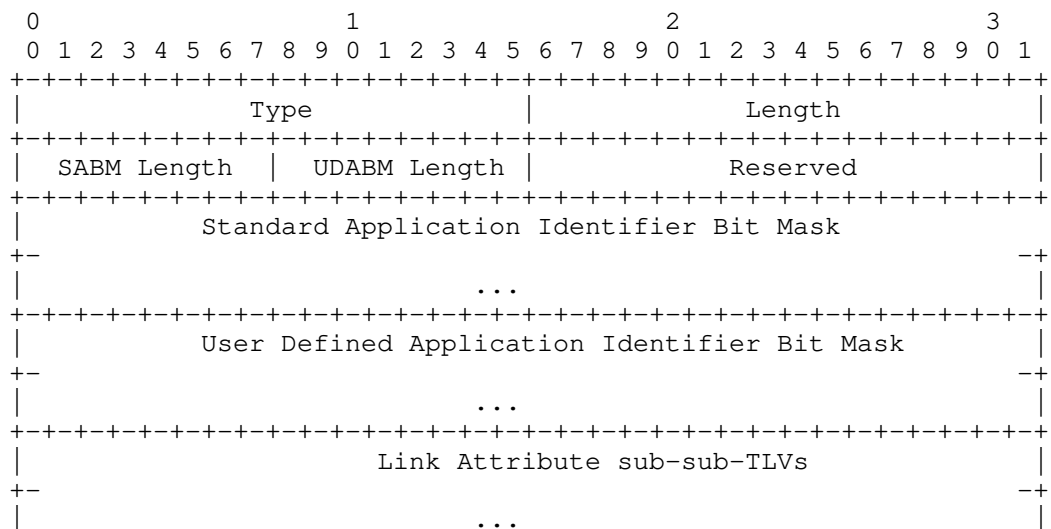
The format of the link attribute TLVs that have been defined for RSVP-TE applications will be kept unchanged even when they are used for non-RSVP-TE applications. Unique code points are allocated for these link attribute TLVs from the OSPFv2 Extended Link TLV Sub-TLV Registry [RFC7684] and from the OSPFv3 Extended-LSA Sub-TLV Registry [RFC8362], as specified in Section 15.

6. Advertisement of Application-Specific Values

To allow advertisement of the application-specific values of the link attribute, a new Application-Specific Link Attributes (ASLA) sub-TLV is defined. The ASLA sub-TLV is a sub-TLV of the OSPFv2 Extended Link TLV [RFC7684] and OSPFv3 Router-Link TLV [RFC8362].

On top of advertising the link attributes for standardized applications, link attributes can be advertised for the purpose of applications that are not standardized. We call such an application a "User Defined Application" or "UDA". These applications are not subject to standardization and are outside of the scope of this specification.

The ASLA sub-TLV is an optional sub-TLV of OSPFv2 Extended Link TLV and OSPFv3 Router-Link TLV. Multiple ASLA sub-TLVs can be present in its parent TLV when different applications want to control different link attributes or when different value of the same attribute needs to be advertised by multiple applications. The ASLA sub-TLV MUST be used for advertisement of the link attributes listed at the end on this section if these are advertised inside OSPFv2 Extended Link TLV and OSPFv3 Router-Link TLV. It has the following format:



where:

Type: 10 (OSPFv2), 11 (OSPFv3)

Length: variable

SABM Length: Standard Application Identifier Bit Mask Length in octets. The value MUST be 0, 4 or 8. If the Standard Application Bit Mask is not present, the Standard Application Bit Mask Length MUST be set to 0.

UDABM Length: User Defined Application Identifier Bit Mask Length in octets. The value MUST be 0, 4 or 8. If the User Defined Application Bit Mask is not present, the User Defined Application Bit Mask Length MUST be set to 0.

Standard Application Identifier Bit Mask: Optional set of bits, where each bit represents a single standard application. Bits are defined in the Link Attribute Application Identifier Registry, which has been defined in [I-D.ietf-isis-te-app]. Current assignments are repeated here for informational purpose:

0	1	2	3	4	5	6	7	...
								...
R	S	F						...
								...

Bit-0 (R-bit): RSVP-TE

Bit-1 (S-bit): Segment Routing Policy

Bit-2 (F-bit): Loop Free Alternate (LFA). Includes all LFA types

User Defined Application Identifier Bit Mask: Optional set of bits, where each bit represents a single user defined application.

If the SABM or UDABM length is other than 0, 4, or 8, the ASLA sub-TLV MUST be ignored by the receiver.

Standard Application Identifier Bits are defined/sent starting with Bit 0. Undefined bits that are transmitted MUST be transmitted as 0 and MUST be ignored on receipt. Bits that are not transmitted MUST be treated as if they are set to 0 on receipt. Bits that are not supported by an implementation MUST be ignored on receipt.

User Defined Application Identifier Bits have no relationship to Standard Application Identifier Bits and are not managed by IANA or any other standards body. It is recommended that bits are used starting with Bit 0 so as to minimize the number of octets required to advertise all UDAs. Undefined bits which are transmitted MUST be transmitted as 0 and MUST be ignored on receipt. Bits that are not transmitted MUST be treated as if they are set to 0 on receipt. Bits that are not supported by an implementation MUST be ignored on receipt.

If the link attribute advertisement is intended to be only used by a specific set of applications, corresponding Bit Masks MUST be present and application-specific bit(s) MUST be set for all applications that use the link attributes advertised in the ASLA sub-TLV.

Application Bit Masks apply to all link attributes that support application-specific values and are advertised in the ASLA sub-TLV.

The advantage of not making the Application Bit Masks part of the attribute advertisement itself is that the format of any previously defined link attributes can be kept and reused when advertising them in the ASLA sub-TLV.

If the same attribute is advertised in more than one ASLA sub-TLVs with the application listed in the Application Bit Masks, the application SHOULD use the first instance of advertisement and ignore any subsequent advertisements of that attribute.

If link attributes are advertised with zero length Application Identifier Bit Masks for both standard applications and user defined applications, then any Standard Application and/or any User Defined

Application is permitted to use that set of link attributes. If support for a new application is introduced on any node in a network in the presence of such advertisements, these advertisements are permitted to be used by the new application. If this is not what is intended, then existing advertisements MUST be readvertised with an explicit set of applications specified before a new application is introduced.

An application-specific advertisement (Application Identifier Bit Mask with a matching Application Identifier Bit set) for an attribute MUST always be preferred over the advertisement of the same attribute with the zero length Application Identifier Bit Masks for both standard applications and user defined applications on the same link.

This document defines the initial set of link attributes that MUST use the ASLA sub-TLV if advertised in the OSPFv2 Extended Link TLV or in the OSPFv3 Router-Link TLV. Documents which define new link attributes MUST state whether the new attributes support application-specific values and as such are advertised in an ASLA sub-TLV. The standard link attributes that are advertised in ASLA sub-TLVs are:

- Shared Risk Link Group [RFC4203]
- Unidirectional Link Delay [RFC7471]
- Min/Max Unidirectional Link Delay [RFC7471]
- Unidirectional Delay Variation [RFC7471]
- Unidirectional Link Loss [RFC7471]
- Unidirectional Residual Bandwidth [RFC7471]
- Unidirectional Available Bandwidth [RFC7471]
- Unidirectional Utilized Bandwidth [RFC7471]
- Administrative Group [RFC3630]
- Extended Administrative Group [RFC7308]
- TE Metric [RFC3630]

7. Reused TE link attributes

This section defines the use case and indicates the code points (Section 15) from the OSPFv2 Extended Link TLV Sub-TLV Registry and

OSPFv3 Extended-LSA Sub-TLV Registry for some of the link attributes that have been originally defined for RSVP-TE or GMPLS.

7.1. Shared Risk Link Group (SRLG)

The SRLG of a link can be used in OSPF calculated IPFRR (IP Fast Reroute) [RFC5714] to compute a backup path that does not share any SRLG group with the protected link.

To advertise the SRLG of the link in the OSPFv2 Extended Link TLV, the same format for the sub-TLV defined in section 1.3 of [RFC4203] is used and TLV type 11 is used. Similarly, for OSPFv3 to advertise the SRLG in the OSPFv3 Router-Link TLV, TLV type 12 is used.

7.2. Extended Metrics

[RFC3630] defines several link bandwidth types. [RFC7471] defines extended link metrics that are based on link bandwidth, delay and loss characteristics. All of these can be used to compute primary and backup paths within an OSPF area to satisfy requirements for bandwidth, delay (nominal or worst case) or loss.

To advertise extended link metrics in the OSPFv2 Extended Link TLV, the same format for the sub-TLVs defined in [RFC7471] is used with the following TLV types:

- 12 - Unidirectional Link Delay
- 13 - Min/Max Unidirectional Link Delay
- 14 - Unidirectional Delay Variation
- 15 - Unidirectional Link Loss
- 16 - Unidirectional Residual Bandwidth
- 17 - Unidirectional Available Bandwidth
- 18 - Unidirectional Utilized Bandwidth

To advertise extended link metrics in the OSPFv3 Extended-LSA Router-Link TLV, the same format for the sub-TLVs defined in [RFC7471] is used with the following TLV types:

- 13 - Unidirectional Link Delay
- 14 - Min/Max Unidirectional Link Delay

- 15 - Unidirectional Delay Variation
- 16 - Unidirectional Link Loss
- 17 - Unidirectional Residual Bandwidth
- 18 - Unidirectional Available Bandwidth
- 19 - Unidirectional Utilized Bandwidth

7.3. Administrative Group

[RFC3630] and [RFC7308] define the Administrative Group and Extended Administrative Group sub-TLVs respectively.

To advertise the Administrative Group and Extended Administrative Group in the OSPFv2 Extended Link TLV, the same format for the sub-TLVs defined in [RFC3630] and [RFC7308] is used with the following TLV types:

- 19 - Administrative Group
- 20 - Extended Administrative Group

To advertise Administrative Group and Extended Administrative Group in the OSPFv3 Router-Link TLV, the same format for the sub-TLVs defined in [RFC3630] and [RFC7308] is used with the following TLV types:

- 20 - Administrative Group
- 21 - Extended Administrative Group

7.4. Traffic Engineering Metric

[RFC3630] defines Traffic Engineering Metric.

To advertise the Traffic Engineering Metric in the OSPFv2 Extended Link TLV, the same format for the sub-TLV defined in section 2.5.5 of [RFC3630] is used and TLV type 22 is used. Similarly, for OSPFv3 to advertise the Traffic Engineering Metric in the OSPFv3 Router-Link TLV, TLV type 22 is used.

8. Maximum Link Bandwidth

Maximum link bandwidth is an application independent attribute of the link that is defined in [RFC3630]. Because it is an application independent attribute, it MUST NOT be advertised in ASLA sub-TLV.

Instead, it MAY be advertised as a sub-TLV of the Extended Link Opaque LSA Extended Link TLV in OSPFv2 [RFC7684] or sub-TLV of OSPFv3 E-Router-LSA Router-Link TLV in OSPFv3 [RFC8362].

To advertise the Maximum link bandwidth in the OSPFv2 Extended Link TLV, the same format for sub-TLV defined in [RFC3630] is used with TLV type 23.

To advertise the Maximum link bandwidth in the OSPFv3 Router-Link TLV, the same format for sub-TLV defined in [RFC3630] is used with TLV type 23.

9. Considerations for Extended TE Metrics

[RFC7471] defines a number of dynamic performance metrics associated with a link. It is conceivable that such metrics could be measured specific to traffic associated with a specific application. Therefore this document includes support for advertising these link attributes specific to a given application. However, in practice it may well be more practical to have these metrics reflect the performance of all traffic on the link regardless of application. In such cases, advertisements for these attributes can be associated with all of the applications utilizing that link. This can be done either by explicitly specifying the applications in the Application Identifier Bit Mask or by using a zero length Application Identifier Bit Mask.

10. Local Interface IPv6 Address Sub-TLV

The Local Interface IPv6 Address Sub-TLV is an application independent attribute of the link that is defined in [RFC5329]. Because it is an application independent attribute, it MUST NOT be advertised in the ASLA sub-TLV. Instead, it MAY be advertised as a sub-TLV of the OSPFv3 E-Router-LSA Router-Link TLV [RFC8362].

To advertise the Local Interface IPv6 Address Sub-TLV in the OSPFv3 Router-Link TLV, the same format for sub-TLV defined in [RFC5329] is used with TLV type 24.

11. Remote Interface IPv6 Address Sub-TLV

The Remote Interface IPv6 Address Sub-TLV is an application independent attribute of the link that is defined in [RFC5329]. Because it is an application independent attribute, it MUST NOT be advertised in the ASLA sub-TLV. Instead, it MAY be advertised as a sub-TLV of the OSPFv3 E-Router-LSA Router-Link TLV [RFC8362].

To advertise the Remote Interface IPv6 Address Sub-TLV in the OSPFv3 Router-Link TLV, the same format for sub-TLV defined in [RFC5329] is used with TLV type 25.

12. Attribute Advertisements and Enablement

This document defines extensions to support the advertisement of application-specific link attributes.

There are applications where the application enablement on the link is relevant - e.g., RSVP-TE - one needs to make sure that RSVP is enabled on the link before sending a RSVP-TE signaling message over it.

There are applications where the enablement of the application on the link is irrelevant and has nothing to do with the fact that some link attributes are advertised for the purpose of such application. An example of this is LFA.

Whether the presence of link attribute advertisements for a given application indicates that the application is enabled on that link depends upon the application. Similarly, whether the absence of link attribute advertisements indicates that the application is not enabled depends upon the application.

In the case of RSVP-TE, the advertisement of application-specific link attributes has no implication of RSVP-TE being enabled on that link. The RSVP-TE enablement is solely derived from the information carried in the OSPFv2 TE Opaque LSA [RFC3630] and OSPFv3 Intra-Area-TE-LSA [RFC5329].

In the case of SR Policy, advertisement of application-specific link attributes does not indicate enablement of SR Policy. The advertisements are only used to support constraints that may be applied when specifying an explicit path. SR Policy is implicitly enabled on all links that are part of the Segment Routing enabled topology independent of the existence of link attribute advertisements

In the case of LFA, advertisement of application-specific link attributes does not indicate enablement of LFA on that link. Enablement is controlled by local configuration.

If, in the future, additional standard applications are defined to use this mechanism, the specification defining this use MUST define the relationship between application-specific link attribute advertisements and enablement for that application.

This document allows the advertisement of application-specific link attributes with no application identifiers i.e., both the Standard Application Identifier Bit Mask and the User Defined Application Identifier Bit Mask are not present (See Section 6). This supports the use of the link attribute by any application. In the presence of an application where the advertisement of link attribute advertisements is used to infer the enablement of an application on that link (e.g., RSVP-TE), the absence of the application identifier leaves ambiguous whether that application is enabled on such a link. This needs to be considered when making use of the "any application" encoding.

13. Deployment Considerations

13.1. Use of Legacy RSVP-TE LSA Advertisements

Bit Identifiers for Standard Applications are defined in Section 6. All of the identifiers defined in this document are associated with applications that were already deployed in some networks prior to the writing of this document. Therefore, such applications have been deployed using the RSVP-TE LSA advertisements. The Standard Applications defined in this document may continue to use RSVP-TE LSA advertisements for a given link so long as at least one of the following conditions is true:

The application is RSVP-TE

The application is SR Policy or LFA and RSVP-TE is not deployed anywhere in the network

The application is SR Policy or LFA, RSVP-TE is deployed in the network, and both the set of links on which SR Policy and/or LFA advertisements are required and the attribute values used by SR Policy and/or LFA on all such links is fully congruent with the links and attribute values used by RSVP-TE

Under the conditions defined above, implementations that support the extensions defined in this document have the choice of using RSVP-TE LSA advertisements or application-specific advertisements in support of SR Policy and/or LFA. This will require implementations to provide controls specifying which type of advertisements are to be sent/ processed on receive for these applications. Further discussion of the associated issues can be found in Section 13.2.

New applications that future documents define to make use of the advertisements defined in this document MUST NOT make use of RSVP-TE LSA advertisements. This simplifies deployment of new applications

by eliminating the need to support multiple ways to advertise attributes for the new applications.

13.2. Interoperability, Backwards Compatibility and Migration Concerns

Existing deployments of RSVP-TE, SR Policy, and/or LFA utilize the legacy advertisements listed in Section 4. Routers which do not support the extensions defined in this document will only process legacy advertisements and are likely to infer that RSVP-TE is enabled on the links for which legacy advertisements exist. It is expected that deployments using the legacy advertisements will persist for a significant period of time. Therefore deployments using the extensions defined in this document in the presence of routers that do not support these extensions need to be able to interoperate with the use of legacy advertisements by the legacy routers. The following sub-sections discuss interoperability and backwards compatibility concerns for a number of deployment scenarios.

13.2.1. Multiple Applications: Common Attributes with RSVP-TE

In cases where multiple applications are utilizing a given link, one of the applications is RSVP-TE, and all link attributes for a given link are common to the set of applications utilizing that link, interoperability is achieved by using legacy advertisements for RSVP-TE. Attributes for applications other than RSVP-TE MUST be advertised using application-specific advertisements. This results in duplicate advertisements for those attributes.

13.2.2. Multiple Applications: Some Attributes Not Shared with RSVP-TE

In cases where one or more applications other than RSVP-TE are utilizing a given link and one or more link attribute values are not shared with RSVP-TE, interoperability is achieved by using legacy advertisements for RSVP-TE. Attributes for applications other than RSVP-TE MUST be advertised using application-specific advertisements. In cases where some link attributes are shared with RSVP-TE, this requires duplicate advertisements for those attributes

13.2.3. Interoperability with Legacy Routers

For the applications defined in this document, routers that do not support the extensions defined in this document will send and receive only legacy link attribute advertisements. So long as there is any legacy router in the network that has any of the applications enabled, all routers MUST continue to advertise link attributes using legacy advertisements. In addition, the link attribute values associated with the set of applications supported by legacy routers (RSVP-TE, SR Policy, and/or LFA) are always shared since legacy

routers have no way of advertising or processing application-specific values. Once all legacy routers have been upgraded, migration from legacy advertisements to application specific advertisements can be achieved via the following steps:

1) Send new application-specific advertisements while continuing to advertise using the legacy advertisement (all advertisements are then duplicated). Receiving routers continue to use legacy advertisements.

2) Enable the use of the application-specific advertisements on all routers

3) Keep legacy advertisements if needed for RSVP-TE purposes.

When the migration is complete, it then becomes possible to advertise incongruent values per application on a given link.

Documents defining new applications that make use of the application-specific advertisements defined in this document MUST discuss interoperability and backwards compatibility issues that could occur in the presence of routers that do not support the new application.

13.2.4. Use of Application-Specific Advertisements for RSVP-TE

The extensions defined in this document support RSVP-TE as one of the supported applications. It is however RECOMMENDED to advertise all link-attributes for RSVP-TE in the existing OSPFv2 TE Opaque LSA [RFC3630] and OSPFv3 Intra-Area-TE-LSA [RFC5329] to maintain backward compatibility. RSVP-TE can eventually utilize the application-specific advertisements for newly defined link attributes, that are defined as application-specific.

Link attributes that are not allowed to be advertised in the ASLA Sub-TLV, such as Maximum Reservable Link Bandwidth and Unreserved Bandwidth MUST use the OSPFv2 TE Opaque LSA [RFC3630] and OSPFv3 Intra-Area-TE-LSA [RFC5329] and MUST NOT be advertised in ASLA Sub-TLV.

14. Security Considerations

Existing security extensions as described in [RFC2328], [RFC5340] and [RFC8362] apply to extensions defined in this document. While OSPF is under a single administrative domain, there can be deployments where potential attackers have access to one or more networks in the OSPF routing domain. In these deployments, stronger authentication mechanisms such as those specified in [RFC5709], [RFC7474], [RFC4552] or [RFC7166] SHOULD be used.

Implementations must assure that malformed TLV and Sub-TLV defined in this document are detected and do not provide a vulnerability for attackers to crash the OSPF router or routing process. Reception of a malformed TLV or Sub-TLV SHOULD be counted and/or logged for further analysis. Logging of malformed TLVs and Sub-TLVs SHOULD be rate-limited to prevent a Denial of Service (DoS) attack (distributed or otherwise) from overloading the OSPF control plane.

This document defines a new way to advertise link attributes. Tampering with the information defined in this document may have an effect on applications using it, including impacting Traffic Engineering that uses various link attributes for its path computation. This is similar in nature to the impacts associated with (for example) [RFC3630]. As the advertisements defined in this document limit the scope to specific applications, the impact of tampering is similarly limited in scope.

15. IANA Considerations

This specifications updates two existing registries:

- OSPFv2 Extended Link TLV Sub-TLVs Registry
- OSPFv3 Extended-LSA Sub-TLV Registry

New values are allocated using the IETF Review procedure as described in [RFC5226].

15.1. OSPFv2

The OSPFv2 Extended Link TLV Sub-TLVs Registry [RFC7684] defines sub-TLVs at any level of nesting for OSPFv2 Extended Link TLVs. IANA has assigned the following Sub-TLV types from the OSPFv2 Extended Link TLV Sub-TLVs Registry:

- 10 - Application-Specific Link Attributes
- 11 - Shared Risk Link Group
- 12 - Unidirectional Link Delay
- 13 - Min/Max Unidirectional Link Delay
- 14 - Unidirectional Delay Variation
- 15 - Unidirectional Link Loss
- 16 - Unidirectional Residual Bandwidth

- 17 - Unidirectional Available Bandwidth
- 18 - Unidirectional Utilized Bandwidth
- 19 - Administrative Group
- 20 - Extended Administrative Group
- 22 - TE Metric
- 23 - Maximum Link Bandwidth

15.2. OSPFv3

The OSPFv3 Extended-LSA Sub-TLV Registry [RFC8362] defines sub-TLVs at any level of nesting for OSPFv3 Extended LSAs. IANA has assigned the following Sub-TLV types from the OSPFv3 Extended-LSA Sub-TLV Registry:

- 11 - Application-Specific Link Attributes
- 12 - Shared Risk Link Group
- 13 - Unidirectional Link Delay
- 14 - Min/Max Unidirectional Link Delay
- 15 - Unidirectional Delay Variation
- 16 - Unidirectional Link Loss
- 17 - Unidirectional Residual Bandwidth
- 18 - Unidirectional Available Bandwidth
- 19 - Unidirectional Utilized Bandwidth
- 20 - Administrative Group
- 21 - Extended Administrative Group
- 22 - TE Metric
- 23 - Maximum Link Bandwidth
- 24 - Local Interface IPv6 Address Sub-TLV
- 25 - Remote Interface IPv6 Address Sub-TLV

16. Contributors

The following people contributed to the content of this document and should be considered as co-authors:

Acee Lindem
Cisco Systems
301 Midenhall Way
Cary, NC 27513
USA

Email: acee@cisco.com

Ketan Talaulikar
Cisco Systems, Inc.
India

Email: ketant@cisco.com

Hannes Gredler
RtBrick Inc.
Austria

Email: hannes@rtbrick.com

17. Acknowledgments

Thanks to Chris Bowers for his review and comments.

Thanks to Alvaro Retana for his detailed review and comments.

18. References

18.1. Normative References

- [I-D.ietf-isis-te-app]
Ginsberg, L., Psenak, P., Previdi, S., Henderickx, W., and J. Drake, "IS-IS Application-Specific Link Attributes", draft-ietf-isis-te-app-19 (work in progress), June 2020.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.

- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", RFC 3630, DOI 10.17487/RFC3630, September 2003, <<https://www.rfc-editor.org/info/rfc3630>>.
- [RFC4203] Kompella, K., Ed. and Y. Rekhter, Ed., "OSPF Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 4203, DOI 10.17487/RFC4203, October 2005, <<https://www.rfc-editor.org/info/rfc4203>>.
- [RFC5329] Ishiguro, K., Manral, V., Davey, A., and A. Lindem, Ed., "Traffic Engineering Extensions to OSPF Version 3", RFC 5329, DOI 10.17487/RFC5329, September 2008, <<https://www.rfc-editor.org/info/rfc5329>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC7308] Osborne, E., "Extended Administrative Groups in MPLS Traffic Engineering (MPLS-TE)", RFC 7308, DOI 10.17487/RFC7308, July 2014, <<https://www.rfc-editor.org/info/rfc7308>>.
- [RFC7471] Giacalone, S., Ward, D., Drake, J., Atlas, A., and S. Previdi, "OSPF Traffic Engineering (TE) Metric Extensions", RFC 7471, DOI 10.17487/RFC7471, March 2015, <<https://www.rfc-editor.org/info/rfc7471>>.
- [RFC7684] Psenak, P., Gredler, H., Shakir, R., Henderickx, W., Tantsura, J., and A. Lindem, "OSPFv2 Prefix/Link Attribute Advertisement", RFC 7684, DOI 10.17487/RFC7684, November 2015, <<https://www.rfc-editor.org/info/rfc7684>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8362] Lindem, A., Roy, A., Goethals, D., Reddy Vallem, V., and F. Baker, "OSPFv3 Link State Advertisement (LSA) Extensibility", RFC 8362, DOI 10.17487/RFC8362, April 2018, <<https://www.rfc-editor.org/info/rfc8362>>.

18.2. Informative References

- [I-D.ietf-spring-segment-routing-policy]
Filsfils, C., Sivabalan, S., Voyer, D., Bogdanov, A., and P. Mattes, "Segment Routing Policy Architecture", draft-ietf-spring-segment-routing-policy-07 (work in progress), May 2020.
- [RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", RFC 3209, DOI 10.17487/RFC3209, December 2001, <<https://www.rfc-editor.org/info/rfc3209>>.
- [RFC4552] Gupta, M. and N. Melam, "Authentication/Confidentiality for OSPFv3", RFC 4552, DOI 10.17487/RFC4552, June 2006, <<https://www.rfc-editor.org/info/rfc4552>>.
- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", RFC 5226, DOI 10.17487/RFC5226, May 2008, <<https://www.rfc-editor.org/info/rfc5226>>.
- [RFC5286] Atlas, A., Ed. and A. Zinin, Ed., "Basic Specification for IP Fast Reroute: Loop-Free Alternates", RFC 5286, DOI 10.17487/RFC5286, September 2008, <<https://www.rfc-editor.org/info/rfc5286>>.
- [RFC5709] Bhatia, M., Manral, V., Fanto, M., White, R., Barnes, M., Li, T., and R. Atkinson, "OSPFv2 HMAC-SHA Cryptographic Authentication", RFC 5709, DOI 10.17487/RFC5709, October 2009, <<https://www.rfc-editor.org/info/rfc5709>>.
- [RFC5714] Shand, M. and S. Bryant, "IP Fast Reroute Framework", RFC 5714, DOI 10.17487/RFC5714, January 2010, <<https://www.rfc-editor.org/info/rfc5714>>.
- [RFC7166] Bhatia, M., Manral, V., and A. Lindem, "Supporting Authentication Trailer for OSPFv3", RFC 7166, DOI 10.17487/RFC7166, March 2014, <<https://www.rfc-editor.org/info/rfc7166>>.
- [RFC7474] Bhatia, M., Hartman, S., Zhang, D., and A. Lindem, Ed., "Security Extension for OSPFv2 When Using Manual Key Management", RFC 7474, DOI 10.17487/RFC7474, April 2015, <<https://www.rfc-editor.org/info/rfc7474>>.

Authors' Addresses

Peter Psenak (editor)
Cisco Systems
Eurovea Centre, Central 3
Pribinova Street 10
Bratislava 81109
Slovakia

Email: ppsenak@cisco.com

Les Ginsberg
Cisco Systems
821 Alder Drive
MILPITAS, CA 95035
USA

Email: ginsberg@cisco.com

Wim Henderickx
Nokia
Copernicuslaan 50
Antwerp, 2018 94089
Belgium

Email: wim.henderickx@nokia.com

Jeff Tantsura
Apstra
US

Email: jefftant.ietf@gmail.com

John Drake
Juniper Networks
1194 N. Mathilda Ave
Sunnyvale, California 94089
USA

Email: jdrake@juniper.net

Internet
Internet-Draft
Intended status: Standards Track
Expires: April 19, 2020

D. Yeung
Arrcus
Y. Qu
Futurewei
J. Zhang
Juniper Networks
I. Chen
The MITRE Corporation
A. Lindem
Cisco Systems
October 17, 2019

YANG Data Model for OSPF Protocol
draft-ietf-ospf-yang-29

Abstract

This document defines a YANG data model that can be used to configure and manage OSPF. The model is based on YANG1.1 as defined in RFC 7950 and conforms to the Network Management Datastore Architecture (NMDA) as described in RFC 8342.

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 <https://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 April 19, 2020.

Copyright Notice

Copyright (c) 2019 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 (<https://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. Overview	2
1.1. Requirements Language	3
1.2. Tree Diagrams	3
2. Design of Data Model	3
2.1. OSPF Operational State	3
2.2. Overview	4
2.3. OSPFv2 and OSPFv3	5
2.4. Optional Features	5
2.5. OSPF Router Configuration/Operational State	7
2.6. OSPF Area Configuration/Operational State	10
2.7. OSPF Interface Configuration/Operational State	16
2.8. OSPF Notifications	19
2.9. OSPF RPC Operations	23
3. OSPF YANG Module	23
4. Security Considerations	120
5. IANA Considerations	123
6. Acknowledgements	123
7. References	124
7.1. Normative References	124
7.2. Informative References	129
Appendix A. Contributors' Addresses	131
Authors' Addresses	131

1. Overview

YANG [RFC6020][RFC7950] is a data definition language used to define the contents of a conceptual data store that allows networked devices to be managed using NETCONF [RFC6241], RESTCONF [RFC8040], and other Network Management protocols. Furthermore, YANG data models can be used as the basis for implementation of other interfaces, such as CLI and programmatic APIs.

This document defines a YANG data model that can be used to configure and manage OSPF and it is an augmentation to the core routing data model. It fully conforms to the Network Management Datastore Architecture (NMDA) [RFC8342]. A core routing data model is defined in [RFC8349], and it provides the basis for the development of data models for routing protocols. The interface data model is defined in [RFC8343] and is used for referencing interfaces from the routing

protocol. The key-chain data model used for OSPF authentication is defined in [RFC8177] and provides both a reference to configured key-chains and an enumeration of cryptographic algorithms.

Both OSPFv2 [RFC2328] and OSPFv3 [RFC5340] are supported. In addition to the core OSPF protocol, features described in other OSPF RFCs are also supported. These includes demand circuit [RFC1793], traffic engineering [RFC3630], multiple address family [RFC5838], graceful restart [RFC3623] [RFC5187], NSSA [RFC3101], and OSPFv2 or OSPFv3 as a PE-CE Protocol [RFC4577], [RFC6565]. These non-core features are optional in the OSPF data model.

1.1. Requirements Language

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.

1.2. Tree Diagrams

This document uses the graphical representation of data models defined in [RFC8340].

2. Design of Data Model

Although the basis of OSPF configuration elements like routers, areas, and interfaces remains the same, the detailed configuration model varies among router vendors. Differences are observed in terms of how the protocol instance is tied to the routing domain and how multiple protocol instances are be instantiated among others.

The goal of this document is to define a data model that provides a common user interface to the OSPFv2 and OSPFv3 protocols. There is very little information that is designated as "mandatory", providing freedom for vendors to adapt this data model to their respective product implementations.

2.1. OSPF Operational State

The OSPF operational state is included in the same tree as OSPF configuration consistent with the Network Management Datastore Architecture [RFC8342]. Consequently, only the routing container in the ietf-routing model [RFC8349] is augmented. The routing-state container is not augmented.

2.2. Overview

The OSPF YANG module defined in this document has all the common building blocks for the OSPF protocol.

The OSPF YANG module augments the /routing/control-plane-protocols/control-plane-protocol path defined in the ietf-routing module. The ietf-ospf model defines a single instance of OSPF which may be instantiated as an OSPFv2 or OSPFv3 instance. Multiple instances are instantiated as multiple control-plane protocols instances.

```

module: ietf-ospf
  augment /rt:routing/rt:control-plane-protocols/
    rt:control-plane-protocol:
      +--rw ospf
        .
        .
        +--rw af?                               identityref
        .
        .
        +--rw areas
          +--rw area* [area-id]
            +--rw area-id                       area-id-type
            .
            .
            +--rw virtual-links
              +--rw virtual-link* [transit-area-id router-id]
              .
              .
            +--rw sham-links {pe-ce-protocol}?
              +--rw sham-link* [local-id remote-id]
              .
              .
            +--rw interfaces
              +--rw interface* [name]
              .
              .
        +--rw topologies {multi-topology}?
          +--rw topology* [name]
          .
          .

```

The ospf container includes one OSPF protocol instance. The instance includes OSPF router level configuration and operational state. Each OSPF instance maps to a control-plane-protocol instance as defined in [RFC8349].

The area and area/interface containers define the OSPF configuration and operational state for OSPF areas and interfaces respectively.

The topologies container defines the OSPF configuration and operational state for OSPF topologies when the multi-topology feature is supported.

2.3. OSPFv2 and OSPFv3

The data model defined herein supports both OSPFv2 and OSPFv3.

The field 'version' is used to indicate the OSPF version and is mandatory. Based on the configured version, the data model varies to accommodate the differences between OSPFv2 and OSPFv3.

2.4. Optional Features

Optional features are beyond the basic OSPF configuration and it is the responsibility of each vendor to decide whether to support a given feature on a particular device.

This model defines the following optional features:

1. multi-topology: Support Multi-Topology Routing (MTR) [RFC4915].
2. multi-area-adj: Support OSPF multi-area adjacency [RFC5185].
3. explicit-router-id: Support explicit per-instance Router-ID specification.
4. demand-circuit: Support OSPF demand circuits [RFC1793].
5. mtu-ignore: Support disabling OSPF Database Description packet MTU mismatch checking specified in section 10.6 of [RFC2328].
6. lls: Support OSPF link-local signaling (LLS) [RFC5613].
7. prefix-suppression: Support OSPF prefix advertisement suppression [RFC6860].
8. ttl-security: Support OSPF Time to Live (TTL) security check support [RFC5082].
9. nsr: Support OSPF Non-Stop Routing (NSR). The OSPF NSR feature allows a router with redundant control-plane capability (e.g., dual Route-Processor (RP) cards) to maintain its state and adjacencies during planned and unplanned control-plane processing restarts. It differs from graceful-restart or Non-

Stop Forwarding (NSF) in that no protocol signaling or assistance from adjacent OSPF neighbors is required to recover control-plane state.

10. graceful-restart: Support Graceful OSPF Restart [RFC3623], [RFC5187].
11. auto-cost: Support OSPF interface cost calculation according to reference bandwidth [RFC2328].
12. max-ecmp: Support configuration of the maximum number of Equal-Cost Multi-Path (ECMP) paths.
13. max-lsa: Support configuration of the maximum number of LSAs the OSPF instance will accept [RFC1765].
14. te-rid: Support configuration of the Traffic Engineering (TE) Router-ID, i.e., the Router Address described in Section 2.4.1 of [RFC3630] or the Router IPv6 Address TLV described in Section 3 of [RFC5329].
15. ldp-igp-sync: Support LDP IGP synchronization [RFC5443].
16. ospfv2-authentication-trailer: Support OSPFv2 Authentication trailer as specified in [RFC5709] or [RFC7474].
17. ospfv3-authentication-ipsec: Support IPsec for OSPFv3 authentication [RFC4552].
18. ospfv3-authentication-trailer: Support OSPFv3 Authentication trailer as specified in [RFC7166].
19. fast-reroute: Support IP Fast Reroute (IP-FRR) [RFC5714].
20. node-flag: Support node-flag for OSPF prefixes. [RFC7684].
21. node-tag: Support node admin tag for OSPF instances [RFC7777].
22. lfa: Support Loop-Free Alternates (LFAs) [RFC5286].
23. remote-lfa: Support Remote Loop-Free Alternates (R-LFA) [RFC7490].
24. stub-router: Support RFC 6987 OSPF Stub Router advertisement [RFC6987].
25. pe-ce-protocol: Support OSPF as a PE-CE protocol [RFC4577], [RFC6565].

- 26. ietf-spf-delay: Support IETF SPF delay algorithm [RFC8405].
- 27. bfd: Support BFD detection of OSPF neighbor reachability [RFC5880], [RFC5881], and [I-D.ietf-bfd-yang].
- 28. hybrid-interface: Support OSPF Hybrid Broadcast and Point-to-Point Interfaces [RFC6845].

It is expected that vendors will support additional features through vendor-specific augmentations.

2.5. OSPF Router Configuration/Operational State

The ospf container is the top-level container in this data model. It represents an OSPF protocol instance and contains the router level configuration and operational state. The operational state includes the instance statistics, IETF SPF delay statistics, AS-Scoped Link State Database, local RIB, SPF Log, and the LSA log.

```

module: ietf-ospf
  augment /rt:routing/rt:control-plane-protocols/
    rt:control-plane-protocol:
      +--rw ospf
      .
      .
      +--rw af iana-rt-types:address-family
      +--rw enable? boolean
      +--rw explicit-router-id? rt-types:router-id
      | {explicit-router-id}?
      +--rw preference
      | +--rw (scope)?
      | | +--:(single-value)
      | | | +--rw all? uint8
      | | +--:(multi-values)
      | | | +--rw (granularity)?
      | | | | +--:(detail)
      | | | | | +--rw intra-area? uint8
      | | | | | +--rw inter-area? uint8
      | | | | +--:(coarse)
      | | | | | +--rw internal? uint8
      | | | +--rw external? uint8
      +--rw nsr {nsr}?
      | +--rw enable? boolean
      +--rw graceful-restart {graceful-restart}?
      | +--rw enable? boolean
      +--rw helper-enable? boolean
      +--rw restart-interval? uint16
      +--rw helper-strict-lsa-checking? boolean
  
```

```

+--rw auto-cost {auto-cost}?
|   +--rw enable?                boolean
|   +--rw reference-bandwidth?   uint32
+--rw spf-control
|   +--rw paths?                  uint16 {max-ecmp}?
|   +--rw ietf-spf-delay {ietf-spf-delay}?
|       +--rw initial-delay?     uint16
|       +--rw short-delay?       uint16
|       +--rw long-delay?        uint16
|       +--rw hold-down?         uint16
|       +--rw time-to-learn?     uint16
|       +--ro current-state?     enumeration
|       +--ro remaining-time-to-learn? uint16
|       +--ro remaining-hold-down? uint16
|       +--ro last-event-received? yang:timestamp
|       +--ro next-spf-time?     yang:timestamp
|       +--ro last-spf-time?     yang:timestamp
+--rw database-control
|   +--rw max-lsa?               uint32 {max-lsa}?
+--rw stub-router {stub-router}?
|   +--rw (trigger)?
|       +--:(always)
|       +--rw always!
+--rw mpls
|   +--rw te-rid {te-rid}?
|       +--rw ipv4-router-id?    inet:ipv4-address
|       +--rw ipv6-router-id?    inet:ipv6-address
|   +--rw ldp
|       +--rw igp-sync?          boolean {ldp-igp-sync}?
+--rw fast-reroute {fast-reroute}?
|   +--rw lfa {lfa}?
+--ro protected-routes
|   +--ro af-stats* [af prefix alternate]
|       +--ro af                iana-rt-types:address-family
|       +--ro prefix              string
|       +--ro alternate          string
|       +--ro alternate-type?    enumeration
|       +--ro best?              boolean
|       +--ro non-best-reason?   string
|       +--ro protection-available? bits
|       +--ro alternate-metric1? uint32
|       +--ro alternate-metric2? uint32
|       +--ro alternate-metric3? uint32
+--ro unprotected-routes
|   +--ro af-stats* [af prefix]
|       +--ro af                iana-rt-types:address-family
|       +--ro prefix              string
+--ro protection-statistics* [frr-protection-method]

```

```

    +--ro frr-protection-method string
    +--ro af-stats* [af]
      +--ro af iana-rt-types:address-family
      +--ro total-routes? uint32
      +--ro unprotected-routes? uint32
      +--ro protected-routes? uint32
      +--ro linkprotected-routes? uint32
      +--ro nodeprotected-routes? uint32
+--rw node-tags {node-tag}?
  +--rw node-tag* [tag]
    +--rw tag uint32
+--ro router-id?
+--ro local-rib
  +--ro route* [prefix]
    +--ro prefix inet:ip-prefix
    +--ro next-hops
      +--ro next-hop* [next-hop]
        +--ro outgoing-interface? if:interface-ref
        +--ro next-hop inet:ip-address
    +--ro metric? uint32
    +--ro route-type? route-type
    +--ro route-tag? uint32
+--ro statistics
  +--ro discontinuity-time yang:date-and-time
  +--ro originate-new-lsa-count? yang:counter32
  +--ro rx-new-lsas-count? yang:counter32
  +--ro as-scope-lsa-count? yang:gauge32
  +--ro as-scope-lsa-chksum-sum? uint32
  +--ro database
    +--ro as-scope-lsa-type*
      +--ro lsa-type? uint16
      +--ro lsa-count? yang:gauge32
      +--ro lsa-cksum-sum? int32
+--ro database
  +--ro as-scope-lsa-type* [lsa-type]
    +--ro as-scope-lsas
      +--ro as-scope-lsa* [lsa-id adv-router]
        +--ro lsa-id union
        +--ro adv-router inet:ipv4-address
        +--ro decoded-completed? boolean
        +--ro raw-data? yang:hex-string
        +--ro (version)?
          +--:(ospfv2)
            +--ro ospfv2
          .
          .
          +--:(ospfv3)
            +--ro ospfv3

```



```

.
.
+--ro spf-log
|   +--ro event* [id]
|   |   +--ro id                uint32
|   |   +--ro spf-type?         enumeration
|   |   +--ro schedule-timestamp? yang:timestamp
|   |   +--ro start-timestamp?  yang:timestamp
|   |   +--ro end-timestamp?    yang:timestamp
|   |   +--ro trigger-lsa*
|   |   |   +--ro area-id?      area-id-type
|   |   |   +--ro link-id?     union
|   |   |   +--ro type?        uint16
|   |   |   +--ro lsa-id?      yang:dotted-quad
|   |   |   +--ro adv-router?  yang:dotted-quad
|   |   |   +--ro seq-num?     uint32
|   +--ro lsa-log
|   |   +--ro event* [id]
|   |   |   +--ro id                uint32
|   |   |   +--ro lsa
|   |   |   |   +--ro area-id?      area-id-type
|   |   |   |   +--ro link-id?     union
|   |   |   |   +--ro type?        uint16
|   |   |   |   +--ro lsa-id?      yang:dotted-quad
|   |   |   |   +--ro adv-router?  yang:dotted-quad
|   |   |   |   +--ro seq-num?     uint32
|   |   +--ro received-timestamp? yang:timestamp
|   +--ro reason?                  identityref
.
.

```

2.6. OSPF Area Configuration/Operational State

The area container contains OSPF area configuration and the list of interface containers representing all the OSPF interfaces in the area. The area operational state includes the area statistics and the Area Link State Database (LSDB).

```

module: ietf-ospf
  augment /rt:routing/rt:control-plane-protocols/
    rt:control-plane-protocol:
      +--rw ospf
      .
      .
      +--rw areas
      |   +--rw area* [area-id]
      |   |   +--rw area-id                area-id-type
      |   |   +--rw area-type?             identityref

```

```

+--rw summary?                               boolean
+--rw default-cost?                           uint32
+--rw ranges
|   +--rw range* [prefix]
|   |   +--rw prefix          inet:ip-prefix
|   |   +--rw advertise?      boolean
|   |   +--rw cost?           uint24
+--rw topologies {ospf:multi-topology}?
|   +--rw topology* [name]
|   |   +--rw name -> ../../../../rt:ribs/rib/name
|   |   +--rw summary?        boolean
|   |   +--rw default-cost?    ospf-metric
|   |   +--rw ranges
|   |   |   +--rw range* [prefix]
|   |   |   |   +--rw prefix          inet:ip-prefix
|   |   |   |   +--rw advertise?      boolean
|   |   |   |   +--rw cost?           ospf-metric
+--ro statistics
|   +--ro discontinuity-time                yang:date-and-time
|   +--ro spf-runs-count?                   yang:counter32
|   +--ro abr-count?                       yang:gauge32
|   +--ro asbr-count?                      yang:gauge32
|   +--ro ar-nssa-translator-event-count?
|   |   +--ro area-scope-lsa-count?         yang:counter32
|   |   +--ro area-scope-lsa-cksum-sum?     int32
+--ro database
|   +--ro area-scope-lsa-type*
|   |   +--ro lsa-type?                     uint16
|   |   +--ro lsa-count?                   yang:gauge32
|   |   +--ro lsa-cksum-sum?               int32
+--ro database
|   +--ro area-scope-lsa-type* [lsa-type]
|   |   +--ro lsa-type                     uint16
|   +--ro area-scope-lsas
|   |   +--ro area-scope-lsa* [lsa-id adv-router]
|   |   |   +--ro lsa-id                     union
|   |   .
|   |   .
|   |   +--ro (version)?
|   |   |   +--:(ospfv2)
|   |   |   |   +--ro ospfv2
|   |   |   |   +--ro header
|   |   .
|   |   .
|   |   +--ro body
|   |   |   +--ro router

```

```

.      .      .      .
|      |      |      +--ro network
.      .      .      .
|      |      |      +--ro summary
.      .      .      .
|      |      |      +--ro external
.      .      .      .
|      |      |      +--ro opaque
.      .      .      .
|      |      |      +--:(ospfv3)
|      |      |      +--ro ospfv3
|      |      |      +--ro header
.      .      .      .
|      |      |      +--ro body
|      |      |      +--ro router
.      .      .      .
|      |      |      +--ro network
.      .      .      .
|      |      |      +--ro inter-area-prefix
.      .      .      .
|      |      |      +--ro inter-area-router
.      .      .      .
|      |      |      +--ro as-external
.      .      .      .
|      |      |      +--ro nssa
.      .      .      .
|      |      |      +--ro link
.      .      .      .
|      |      |      +--ro intra-area-prefix
.      .      .      .
|      |      |      +--ro router-information
.      .      .      .
|      |      |      +--rw virtual-links

```

```

+--rw virtual-link* [transit-area-id router-id]
  +--rw transit-area-id      -> ../../../../
                               area/area-id
  +--rw router-id            rt-types:router-id
  +--rw hello-interval?      uint16
  +--rw dead-interval?       uint32
  +--rw retransmit-interval? uint16
  +--rw transmit-delay?      uint16
  +--rw lls?                  boolean {lls}?
  +--rw ttl-security {ttl-security}?
    | +--rw enable?          boolean
    | +--rw hops?            uint8
  +--rw enable?              boolean
  +--rw authentication
    +--rw (auth-type-selection)?
      +--:(ospfv2-auth)
        | +--rw ospfv2-auth-trailer-rfc?
        | | ospfv2-auth-trailer-rfc-version
        | | {ospfv2-authentication-trailer}?
        +--rw (ospfv2-auth-specification)?
          +--:(auth-key-chain) {key-chain}?
            | +--rw ospfv2-key-chain?
            | | key-chain:key-chain-ref
            +--:(auth-key-explicit)
              +--rw ospfv2-key-id?      uint32
              +--rw ospfv2-key?        string
              +--rw ospfv2-crypto-algorithm?
              | identityref
          +--:(ospfv3-auth-ipsec)
            | {ospfv3-authentication-ipsec}?
            +--rw sa?                  string
          +--:(ospfv3-auth-trailer)
            | {ospfv3-authentication-trailer}?
            +--rw (ospfv3-auth-specification)?
              +--:(auth-key-chain) {key-chain}?
                | +--rw ospfv3-key-chain?
                | | key-chain:key-chain-ref
                +--:(auth-key-explicit)
                  +--rw ospfv3-sa-id?      uint16
                  +--rw ospfv3-key?        string
                  +--rw ospfv3-crypto-algorithm?
                  | identityref
      +--ro cost?                    uint16
      +--ro state?                   if-state-type
      +--ro hello-timer?             rt-types:
      |                               rtimer-value-seconds16
      +--ro wait-timer?              rt-types:
      |                               rtimer-value-seconds16

```

```

+--ro dr-router-id?          rt-types:router-id
+--ro dr-ip-addr?            inet:ip-address
+--ro bdr-router-id?         rt-types:router-id
+--ro bdr-ip-addr?           inet:ip-address
+--ro statistics
|   +--ro discontinuity-time    yang:date-and-time
|   +--ro if-event-count?      yang:counter32
|   +--ro link-scope-lsa-count? yang:gauge32
|   +--ro link-scope-lsa-cksum-sum?
|                                   uint32
|   +--ro database
|       +--ro link-scope-lsa-type*
|           +--ro lsa-type?      uint16
|           +--ro lsa-count?     yang:gauge32
|           +--ro lsa-cksum-sum? int32
+--ro neighbors
|   +--ro neighbor* [neighbor-router-id]
|       +--ro neighbor-router-id
|                                   rt-types:router-id
|       +--ro address?          inet:ip-address
|       +--ro dr-router-id?     rt-types:router-id
|       +--ro dr-ip-addr?       inet:ip-address
|       +--ro bdr-router-id?    rt-types:router-id
|       +--ro bdr-ip-addr?      inet:ip-address
|       +--ro state?            nbr-state-type
|       +--ro dead-timer? rt-types:
|           |                   rtimer-value-seconds16
|       +--ro statistics
|           +--ro discontinuity-time
|                                   yang:date-and-time
|           +--ro nbr-event-count?
|                                   yang:counter32
|           +--ro nbr-retrans-qlen?
|                                   yang:gauge32
+--ro database
|   +--ro link-scope-lsa-type* [lsa-type]
|       +--ro lsa-type          uint16
|       +--ro link-scope-lsas
|
+--rw sham-links {pe-ce-protocol}?
|   +--rw sham-link* [local-id remote-id]
|       +--rw local-id          inet:ip-address
|       +--rw remote-id         inet:ip-address
|       +--rw hello-interval?   uint16
|       +--rw dead-interval?    uint32
|       +--rw retransmit-interval?
|                               uint16
|       +--rw transmit-delay?   uint16

```

```

+--rw lls?                               boolean {lls}?
+--rw ttl-security {ttl-security}?
|   +--rw enable?    boolean
|   +--rw hops?      uint8
+--rw enable?                boolean
+--rw authentication
|   +--rw (auth-type-selection)?
|   |   +--:(ospfv2-auth)
|   |   |   +--rw ospfv2-auth-trailer-rfc?
|   |   |   |   ospfv2-auth-trailer-rfc-version
|   |   |   |   {ospfv2-authentication-trailer}?
|   |   +--rw (ospfv2-auth-specification)?
|   |   |   +--:(auth-key-chain) {key-chain}?
|   |   |   |   +--rw ospfv2-key-chain?
|   |   |   |   |   key-chain:key-chain-ref
|   |   |   +--:(auth-key-explicit)
|   |   |   |   +--rw ospfv2-key-id?      uint32
|   |   |   |   +--rw ospfv2-key?        string
|   |   |   |   +--rw ospfv2-crypto-algorithm?
|   |   |   |   |   identityref
|   |   +--:(ospfv3-auth-ipsec)
|   |   |   {ospfv3-authentication-ipsec}?
|   |   |   +--rw sa?                      string
|   |   +--:(ospfv3-auth-trailer)
|   |   |   {ospfv3-authentication-trailer}?
|   |   +--rw (ospfv3-auth-specification)?
|   |   |   +--:(auth-key-chain) {key-chain}?
|   |   |   |   +--rw ospfv3-key-chain?
|   |   |   |   |   key-chain:key-chain-ref
|   |   |   +--:(auth-key-explicit)
|   |   |   |   +--rw ospfv3-sa-id?        uint16
|   |   |   |   +--rw ospfv3-key?          string
|   |   |   |   +--rw ospfv3-crypto-algorithm?
|   |   |   |   |   identityref
+--rw cost?                      uint16
+--rw mtu-ignore?                boolean
|   {mtu-ignore}?
+--rw prefix-suppression?        boolean
|   {prefix-suppression}?
+--ro state?                     if-state-type
+--ro hello-timer?              rt-types:
|   rtimer-value-seconds16
+--ro wait-timer?               rt-types:
|   rtimer-value-seconds16
+--ro dr-router-id?             rt-types:router-id
+--ro dr-ip-addr?               inet:ip-address
+--ro bdr-router-id?            rt-types:router-id
+--ro bdr-ip-addr?              inet:ip-address

```

```

+--ro statistics
  +--ro discontinuity-time      yang:date-and-time
  +--ro if-event-count?        yang:counter32
  +--ro link-scope-lsa-count?   yang:gauge32
  +--ro link-scope-lsa-cksum-sum?
                                uint32
  +--ro database
    +--ro link-scope-lsa-type*
      +--ro lsa-type?           uint16
      +--ro lsa-count?          yang:gauge32
      +--ro lsa-cksum-sum?      int32
+--ro neighbors
  +--ro neighbor* [neighbor-router-id]
    +--ro neighbor-router-id
                                rt-types:router-id
    +--ro address?              inet:ip-address
    +--ro dr-router-id?         rt-types:router-id
    +--ro dr-ip-addr?           inet:ip-address
    +--ro bdr-router-id?       rt-types:router-id
    +--ro bdr-ip-addr?         inet:ip-address
    +--ro state?                nbr-state-type
    +--ro cost?                 uint32
    +--ro dead-timer? rt-types:
      |                         rtimer-value-seconds16
    +--ro statistics
      +--ro nbr-event-count?
                                yang:counter32
      +--ro nbr-retrans-qlen?
                                yang:gauge32
+--ro database
  +--ro link-scope-lsa-type* [lsa-type]
    +--ro lsa-type              uint16
    +--ro link-scope-lsas

```

2.7. OSPF Interface Configuration/Operational State

The interface container contains OSPF interface configuration and operational state. The interface operational state includes the statistics, list of neighbors, and Link-Local Link State Database (LSDB).

```

module: ietf-ospf
  augment /rt:routing/rt:control-plane-protocols/
    rt:control-plane-protocol:
      +--rw ospf
      .

```

```

.
+--rw areas
|   +--rw area* [area-id]
|   |   .
|   |   .
|   |   +--rw interfaces
|   |   |   +--rw interface* [name]
|   |   |   |   +--rw name                if:interface-ref
|   |   |   |   +--rw interface-type?     enumeration
|   |   |   |   +--rw passive?            boolean
|   |   |   |   +--rw demand-circuit?     boolean
|   |   |   |   |   {demand-circuit}?
|   |   |   |   +--rw priority?           uint8
|   |   |   +--rw multi-areas {multi-area-adj}?
|   |   |   |   +--rw multi-area* [multi-area-id]
|   |   |   |   |   +--rw multi-area-id     area-id-type
|   |   |   |   |   +--rw cost?            uint16
|   |   |   +--rw static-neighbors
|   |   |   |   +--rw neighbor* [identifier]
|   |   |   |   |   +--rw identifier        inet:ip-address
|   |   |   |   |   +--rw cost?            uint16
|   |   |   |   |   +--rw poll-interval?   uint16
|   |   |   |   |   +--rw priority?        uint8
|   |   |   +--rw node-flag?              boolean
|   |   |   |   {node-flag}?
|   |   +--rw bfd {bfd}?
|   |   |   +--rw enable?    boolean
|   |   +--rw fast-reroute {fast-reroute}?
|   |   |   +--rw lfa {lfa}?
|   |   |   |   +--rw candidate-enable?    boolean
|   |   |   |   +--rw enable?              boolean
|   |   |   |   +--rw remote-lfa {remote-lfa}?
|   |   |   |   |   +--rw enable?    boolean
|   |   +--rw hello-interval?    uint16
|   |   +--rw dead-interval?     uint32
|   |   +--rw retransmit-interval? uint16
|   |   +--rw transmit-delay?     uint16
|   |   +--rw lls?                boolean {lls}?
|   |   +--rw ttl-security {ttl-security}?
|   |   |   +--rw enable?    boolean
|   |   |   +--rw hops?      uint8
|   |   +--rw enable?        boolean
|   |   +--rw authentication
|   |   |   +--rw (auth-type-selection)?
|   |   |   |   +--:(ospfv2-auth)
|   |   |   |   |   +--rw ospfv2-auth-trailer-rfc?
|   |   |   |   |   |   ospfv2-auth-trailer-rfc-version
|   |   |   |   |   |   {ospfv2-authentication-trailer}?

```



```

    +--rw (ospfv2-auth-specification)?
      +--:(auth-key-chain) {key-chain}?
        |   +--rw ospfv2-key-chain?
        |       key-chain:key-chain-ref
      +--:(auth-key-explicit)
        +--rw ospfv2-key-id?      uint32
        +--rw ospfv2-key?        string
        +--rw ospfv2-crypto-algorithm?
            identityref
    +--:(ospfv3-auth-ipsec)
      |   {ospfv3-authentication-ipsec}?
      |   +--rw sa?                string
    +--:(ospfv3-auth-trailer)
      |   {ospfv3-authentication-trailer}?
    +--rw (ospfv3-auth-specification)?
      +--:(auth-key-chain) {key-chain}?
        |   +--rw ospfv3-key-chain?
        |       key-chain:key-chain-ref
      +--:(auth-key-explicit)
        +--rw ospfv3-sa-id?        uint16
        +--rw ospfv3-key?          string
        +--rw ospfv3-crypto-algorithm?
            identityref
    +--rw cost?                    uint16
    +--rw mtu-ignore?              boolean
    |   {mtu-ignore}?
    +--rw prefix-suppression?     boolean
    |   {prefix-suppression}?
    +--ro state?                  if-state-type
    +--ro hello-timer?            rt-types:
    |   rtimer-value-seconds16
    +--ro wait-timer?            rt-types:
    |   rtimer-value-seconds16
    +--ro dr-router-id?           rt-types:router-id
    +--ro dr-ip-addr?             inet:ip-address
    +--ro bdr-router-id?          rt-types:router-id
    +--ro bdr-ip-addr?           inet:ip-address
    +--ro statistics
      +--ro if-event-count?        yang:counter32
      +--ro link-scope-lsa-count?  yang:gauge32
      +--ro link-scope-lsa-cksum-sum?
          uint32
      +--ro database
        +--ro link-scope-lsa-type*
          +--ro lsa-type?          uint16
          +--ro lsa-count?         yang:gauge32
          +--ro lsa-cksum-sum?    int32
    +--ro neighbors

```

```

|
|
|      +---ro neighbor* [neighbor-router-id]
|      |      +---ro neighbor-router-id
|      |      |
|      |      |      rt-types:router-id
|      |      +---ro address?      inet:ip-address
|      |      +---ro dr-router-id?  rt-types:router-id
|      |      +---ro dr-ip-addr?    inet:ip-address
|      |      +---ro bdr-router-id? rt-types:router-id
|      |      +---ro bdr-ip-addr?   inet:ip-address
|      |      +---ro state?         nbr-state-type
|      |      +---ro dead-timer?    rt-types:
|      |      |      rtimer-value-seconds16
|      |      +---ro statistics
|      |      |      +---ro nbr-event-count?
|      |      |      |      yang:counter32
|      |      |      +---ro nbr-retrans-qlen?
|      |      |      |      yang:gauge32
|      +---ro database
|      .   +---ro link-scope-lsa-type* [lsa-type]
|      .   +---ro lsa-type      uint16
|      .   +---ro link-scope-lsas
|      .
|      .
|      +---rw topologies {ospf:multi-topology}?
|      |      +---rw topology* [name]
|      |      |      +---rw name -> ../../../../rt:ribs/rib/name
|      |      |      |
|      |      |      +---rw cost? uint32
|      +---rw instance-id?      uint8
|
|
|

```

2.8. OSPF Notifications

This YANG model defines a list of notifications that inform YANG clients of important events detected during protocol operation. The defined notifications cover the common set of traps from the OSPFv2 MIB [RFC4750] and OSPFv3 MIB [RFC5643].

```

notifications:
  +---n if-state-change
  |   +---ro routing-protocol-name?
  |   +   -> /rt:routing/control-plane-protocols/
  |   +   control-plane-protocol/name
  |   +---ro af?
  |   +   -> /rt:routing/control-plane-protocols/
  |   +   control-plane-protocol
  |   +   [rt:name=current()/../routing-protocol-name]/
  |   +   ospf:ospf/af

```

```

+--ro (if-link-type-selection)?
+--:(interface)
+--ro interface
+--ro interface?   if:interface-ref
+--:(virtual-link)
+--ro virtual-link
+--ro transit-area-id?   area-id-type
+--ro neighbor-router-id? rt-types:router-id
+--:(sham-link)
+--ro sham-link
+--ro area-id?   area-id-type
+--ro local-ip-addr?   inet:ip-address
+--ro remote-ip-addr?   inet:ip-address
+--ro state?   if-state-type
+---n if-config-error
+--ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol/name
+--ro af?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol
+       [rt:name=current()/../routing-protocol-name]/
+       ospf:ospf/af
+--ro (if-link-type-selection)?
+--:(interface)
+--ro interface
+--ro interface?   if:interface-ref
+--:(virtual-link)
+--ro virtual-link
+--ro transit-area-id?   area-id-type
+--ro neighbor-router-id? rt-types:router-id
+--:(sham-link)
+--ro sham-link
+--ro area-id?   area-id-type
+--ro local-ip-addr?   inet:ip-address
+--ro remote-ip-addr?   inet:ip-address
+--ro packet-source?   yang:dotted-quad
+--ro packet-type?   packet-type
+--ro error?   enumeration
+---n nbr-state-change
+--ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol/name
+--ro af?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol
+       [rt:name=current()/../routing-protocol-name]/
+       ospf:ospf/af

```

```

+---ro (if-link-type-selection)?
+---:(interface)
+---ro interface
+---ro interface?    if:interface-ref
+---:(virtual-link)
+---ro virtual-link
+---ro transit-area-id?    area-id-type
+---ro neighbor-router-id?  rt-types:router-id
+---:(sham-link)
+---ro sham-link
+---ro area-id?    area-id-type
+---ro local-ip-addr?    inet:ip-address
+---ro remote-ip-addr?    inet:ip-address
+---ro neighbor-router-id?    rt-types:router-id
+---ro neighbor-ip-addr?    yang:dotted-quad
+---ro state?    nbr-state-type
+---n nbr-restart-helper-status-change
+---ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol/name
+---ro af?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol
+       [rt:name=current()/../routing-protocol-name]/
+       ospf:ospf/af
+---ro (if-link-type-selection)?
+---:(interface)
+---ro interface
+---ro interface?    if:interface-ref
+---:(virtual-link)
+---ro virtual-link
+---ro transit-area-id?    area-id-type
+---ro neighbor-router-id?  rt-types:router-id
+---:(sham-link)
+---ro sham-link
+---ro area-id?    area-id-type
+---ro local-ip-addr?    inet:ip-address
+---ro remote-ip-addr?    inet:ip-address
+---ro neighbor-router-id?    rt-types:router-id
+---ro neighbor-ip-addr?    yang:dotted-quad
+---ro status?    restart-helper-status-type
+---ro age?    uint32
+---ro exit-reason?    restart-exit-reason-type
+---n if-rx-bad-packet
+---ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+       control-plane-protocol/name
+---ro af?

```

```

+      -> /rt:routing/control-plane-protocols/
+      control-plane-protocol
+      [rt:name=current()/../routing-protocol-name]/
+      ospf:ospf/af
+---ro (if-link-type-selection)?
+   +---:(interface)
+   |   +---ro interface
+   |   |   +---ro interface?    if:interface-ref
+   |   +---:(virtual-link)
+   |   |   +---ro virtual-link
+   |   |   |   +---ro transit-area-id?    area-id-type
+   |   |   |   +---ro neighbor-router-id? rt-types:router-id
+   |   +---:(sham-link)
+   |   |   +---ro sham-link
+   |   |   |   +---ro area-id?    area-id-type
+   |   |   |   +---ro local-ip-addr?    inet:ip-address
+   |   |   |   +---ro remote-ip-addr?    inet:ip-address
+   +---ro packet-source?    yang:dotted-quad
+   +---ro packet-type?    packet-type
+---n lsdb-approaching-overflow
+---ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol/name
+---ro af?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol
+   [rt:name=current()/../routing-protocol-name]/
+   ospf:ospf/af
+---ro ext-lsdb-limit?    uint32
+---n lsdb-overflow
+---ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol/name
+---ro af?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol
+   [rt:name=current()/../routing-protocol-name]/
+   ospf:ospf/af
+---ro ext-lsdb-limit?    uint32
+---n nssa-translator-status-change
+---ro routing-protocol-name?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol/name
+---ro af?
+   -> /rt:routing/control-plane-protocols/
+   control-plane-protocol
+   [rt:name=current()/../routing-protocol-name]/
+   ospf:ospf/af

```

```

|   +---ro area-id?                area-id-type
|   +---ro status?                nssa-translator-state-type
+---n restart-status-change
|   +---ro routing-protocol-name?
|   +       -> /rt:routing/control-plane-protocols/
|   +       control-plane-protocol/name
+---ro af?
|   +       -> /rt:routing/control-plane-protocols/
|   +       control-plane-protocol
|   +       [rt:name=current()/../routing-protocol-name]/
|   +       ospf:ospf/af
+---ro status?                    restart-status-type
+---ro restart-interval?          uint16
+---ro exit-reason?              restart-exit-reason-type

```

2.9. OSPF RPC Operations

The "ietf-ospf" module defines two RPC operations:

- o clear-database: reset the content of a particular OSPF Link State Database.
- o clear-neighbor: Reset a particular OSPF neighbor or group of neighbors associated with an OSPF interface.

```

rpcs:
+---x clear-neighbor
|   +---w input
|   |   +---w routing-protocol-name
|   |   +       -> /rt:routing/control-plane-protocols/
|   |   +       control-plane-protocol/name
|   |   +---w interface?            if:interface-ref
+---x clear-database
|   +---w input
|   |   +---w routing-protocol-name
|   |   |       -> /rt:routing/control-plane-protocols/
|   |   |       control-plane-protocol/name

```

3. OSPF YANG Module

The following RFCs and drafts are not referenced in the document text but are referenced in the ietf-ospf.yang module: [RFC0905], [RFC4576], [RFC4973], [RFC5250], [RFC5309], [RFC5642], [RFC5881], [RFC6991], [RFC7770], [RFC7884], [RFC8294], and [RFC8476].

```

<CODE BEGINS> file "ietf-ospf@2019-10-17.yang"
module ietf-ospf {
  yang-version 1.1;

```

```
namespace "urn:ietf:params:xml:ns:yang:ietf-ospf";

prefix ospf;

import ietf-inet-types {
  prefix "inet";
  reference "RFC 6991: Common YANG Data Types";
}

import ietf-yang-types {
  prefix "yang";
  reference "RFC 6991: Common YANG Data Types";
}

import ietf-interfaces {
  prefix "if";
  reference "RFC 8343: A YANG Data Model for Interface
            Management (NMDA Version)";
}

import ietf-routing-types {
  prefix "rt-types";
  reference "RFC 8294: Common YANG Data Types for the
            Routing Area";
}

import iana-routing-types {
  prefix "iana-rt-types";
  reference "RFC 8294: Common YANG Data Types for the
            Routing Area";
}

import ietf-routing {
  prefix "rt";
  reference "RFC 8349: A YANG Data Model for Routing
            Management (NMDA Version)";
}

import ietf-key-chain {
  prefix "key-chain";
  reference "RFC 8177: YANG Data Model for Key Chains";
}

import ietf-bfd-types {
  prefix "bfd-types";
  reference "RFC YYYY: YANG Data Model for Bidirectional
            Forwarding Detection (BFD). Please replace YYYY with
            published RFC number for draft-ietf-bfd-yang.";
```

```
}

organization
  "IETF LSR - Link State Routing Working Group";

contact
  "WG Web:  <https://datatracker.ietf.org/group/lsr/>
  WG List:  <mailto:lsr@ietf.org>

  Editor:    Derek Yeung
             <mailto:derek@arrcus.com>
  Author:    Acee Lindem
             <mailto:acee@cisco.com>
  Author:    Yingzhen Qu
             <mailto:yingzhen.qu@futurewei.com>
  Author:    Salih K A
             <mailto:salih@juniper.net>
  Author:    Ing-Wher Chen
             <mailto:ingwherchen@mitre.org>;

description
  "This YANG module defines the generic configuration and
  operational state for the OSPF protocol common to all
  vendor implementations. It is intended that the module
  will be extended by vendors to define vendor-specific
  OSPF configuration parameters and policies,
  for example, route maps or route policies.

  This YANG model conforms to the Network Management
  Datastore Architecture (NMDA) as described in RFC 8242.

  Copyright (c) 2018 IETF Trust and the persons identified as
  authors of the code. All rights reserved.

  Redistribution and use in source and binary forms, with or
  without modification, is permitted pursuant to, and subject to
  the license terms contained in, the Simplified BSD License set
  forth in Section 4.c of the IETF Trust's Legal Provisions
  Relating to IETF Documents
  (https://trustee.ietf.org/license-info).

  This version of this YANG module is part of RFC XXXX
  (https://www.rfc-editor.org/info/rfcXXXX); see the RFC itself
  for full legal notices.

  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 (RFC 2119) (RFC 8174) when, and only when, they appear in all capitals, as shown here.

This version of this YANG module is part of RFC XXXX; see the RFC itself for full legal notices.";

```
revision 2019-10-17 {
  description
    "Initial revision.";
  reference
    "RFC XXXX: A YANG Data Model for OSPF.";
}

feature multi-topology {
  description
    "Support Multiple-Topology Routing (MTR).";
  reference "RFC 4915: Multi-Topology Routing";
}

feature multi-area-adj {
  description
    "OSPF multi-area adjacency support as in RFC 5185.";
  reference "RFC 5185: Multi-Area Adjacency";
}

feature explicit-router-id {
  description
    "Set Router-ID per instance explicitly.";
}

feature demand-circuit {
  description
    "OSPF demand circuit support as in RFC 1793.";
  reference "RFC 1793: OSPF Demand Circuits";
}

feature mtu-ignore {
  description
    "Disable OSPF Database Description packet MTU
     mismatch checking specified in the OSPF
     protocol specification.";
  reference "RFC 2328: OSPF Version 2, section 10.6";
}

feature lls {
  description
    "OSPF link-local signaling (LLS) as in RFC 5613.";
  reference "RFC 5613: OSPF Link-Local Signaling";
}
```

```
feature prefix-suppression {
  description
    "OSPF prefix suppression support as in RFC 6860.";
  reference "RFC 6860: Hide Transit-Only Networks in OSPF";
}

feature ttl-security {
  description
    "OSPF Time to Live (TTL) security check support.";
  reference "RFC 5082: The Generalized TTL Security
            Mechanism (GTSM)";
}

feature nsr {
  description
    "Non-Stop-Routing (NSR) support. The OSPF NSR feature
     allows a router with redundant control-plane capability
     (e.g., dual Route-Processor (RP) cards) to maintain its
     state and adjacencies during planned and unplanned
     OSPF instance restarts. It differs from graceful-restart
     or Non-Stop Forwarding (NSF) in that no protocol signaling
     or assistance from adjacent OSPF neighbors is required to
     recover control-plane state.";
}

feature graceful-restart {
  description
    "Graceful OSPF Restart as defined in RFC 3623 and
     RFC 5187.";
  reference "RFC 3623: Graceful OSPF Restart
            RFC 5187: OSPFv3 Graceful Restart";
}

feature auto-cost {
  description
    "Calculate OSPF interface cost according to
     reference bandwidth.";
  reference "RFC 2328: OSPF Version 2";
}

feature max-ecmp {
  description
    "Setting maximum number of ECMP paths.";
}

feature max-lsa {
  description
    "Setting the maximum number of LSAs the OSPF instance
```

```
        will accept.";
        reference "RFC 1765: OSPF Database Overload";
    }

    feature te-rid {
        description
            "Support configuration of the Traffic Engineering (TE)
            Router-ID, i.e., the Router Address described in Section
            2.4.1 of RFC3630 or the Router IPv6 Address TLV described
            in Section 3 of RFC5329.";
        reference "RFC 3630: Traffic Engineering (TE) Extensions
            to OSPF Version 2
            RFC 5329: Traffic Engineering (TE) Extensions
            to OSPF Version 3";
    }

    feature ldp-igp-sync {
        description
            "LDP IGP synchronization.";
        reference "RFC 5443: LDP IGP Synchronization";
    }

    feature ospfv2-authentication-trailer {
        description
            "Support OSPFv2 authentication trailer for OSPFv2
            authentication.";
        reference "RFC 5709: Supporting Authentication
            Trailer for OSPFv2
            RFC 7474: Security Extension for OSPFv2 When
            Using Manual Key Management";
    }

    feature ospfv3-authentication-ipsec {
        description
            "Support IPsec for OSPFv3 authentication.";
        reference "RFC 4552: Authentication/Confidentiality
            for OSPFv3";
    }

    feature ospfv3-authentication-trailer {
        description
            "Support OSPFv3 authentication trailer for OSPFv3
            authentication.";
        reference "RFC 7166: Supporting Authentication
            Trailer for OSPFv3";
    }

    feature fast-reroute {
```

```
    description
      "Support for IP Fast Reroute (IP-FRR).";
    reference "RFC 5714: IP Fast Reroute Framework";
  }

  feature key-chain {
    description
      "Support of keychain for authentication.";
    reference "RFC8177: YANG Data Model for Key Chains";
  }

  feature node-flag {
    description
      "Support for node-flag for OSPF prefixes.";
    reference "RFC 7684: OSPFv2 Prefix/Link Advertisement";
  }

  feature node-tag {
    description
      "Support for node admin tag for OSPF routing instances.";
    reference "RFC 7777: Advertising Node Administrative
              Tags in OSPF";
  }

  feature lfa {
    description
      "Support for Loop-Free Alternates (LFAs).";
    reference "RFC 5286: Basic Specification for IP Fast
              Reroute: Loop-Free Alternates";
  }

  feature remote-lfa {
    description
      "Support for Remote Loop-Free Alternates (R-LFA).";
    reference "RFC 7490: Remote Loop-Free Alternate (LFA)
              Fast Reroute (FRR)";
  }

  feature stub-router {
    description
      "Support for RFC 6987 OSPF Stub Router Advertisement.";
    reference "RFC 6987: OSPF Stub Router Advertisement";
  }

  feature pe-ce-protocol {
    description
      "Support for OSPF as a PE-CE protocol";
    reference "RFC 4577: OSPF as the Provider/Customer Edge
```

```
        Protocol for BGP/MPLS IP Virtual Private
        Networks (VPNs)
        RFC 6565: OSPFv3 as a Provider Edge to Customer
        Edge (PE-CE) Routing Protocol";
    }

    feature ietf-spf-delay {
        description
            "Support for IETF SPF delay algorithm.";
        reference "RFC 8405: SPF Back-off algorithm for link
            state IGP";
    }

    feature bfd {
        description
            "Support for BFD detection of OSPF neighbor reachability.";
        reference "RFC 5880: Bidirectional Forwarding Detection (BFD)
            RFC 5881: Bidirectional Forwarding Detection
            (BFD) for IPv4 and IPv6 (Single Hop)";
    }

    feature hybrid-interface {
        description
            "Support for OSPF Hybrid interface type.";
        reference "RFC 6845: OSPF Hybrid Broadcast and
            Point-to-Multipoint Interface Type";
    }

    identity ospf {
        base "rt:routing-protocol";
        description "Any OSPF protocol version";
    }

    identity ospfv2 {
        base "ospf";
        description "OSPFv2 protocol";
    }

    identity ospfv3 {
        base "ospf";
        description "OSPFv3 protocol";
    }

    identity area-type {
        description "Base identity for OSPF area type.";
    }

    identity normal-area {
```

```
    base area-type;
    description "OSPF normal area.";
}

identity stub-nssa-area {
    base area-type;
    description "OSPF stub or NSSA area.";
}

identity stub-area {
    base stub-nssa-area;
    description "OSPF stub area.";
}

identity nssa-area {
    base stub-nssa-area;
    description "OSPF Not-So-Stubby Area (NSSA).";
    reference "RFC 3101: The OSPF Not-So-Stubby Area
              (NSSA) Option";
}

identity ospf-lsa-type {
    description
        "Base identity for OSPFv2 and OSPFv3
         Link State Advertisement (LSA) types";
}

identity ospfv2-lsa-type {
    base ospf-lsa-type;
    description
        "OSPFv2 LSA types";
}

identity ospfv2-router-lsa {
    base ospfv2-lsa-type;
    description
        "OSPFv2 Router LSA - Type 1";
}

identity ospfv2-network-lsa {
    base ospfv2-lsa-type;
    description
        "OSPFv2 Network LSA - Type 2";
}

identity ospfv2-summary-lsa-type {
    base ospfv2-lsa-type;
    description
```

```
    "OSPFv2 Summary LSA types";
}

identity ospfv2-network-summary-lsa {
    base ospfv2-summary-lsa-type;
    description
        "OSPFv2 Network Summary LSA - Type 3";
}

identity ospfv2-asbr-summary-lsa {
    base ospfv2-summary-lsa-type;
    description
        "OSPFv2 AS Boundary Router (ASBR) Summary LSA - Type 4";
}

identity ospfv2-external-lsa-type {
    base ospfv2-lsa-type;
    description
        "OSPFv2 External LSA types";
}

identity ospfv2-as-external-lsa {
    base ospfv2-external-lsa-type;
    description
        "OSPFv2 AS External LSA - Type 5";
}

identity ospfv2-nssa-lsa {
    base ospfv2-external-lsa-type;
    description
        "OSPFv2 Not-So-Stubby-Area (NSSA) LSA - Type 7";
}

identity ospfv2-opaque-lsa-type {
    base ospfv2-lsa-type;
    description
        "OSPFv2 Opaque LSA types";
}

identity ospfv2-link-scope-opaque-lsa {
    base ospfv2-opaque-lsa-type;
    description
        "OSPFv2 Link-Scoped Opaque LSA - Type 9";
}

identity ospfv2-area-scope-opaque-lsa {
    base ospfv2-opaque-lsa-type;
    description
```

```
        "OSPFv2 Area-Scoped Opaque LSA - Type 10";
    }

    identity ospfv2-as-scope-opaque-lsa {
        base ospfv2-opaque-lsa-type;
        description
            "OSPFv2 AS-Scoped Opaque LSA - Type 11";
    }

    identity ospfv2-unknown-lsa-type {
        base ospfv2-lsa-type;
        description
            "OSPFv2 Unknown LSA type";
    }

    identity ospfv3-lsa-type {
        base ospf-lsa-type;
        description
            "OSPFv3 LSA types.";
    }

    identity ospfv3-router-lsa {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Router LSA - Type 0x2001";
    }

    identity ospfv3-network-lsa {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Network LSA - Type 0x2002";
    }

    identity ospfv3-summary-lsa-type {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Summary LSA types";
    }

    identity ospfv3-inter-area-prefix-lsa {
        base ospfv3-summary-lsa-type;
        description
            "OSPFv3 Inter-area Prefix LSA - Type 0x2003";
    }

    identity ospfv3-inter-area-router-lsa {
        base ospfv3-summary-lsa-type;
        description
```



```
        "OSPFv3 Inter-area Router LSA - Type 0x2004";
    }

    identity ospfv3-external-lsa-type {
        base ospfv3-lsa-type;
        description
            "OSPFv3 External LSA types";
    }

    identity ospfv3-as-external-lsa {
        base ospfv3-external-lsa-type;
        description
            "OSPFv3 AS-External LSA - Type 0x4005";
    }

    identity ospfv3-nssa-lsa {
        base ospfv3-external-lsa-type;
        description
            "OSPFv3 Not-So-Stubby-Area (NSSA) LSA - Type 0x2007";
    }

    identity ospfv3-link-lsa {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Link LSA - Type 0x0008";
    }

    identity ospfv3-intra-area-prefix-lsa {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Intra-area Prefix LSA - Type 0x2009";
    }

    identity ospfv3-router-information-lsa {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Router Information LSA - Types 0x800C,
            0xA00C, and 0xC00C";
    }

    identity ospfv3-unknown-lsa-type {
        base ospfv3-lsa-type;
        description
            "OSPFv3 Unknown LSA type";
    }

    identity lsa-log-reason {
        description
```

```
    "Base identity for an LSA log reason.";
}

identity lsa-refresh {
  base lsa-log-reason;
  description
    "Identity used when the LSA is logged
     as a result of receiving a refresh LSA.";
}

identity lsa-content-change {
  base lsa-log-reason;
  description
    "Identity used when the LSA is logged
     as a result of a change in the content
     of the LSA.";
}

identity lsa-purge {
  base lsa-log-reason;
  description
    "Identity used when the LSA is logged
     as a result of being purged.";
}

identity informational-capability {
  description
    "Base identity for router informational capabilities.";
}

identity graceful-restart {
  base informational-capability;
  description
    "When set, the router is capable of restarting
     gracefully.";
  reference "RFC 3623: Graceful OSPF Restart
            RFC 5187: OSPFv3 Graceful Restart";
}

identity graceful-restart-helper {
  base informational-capability;
  description
    "When set, the router is capable of acting as
     a graceful restart helper.";
  reference "RFC 3623: Graceful OSPF Restart
            RFC 5187: OSPFv3 Graceful Restart";
}
```

```
identity stub-router {
  base informational-capability;
  description
    "When set, the router is capable of acting as
    an OSPF Stub Router.";
  reference "RFC 6987: OSPF Stub Router Advertisement";
}

identity traffic-engineering {
  base informational-capability;
  description
    "When set, the router is capable of OSPF traffic
    engineering.";
  reference "RFC 3630: Traffic Engineering (TE) Extensions
    to OSPF Version 2
    RFC 5329: Traffic Engineering (TE) Extensions
    to OSPF Version 3";
}

identity p2p-over-lan {
  base informational-capability;
  description
    "When set, the router is capable of OSPF Point-to-Point
    over LAN.";
  reference "RFC 5309: Point-to-Point Operation over LAN
    in Link State Routing Protocols";
}

identity experimental-te {
  base informational-capability;
  description
    "When set, the router is capable of OSPF experimental
    traffic engineering.";
  reference
    "RFC 4973: OSPF-xTE OSPF Experimental Traffic
    Engineering";
}

identity router-lsa-bit {
  description
    "Base identity for Router-LSA bits.";
}

identity vlink-end-bit {
  base router-lsa-bit;
  description
    "V bit, when set, the router is an endpoint of one or
    more virtual links.";
```

```
}

identity asbr-bit {
  base router-lsa-bit;
  description
    "E bit, when set, the router is an AS Boundary
    Router (ASBR).";
}

identity abr-bit {
  base router-lsa-bit;
  description
    "B bit, when set, the router is an Area Border
    Router (ABR).";
}

identity nssa-bit {
  base router-lsa-bit;
  description
    "Nt bit, when set, the router is an NSSA border router
    that is unconditionally translating NSSA LSAs into
    AS-external LSAs.";
}

identity ospfv3-lsa-option {
  description
    "Base identity for OSPF LSA options flags.";
}

identity af-bit {
  base ospfv3-lsa-option;
  description
    "AF bit, when set, the router supports OSPFv3 Address
    Families as in RFC5838.";
}

identity dc-bit {
  base ospfv3-lsa-option;
  description
    "DC bit, when set, the router supports demand circuits.";
}

identity r-bit {
  base ospfv3-lsa-option;
  description
    "R bit, when set, the originator is an active router.";
}
```

```
identity n-bit {
  base ospfv3-lsa-option;
  description
    "N bit, when set, the router is attached to an NSSA";
}

identity e-bit {
  base ospfv3-lsa-option;
  description
    "E bit, this bit describes the way AS-external LSAs
    are flooded";
}

identity v6-bit {
  base ospfv3-lsa-option;
  description
    "V6 bit, if clear, the router/link should be excluded
    from IPv6 routing calculation";
}

identity ospfv3-prefix-option {
  description
    "Base identity for OSPFv3 Prefix Options.";
}

identity nu-bit {
  base ospfv3-prefix-option;
  description
    "NU Bit, when set, the prefix should be excluded
    from IPv6 unicast calculations.";
}

identity la-bit {
  base ospfv3-prefix-option;
  description
    "LA bit, when set, the prefix is actually an IPv6
    interface address of the Advertising Router.";
}

identity p-bit {
  base ospfv3-prefix-option;
  description
    "P bit, when set, the NSSA area prefix should be
    translated to an AS External LSA and advertised
    by the translating NSSA Border Router.";
}

identity dn-bit {
```

```
    base ospfv3-prefix-option;
    description
      "DN bit, when set, the inter-area-prefix LSA or
      AS-external LSA prefix has been advertised as an
      L3VPN prefix.";
  }

  identity ospfv2-lsa-option {
    description
      "Base identity for OSPFv2 LSA option flags.";
  }

  identity mt-bit {
    base ospfv2-lsa-option;
    description
      "MT bit, When set, the router supports multi-topology as
      in RFC 4915.";
  }

  identity v2-dc-bit {
    base ospfv2-lsa-option;
    description
      "DC bit, When set, the router supports demand circuits.";
  }

  identity v2-p-bit {
    base ospfv2-lsa-option;
    description
      "P bit, wnlly used in type-7 LSA. When set, an NSSA
      border router should translate the type-7 LSA
      to a type-5 LSA.";
  }

  identity mc-flag {
    base ospfv2-lsa-option;
    description
      "MC Bit, when set, the router supports MOSPF.";
  }

  identity v2-e-flag {
    base ospfv2-lsa-option;
    description
      "E Bit, this bit describes the way AS-external LSAs
      are flooded.";
  }

  identity o-bit {
    base ospfv2-lsa-option;
```

```
    description
      "O bit, when set, the router is opaque-capable as in
       RFC 5250.";
  }

  identity v2-dn-bit {
    base ospfv2-lsa-option;
    description
      "DN bit, when a type 3, 5 or 7 LSA is sent from a PE
       to a CE, the DN bit must be set. See RFC 4576.";
  }

  identity ospfv2-extended-prefix-flag {
    description
      "Base identity for extended prefix TLV flag.";
  }

  identity a-flag {
    base ospfv2-extended-prefix-flag;
    description
      "Attach flag, when set it indicates that the prefix
       corresponds and a route what is directly connected to
       the advertising router..";
  }

  identity node-flag {
    base ospfv2-extended-prefix-flag;
    description
      "Node flag, when set, it indicates that the prefix is
       used to represent the advertising node, e.g., a loopback
       address.";
  }

  typedef ospf-metric {
    type uint32 {
      range "0 .. 16777215";
    }
    description
      "OSPF Metric - 24-bit unsigned integer.";
  }

  typedef ospf-link-metric {
    type uint16 {
      range "0 .. 65535";
    }
    description
      "OSPF Link Metric - 16-bit unsigned integer.";
  }
```

```
typedef opaque-id {
  type uint32 {
    range "0 .. 16777215";
  }
  description
    "Opaque ID - 24-bit unsigned integer.";
}

typedef area-id-type {
  type yang:dotted-quad;
  description
    "Area ID type.";
}

typedef route-type {
  type enumeration {
    enum intra-area {
      description "OSPF intra-area route.";
    }
    enum inter-area {
      description "OSPF inter-area route.";
    }
    enum external-1 {
      description "OSPF type 1 external route.";
    }
    enum external-2 {
      description "OSPF type 2 external route.";
    }
    enum nssa-1 {
      description "OSPF type 1 NSSA route.";
    }
    enum nssa-2 {
      description "OSPF type 2 NSSA route.";
    }
  }
  description "OSPF route type.";
}

typedef if-state-type {
  type enumeration {
    enum down {
      value "1";
      description
        "Interface down state.";
    }
    enum loopback {
      value "2";
      description

```



```
        "Interface loopback state.";
    }
    enum waiting {
        value "3";
        description
            "Interface waiting state.";
    }
    enum point-to-point {
        value "4";
        description
            "Interface point-to-point state.";
    }
    enum dr {
        value "5";
        description
            "Interface Designated Router (DR) state.";
    }
    enum bdr {
        value "6";
        description
            "Interface Backup Designated Router (BDR) state.";
    }
    enum dr-other {
        value "7";
        description
            "Interface Other Designated Router state.";
    }
}
description
    "OSPF interface state type.";
}

typedef router-link-type {
    type enumeration {
        enum point-to-point-link {
            value "1";
            description
                "Point-to-Point link to Router";
        }
        enum transit-network-link {
            value "2";
            description
                "Link to transit network identified by
                Designated-Router (DR) ";
        }
        enum stub-network-link {
            value "3";
            description
```

```
        "Link to stub network identified by subnet";
    }
    enum virtual-link {
        value "4";
        description
            "Virtual link across transit area";
    }
}
description
    "OSPF Router Link Type.";
}

typedef nbr-state-type {
    type enumeration {
        enum down {
            value "1";
            description
                "Neighbor down state.";
        }
        enum attempt {
            value "2";
            description
                "Neighbor attempt state.";
        }
        enum init {
            value "3";
            description
                "Neighbor init state.";
        }
        enum 2-way {
            value "4";
            description
                "Neighbor 2-Way state.";
        }
        enum exstart {
            value "5";
            description
                "Neighbor exchange start state.";
        }
        enum exchange {
            value "6";
            description
                "Neighbor exchange state.";
        }
        enum loading {
            value "7";
            description
                "Neighbor loading state.";
        }
    }
}
```

```
    }
    enum full {
        value "8";
        description
            "Neighbor full state.";
    }
}
description
    "OSPF neighbor state type.";
}

typedef restart-helper-status-type {
    type enumeration {
        enum not-helping {
            value "1";
            description
                "Restart helper status not helping.";
        }
        enum helping {
            value "2";
            description
                "Restart helper status helping.";
        }
    }
    description
        "Restart helper status type.";
}

typedef restart-exit-reason-type {
    type enumeration {
        enum none {
            value "1";
            description
                "Restart not attempted.";
        }
        enum in-progress {
            value "2";
            description
                "Restart in progress.";
        }
        enum completed {
            value "3";
            description
                "Restart successfully completed.";
        }
        enum timed-out {
            value "4";
            description
```

```
        "Restart timed out.";
    }
    enum topology-changed {
        value "5";
        description
            "Restart aborted due to topology change.";
    }
}
description
    "Describes the outcome of the last attempt at a
    graceful restart, either by itself or acting
    as a helper.";
}

typedef packet-type {
    type enumeration {
        enum hello {
            value "1";
            description
                "OSPF Hello packet.";
        }
        enum database-description {
            value "2";
            description
                "OSPF Database Description packet.";
        }
        enum link-state-request {
            value "3";
            description
                "OSPF Link State Request packet.";
        }
        enum link-state-update {
            value "4";
            description
                "OSPF Link State Update packet.";
        }
        enum link-state-ack {
            value "5";
            description
                "OSPF Link State Acknowledgement packet.";
        }
    }
}
description
    "OSPF packet type.";
}

typedef nssa-translator-state-type {
    type enumeration {
```

```
    enum enabled {
      value "1";
      description
        "NSSA translator enabled state.";
    }
    enum elected {
      value "2";
      description
        "NSSA translator elected state.";
    }
    enum disabled {
      value "3";
      description
        "NSSA translator disabled state.";
    }
  }
  description
    "OSPF NSSA translator state type.";
}

typedef restart-status-type {
  type enumeration {
    enum not-restarting {
      value "1";
      description
        "Router is not restarting.";
    }
    enum planned-restart {
      value "2";
      description
        "Router is going through planned restart.";
    }
    enum unplanned-restart {
      value "3";
      description
        "Router is going through unplanned restart.";
    }
  }
  description
    "OSPF graceful restart status type.";
}

typedef fletcher-checksum16-type {
  type string {
    pattern '(0x)?[0-9a-fA-F]{4}';
  }
  description
    "Fletcher 16-bit checksum in hex-string format 0XXXXX.";
```

```
        reference "RFC 905: ISO Transport Protocol specification
                  ISO DP 8073";
    }

    typedef ospfv2-auth-trailer-rfc-version {
        type enumeration {
            enum rfc5709 {
                description
                    "Support OSPF Authentication Trailer as
                     described in RFC 5709";
                reference "RFC 5709: OSPFv2 HMAC-SHA Cryptographic
                          Authentication";
            }

            enum rfc7474 {
                description
                    "Support OSPF Authentication Trailer as
                     described in RFC 7474";
                reference
                    "RFC 7474: Security Extension for OSPFv2
                     When Using Manual Key Management Authentication";
            }
        }
        description
            "OSPFv2 Authentication Trailer Support";
    }

    grouping tlv {
        description
            "Type-Length-Value (TLV)";
        leaf type {
            type uint16;
            description "TLV type.";
        }
        leaf length {
            type uint16;
            description "TLV length (octets).";
        }
        leaf value {
            type yang:hex-string;
            description "TLV value.";
        }
    }

    grouping unknown-tlvs {
        description
            "Unknown TLVs grouping - Used for unknown TLVs or
```

```
        unknown sub-TLVs.";
    container unknown-tlvs {
        description "All unknown TLVs.";
        list unknown-tlv {
            description "Unknown TLV.";
            uses tlv;
        }
    }
}

grouping node-tag-tlv {
    description "OSPF Node Admin Tag TLV grouping.";
    list node-tag {
        leaf tag {
            type uint32;
            description
                "Node admin tag value.";
        }
        description
            "List of tags.";
    }
}

grouping router-capabilities-tlv {
    description "OSPF Router Capabilities TLV grouping.";
    reference "RFC 7770: OSPF Router Capabilities";
    container router-informational-capabilities {
        leaf-list informational-capabilities {
            type identityref {
                base informational-capability;
            }
            description
                "Informational capability list. This list will
                contains the identities for the informational
                capabilities supported by router.";
        }
        description
            "OSPF Router Informational Flag Definitions.";
    }
    list informational-capabilities-flags {
        leaf informational-flag {
            type uint32;
            description
                "Individual informational capability flag.";
        }
        description
            "List of informational capability flags. This will
            return all the 32-bit informational flags irrespective
```

```
        of whether or not they are known to the device.";
    }
    list functional-capabilities {
        leaf functional-flag {
            type uint32;
            description
                "Individual functional capability flag.";
        }
        description
            "List of functional capability flags. This will
            return all the 32-bit functional flags irrespective
            of whether or not they are known to the device.";
    }
}

grouping dynamic-hostname-tlv {
    description "Dynamic Hostname TLV";
    reference "RFC 5642: Dynamic Hostnames for OSPF";
    leaf hostname {
        type string {
            length "1..255";
        }
        description "Dynamic Hostname";
    }
}

grouping sbfd-discriminator-tlv {
    description "Seamless BFD Discriminator TLV";
    reference "RFC 7884: S-BFD Discriminators in OSPF";
    list sbfd-discriminators {
        leaf sbfd-discriminator {
            type uint32;
            description "Individual S-BFD Discriminator.";
        }
        description
            "List of S-BFD Discriminators";
    }
}

grouping maximum-sid-depth-tlv {
    description "Maximum SID Depth (MSD) TLV";
    reference
        "RFC 8476: Signaling Maximum Segment Depth (MSD)
        using OSPF";
    list msd-type {
        leaf msd-type {
            type uint8;
            description "Maximum Segment Depth (MSD) type";
        }
    }
}
```



```
    }
    leaf msd-value {
      type uint8;
      description
        "Maximum Segment Depth (MSD) value for the type";
    }
    description
      "List of Maximum Segment Depth (MSD) tuples";
  }
}

grouping ospf-router-lsa-bits {
  container router-bits {
    leaf-list rtr-lsa-bits {
      type identityref {
        base router-lsa-bit;
      }
      description
        "Router LSA bits list. This list will contain
        identities for the bits which are set in the
        Router-LSA bits.";
    }
    description "Router LSA Bits.";
  }
  description
    "Router LSA Bits - Currently common for OSPFv2 and
    OSPFv3 but it may diverge with future augmentations.";
}

grouping ospfv2-router-link {
  description "OSPFv2 router link.";
  leaf link-id {
    type union {
      type inet:ipv4-address;
      type yang:dotted-quad;
    }
    description "Router-LSA Link ID";
  }
  leaf link-data {
    type union {
      type inet:ipv4-address;
      type uint32;
    }
    description "Router-LSA Link data.";
  }
  leaf type {
    type router-link-type;
    description "Router-LSA Link type.";
  }
}
```

```
    }  
  }  
  
  grouping ospfv2-lsa-body {  
    description "OSPFv2 LSA body.";  
    container router {  
      when "derived-from-or-self ../../header/type, "  
        + "'ospfv2-router-lsa'" {  
        description  
          "Only applies to Router-LSAs.";  
      }  
      description  
        "Router LSA.";  
      uses ospf-router-lsa-bits;  
      leaf num-of-links {  
        type uint16;  
        description "Number of links in Router LSA.";  
      }  
      container links {  
        description "All router Links.";  
        list link {  
          description "Router LSA link.";  
          uses ospfv2-router-link;  
          container topologies {  
            description "All topologies for the link.";  
            list topology {  
              description  
                "Topology specific information.";  
              leaf mt-id {  
                type uint8;  
                description  
                  "The MT-ID for the topology enabled on  
                  the link.";  
              }  
              leaf metric {  
                type uint16;  
                description "Metric for the topology.";  
              }  
            }  
          }  
        }  
      }  
    }  
  }  
  container network {  
    when "derived-from-or-self ../../header/type, "  
      + "'ospfv2-network-lsa'" {  
      description  
        "Only applies to Network LSAs.";  
    }  
  }  
}
```

```
    }
    description
      "Network LSA.";
    leaf network-mask {
      type yang:dotted-quad;
      description
        "The IP address mask for the network.";
    }
    container attached-routers {
      description "All attached routers.";
      leaf-list attached-router {
        type inet:ipv4-address;
        description
          "List of the routers attached to the network.";
      }
    }
  }
  container summary {
    when "derived-from ../../header/type, "
      + "'ospfv2-summary-lsa-type'" {
      description
        "Only applies to Summary LSAs.";
    }
    description
      "Summary LSA.";
    leaf network-mask {
      type inet:ipv4-address;
      description
        "The IP address mask for the network";
    }
    container topologies {
      description "All topologies for the summary LSA.";
      list topology {
        description
          "Topology specific information.";
        leaf mt-id {
          type uint8;
          description
            "The MT-ID for the topology enabled for
              the summary.";
        }
        leaf metric {
          type ospf-metric;
          description "Metric for the topology.";
        }
      }
    }
  }
}
```

```
container external {
  when "derived-from ../../header/type, "
    + "'ospfv2-external-lsa-type'" {
    description
      "Only applies to AS-external LSAs and NSSA LSAs.";
  }
  description
    "External LSA.";
  leaf network-mask {
    type inet:ipv4-address;
    description
      "The IP address mask for the network";
  }
  container topologies {
    description "All topologies for the external.";
    list topology {
      description
        "Topology specific information.";
      leaf mt-id {
        type uint8;
        description
          "The MT-ID for the topology enabled for the
            external or NSSA prefix.";
      }
      leaf flags {
        type bits {
          bit E {
            description
              "When set, the metric specified is a Type 2
                external metric.";
          }
        }
        description "Flags.";
      }
      leaf metric {
        type ospf-metric;
        description "Metric for the topology.";
      }
      leaf forwarding-address {
        type inet:ipv4-address;
        description
          "Forwarding address.";
      }
      leaf external-route-tag {
        type uint32;
        description
          "Route tag for the topology.";
      }
    }
  }
}
```

```
    }
  }
}
container opaque {
  when "derived-from ../../header/type, "
    + "'ospfv2-opaque-lsa-type'" {
    description
      "Only applies to Opaque LSAs.";
  }
  description
    "Opaque LSA.";

  container ri-opaque {
    description "OSPF Router Information (RI) opaque LSA.";
    reference "RFC 7770: OSPF Router Capabilities";

    container router-capabilities-tlv {
      description
        "Informational and functional router capabilities";
      uses router-capabilities-tlv;
    }

    container node-tag-tlvs {
      description
        "All node tag TLVs.";
      list node-tag-tlv {
        description
          "Node tag TLV.";
        uses node-tag-tlv;
      }
    }

    container dynamic-hostname-tlv {
      description "OSPF Dynamic Hostname";
      uses dynamic-hostname-tlv;
    }

    container sbfd-discriminator-tlv {
      description "OSPF S-BFD Discriminators";
      uses sbfd-discriminator-tlv;
    }

    container maximum-sid-depth-tlv {
      description "OSPF Maximum SID Depth (MSD) values";
      uses maximum-sid-depth-tlv;
    }
    uses unknown-tlvs;
  }
}
```

```
container te-opaque {
  description "OSPFv2 Traffic Engineering (TE) opaque LSA.";
  reference "RFC 3630: Traffic Engineering (TE)
    Extensions to OSPFv2";

  container router-address-tlv {
    description
      "Router address TLV.";
    leaf router-address {
      type inet:ipv4-address;
      description
        "Router address.";
    }
  }
}

container link-tlv {
  description "Describes a single link, and it is constructed
    of a set of Sub-TLVs.";
  leaf link-type {
    type router-link-type;
    mandatory true;
    description "Link type.";
  }
  leaf link-id {
    type union {
      type inet:ipv4-address;
      type yang:dotted-quad;
    }
    mandatory true;
    description "Link ID.";
  }
  container local-if-ipv4-addrs {
    description "All local interface IPv4 addresses.";
    leaf-list local-if-ipv4-addr {
      type inet:ipv4-address;
      description
        "List of local interface IPv4 addresses.";
    }
  }
  container remote-if-ipv4-addrs {
    description "All remote interface IPv4 addresses.";
    leaf-list remote-if-ipv4-addr {
      type inet:ipv4-address;
      description
        "List of remote interface IPv4 addresses.";
    }
  }
  leaf te-metric {
```

```
        type uint32;
        description "TE metric.";
    }
    leaf max-bandwidth {
        type rt-types:bandwidth-ieee-float32;
        description "Maximum bandwidth.";
    }
    leaf max-reservable-bandwidth {
        type rt-types:bandwidth-ieee-float32;
        description "Maximum reservable bandwidth.";
    }
    container unreserved-bandwidths {
        description "All unreserved bandwidths.";
        list unreserved-bandwidth {
            leaf priority {
                type uint8 {
                    range "0 .. 7";
                }
                description "Priority from 0 to 7.";
            }
            leaf unreserved-bandwidth {
                type rt-types:bandwidth-ieee-float32;
                description "Unreserved bandwidth.";
            }
            description
                "List of unreserved bandwidths for different
                priorities.";
        }
    }
    leaf admin-group {
        type uint32;
        description
            "Administrative group/Resource Class/Color.";
    }
    uses unknown-tlvs;
}

container extended-prefix-opaque {
    description "All extended prefix TLVs in the LSA.";
    list extended-prefix-tlv {
        description "Extended prefix TLV.";
        leaf route-type {
            type enumeration {
                enum unspecified {
                    value "0";
                    description "Unspecified.";
                }
            }
        }
    }
}
```

```
    enum intra-area {
      value "1";
      description "OSPF intra-area route.";
    }
    enum inter-area {
      value "3";
      description "OSPF inter-area route.";
    }
    enum external {
      value "5";
      description "OSPF External route.";
    }
    enum nssa {
      value "7";
      description "OSPF NSSA external route.";
    }
  }
  description "Route type.";
}
container flags {
  leaf-list extended-prefix-flags {
    type identityref {
      base ospfv2-extended-prefix-flag;
    }
    description
      "Extended prefix TLV flags list. This list will
       contain identities for the prefix flags that
       are set in the extended prefix flags.";
  }
  description "Prefix Flags.";
}
leaf prefix {
  type inet:ip-prefix;
  description "Address prefix.";
}
uses unknown-tlvs;
}

container extended-link-opaque {
  description "All extended link TLVs in the LSA.";
  container extended-link-tlv {
    description "Extended link TLV.";
    uses ospfv2-router-link;
    container maximum-sid-depth-tlv {
      description "OSPF Maximum SID Depth (MSD) values";
      uses maximum-sid-depth-tlv;
    }
  }
}
```



```
        uses unknown-tlvs;
    }
}

grouping ospfv3-lsa-options {
    description "OSPFv3 LSA options";
    container lsa-options {
        leaf-list lsa-options {
            type identityref {
                base ospfv3-lsa-option;
            }
            description
                "OSPFv3 LSA Option flags list. This list will contain
                 the identities for the OSPFv3 LSA options that are
                 set for the LSA.";
        }
        description "OSPFv3 LSA options.";
    }
}

grouping ospfv3-lsa-prefix {
    description
        "OSPFv3 LSA prefix.";

    leaf prefix {
        type inet:ip-prefix;
        description
            "LSA Prefix.";
    }
    container prefix-options {
        leaf-list prefix-options {
            type identityref {
                base ospfv3-prefix-option;
            }
            description
                "OSPFv3 prefix option flag list. This list will
                 contain the identities for the OSPFv3 options
                 that are set for the OSPFv3 prefix.";
        }
        description "Prefix options.";
    }
}

grouping ospfv3-lsa-external {
    description
        "AS-External and NSSA LSA.";
```

```
leaf metric {
  type ospf-metric;
  description "Metric";
}
leaf flags {
  type bits {
    bit E {
      description
        "When set, the metric specified is a Type 2
        external metric.";
    }
    bit F {
      description
        "When set, a Forwarding Address is included
        in the LSA.";
    }
    bit T {
      description
        "When set, an External Route Tag is included
        in the LSA.";
    }
  }
  description "Flags.";
}

leaf referenced-ls-type {
  type identityref {
    base ospfv3-lsa-type;
  }
  description "Referenced Link State type.";
}
leaf unknown-referenced-ls-type {
  type uint16;
  description
    "Value for an unknown Referenced Link State type.";
}

uses ospfv3-lsa-prefix;

leaf forwarding-address {
  type inet:ipv6-address;
  description
    "Forwarding address.";
}

leaf external-route-tag {
  type uint32;
  description
```

```
        "Route tag.";
    }
    leaf referenced-link-state-id {
        type uint32;
        description
            "Referenced Link State ID.";
    }
}

grouping ospfv3-lsa-body {
    description "OSPFv3 LSA body.";
    container router {
        when "derived-from-or-self ../../header/type, "
            + "'ospfv3-router-lsa'" {
            description
                "Only applies to Router LSAs.";
        }
        description "Router LSA.";
        uses ospf-router-lsa-bits;
        uses ospfv3-lsa-options;
    }
    container links {
        description "All router link.";
        list link {
            description "Router LSA link.";
            leaf interface-id {
                type uint32;
                description "Interface ID for link.";
            }
            leaf neighbor-interface-id {
                type uint32;
                description "Neighbor's Interface ID for link.";
            }
            leaf neighbor-router-id {
                type rt-types:router-id;
                description "Neighbor's Router ID for link.";
            }
            leaf type {
                type router-link-type;
                description "Link type: 1 - Point-to-Point Link
                               2 - Transit Network Link
                               3 - Stub Network Link
                               4 - Virtual Link";
            }
            leaf metric {
                type uint16;
                description "Link Metric.";
            }
        }
    }
}
```

```
    }
  }
}
container network {
  when "derived-from-or-self ../../header/type, "
    + "'ospfv3-network-lsa'" {
    description
      "Only applies to Network LSAs.";
  }
  description "Network LSA.";

  uses ospfv3-lsa-options;

  container attached-routers {
    description "All attached routers.";
    leaf-list attached-router {
      type rt-types:router-id;
      description
        "List of the routers attached to the network.";
    }
  }
}
container inter-area-prefix {
  when "derived-from-or-self ../../header/type, "
    + "'ospfv3-inter-area-prefix-lsa'" {
    description
      "Only applies to Inter-Area-Prefix LSAs.";
  }
  leaf metric {
    type ospf-metric;
    description "Inter-Area Prefix Metric";
  }
  uses ospfv3-lsa-prefix;
  description "Prefix LSA.";
}
container inter-area-router {
  when "derived-from-or-self ../../header/type, "
    + "'ospfv3-inter-area-router-lsa'" {
    description
      "Only applies to Inter-Area-Router LSAs.";
  }
  uses ospfv3-lsa-options;
  leaf metric {
    type ospf-metric;
    description "AS Boundary Router (ASBR) Metric.";
  }
  leaf destination-router-id {
    type rt-types:router-id;
  }
}
```

```
        description
            "The Router ID of the ASBR described by the LSA.";
    }
    description "Inter-Area-Router LSA.";
}
container as-external {
    when "derived-from-or-self ../../header/type, "
        + "'ospfv3-as-external-lsa'" {
        description
            "Only applies to AS-external LSAs.";
    }

    uses ospfv3-lsa-external;

    description "AS-External LSA.";
}
container nssa {
    when "derived-from-or-self ../../header/type, "
        + "'ospfv3-nssa-lsa'" {
        description
            "Only applies to NSSA LSAs.";
    }
    uses ospfv3-lsa-external;

    description "NSSA LSA.";
}
container link {
    when "derived-from-or-self ../../header/type, "
        + "'ospfv3-link-lsa'" {
        description
            "Only applies to Link LSAs.";
    }
}
leaf rtr-priority {
    type uint8;
    description
        "Router priority for DR election. A router with a
        higher priority will be preferred in the election
        and a value of 0 indicates the router is not
        eligible to become Designated Router or Backup
        Designated Router (BDR).";
}
uses ospfv3-lsa-options;

leaf link-local-interface-address {
    type inet:ipv6-address;
    description
        "The originating router's link-local
        interface address for the link.";
```

```
    }

    leaf num-of-prefixes {
      type uint32;
      description "Number of prefixes.";
    }

    container prefixes {
      description "All prefixes for the link.";
      list prefix {
        description
          "List of prefixes associated with the link.";
        uses ospfv3-lsa-prefix;
      }
    }
    description "Link LSA.";
  }
  container intra-area-prefix {
    when "derived-from-or-self ../../header/type, "
      + "'ospfv3-intra-area-prefix-lsa'" {
      description
        "Only applies to Intra-Area-Prefix LSAs.";
    }
    description "Intra-Area-Prefix LSA.";

    leaf referenced-ls-type {
      type identityref {
        base ospfv3-lsa-type;
      }
      description "Referenced Link State type.";
    }
    leaf unknown-referenced-ls-type {
      type uint16;
      description
        "Value for an unknown Referenced Link State type.";
    }
    leaf referenced-link-state-id {
      type uint32;
      description
        "Referenced Link State ID.";
    }
    leaf referenced-adv-router {
      type rt-types:router-id;
      description
        "Referenced Advertising Router.";
    }
  }

  leaf num-of-prefixes {
```

```
        type uint16;
        description "Number of prefixes.";
    }
    container prefixes {
        description "All prefixes in this LSA.";
        list prefix {
            description "List of prefixes in this LSA.";
            uses ospfv3-lsa-prefix;
            leaf metric {
                type ospf-metric;
                description "Prefix Metric.";
            }
        }
    }
}
container router-information {
    when "derived-from-or-self ../../header/type, "
        + "'ospfv3-router-information-lsa'" {
        description
            "Only applies to Router Information LSAs (RFC7770).";
    }
    container router-capabilities-tlv {
        description
            "Informational and functional router capabilities";
        uses router-capabilities-tlv;
    }
    container node-tag-tlvs {
        description
            "All node tag tlvs.";
        list node-tag-tlv {
            description
                "Node tag tlv.";
            uses node-tag-tlv;
        }
    }
    container dynamic-hostname-tlv {
        description "OSPF Dynamic Hostname";
        uses dynamic-hostname-tlv;
    }
    container sbfd-discriminator-tlv {
        description "OSPF S-BFD Discriminators";
        uses sbfd-discriminator-tlv;
    }
    description "Router Information LSA.";
    reference "RFC 7770: Extensions for Advertising Router
        Capabilities";
}
}
```

```
grouping lsa-header {
  description
    "Common LSA for OSPFv2 and OSPFv3";
  leaf age {
    type uint16;
    mandatory true;
    description "LSA age.";
  }
  leaf type {
    type identityref {
      base ospf-lsa-type;
    }
    mandatory true;
    description "LSA type";
  }
  leaf adv-router {
    type rt-types:router-id;
    mandatory true;
    description "LSA advertising router.";
  }
  leaf seq-num {
    type uint32;
    mandatory true;
    description "LSA sequence number.";
  }
  leaf checksum {
    type fletcher-checksum16-type;
    mandatory true;
    description "LSA checksum.";
  }
  leaf length {
    type uint16;
    mandatory true;
    description "LSA length including the header.";
  }
}

grouping ospfv2-lsa {
  description
    "OSPFv2 LSA - LSAs are uniquely identified by
    the <LSA Type, Link-State ID, Advertising Router>
    tuple with the sequence number differentiating
    LSA instances.";
  container header {
    must "(derived-from(type, "
      + "'ospfv2-opaque-lsa-type') and "
      + "opaque-id and opaque-type) or "
      + "(not(derived-from(type, "
```



```
        + "'ospfv2-opaque-lsa-type')) "
        + "and not(opaque-id) and not(opaque-type))" {
    description
        "Opaque type and ID only apply to Opaque LSAs.";
}
description
    "Decoded OSPFv2 LSA header data.";

container lsa-options {
    leaf-list lsa-options {
        type identityref {
            base ospfv2-lsa-option;
        }
        description
            "LSA option flags list. This list will contain
             the identities for the identities for the OSPFv2
             LSA options that are set.";
    }
    description
        "LSA options.";
}

leaf lsa-id {
    type yang:dotted-quad;
    mandatory true;
    description "Link-State ID.";
}

leaf opaque-type {
    type uint8;
    description "Opaque type.";
}

leaf opaque-id {
    type opaque-id;
    description "Opaque ID.";
}

uses lsa-header;
}
container body {
    description
        "Decoded OSPFv2 LSA body data.";
    uses ospfv2-lsa-body;
}
}

grouping ospfv3-lsa {
```

```
description
    "Decoded OSPFv3 LSA.";
container header {
    description
        "Decoded OSPFv3 LSA header data.";
    leaf lsa-id {
        type uint32;
        mandatory true;
        description "OSPFv3 LSA ID.";
    }
    uses lsa-header;
}
container body {
    description
        "Decoded OSPF LSA body data.";
    uses ospfv3-lsa-body;
}
}
grouping lsa-common {
    description
        "Common fields for OSPF LSA representation.";
    leaf decode-completed {
        type boolean;
        description
            "The OSPF LSA body was successfully decoded other than
            unknown TLVs. Unknown LSAs types and OSPFv2 unknown
            opaque LSA types are not decoded. Additionally,
            malformed LSAs are generally not accepted and will
            not be in the Link State Database.";
    }
    leaf raw-data {
        type yang:hex-string;
        description
            "The complete LSA in network byte
            order hexadecimal as received or originated.";
    }
}
}
grouping lsa {
    description
        "OSPF LSA.";
    uses lsa-common;
    choice version {
        description
            "OSPFv2 or OSPFv3 LSA body.";
        container ospfv2 {
            description "OSPFv2 LSA";
            uses ospfv2-lsa;
        }
    }
}
```

```
    }
    container ospfv3 {
      description "OSPFv3 LSA";
      uses ospfv3-lsa;
    }
  }
}

grouping lsa-key {
  description
    "OSPF LSA key - the database key for each LSA of a given
    type in the Link State DataBase (LSDB).";
  leaf lsa-id {
    type union {
      type yang:dotted-quad;
      type uint32;
    }
    description
      "Link-State ID.";
  }
  leaf adv-router {
    type rt-types:router-id;
    description
      "Advertising router.";
  }
}

grouping instance-stat {
  description "Per-instance statistics";
  leaf discontinuity-time {
    type yang:date-and-time;
    description
      "The time on the most recent occasion at which any one or
      more of this OSPF instance's counters suffered a
      discontinuity. If no such discontinuities have occurred
      since the OSPF instance was last re-initialized, then
      this node contains the time the OSPF instance was
      re-initialized which normally occurs when it was
      created.";
  }
  leaf originate-new-lsa-count {
    type yang:counter32;
    description
      "The number of new LSAs originated. Discontinuities in the
      value of this counter can occur when the OSPF instance is
      re-initialized.";
  }
  leaf rx-new-lsas-count {
```

```
    type yang:counter32;
    description
      "The number of new LSAs received. Discontinuities in the
       value of this counter can occur when the OSPF instance is
       re-initialized.";
  }
  leaf as-scope-lsa-count {
    type yang:gauge32;
    description "The number of AS-scope LSAs.";
  }
  leaf as-scope-lsa-chksum-sum {
    type uint32;
    description
      "The module 2**32 sum of the LSA checksums
       for AS-scope LSAs. The value should be treated as
       unsigned when comparing two sums of checksums. While
       differing checksums indicate a different combination
       of LSAs, equivalent checksums don't guarantee that the
       LSAs are the same given that multiple combinations of
       LSAs can result in the same checksum.";
  }
  container database {
    description "Container for per AS-scope LSA statistics.";
    list as-scope-lsa-type {
      description "List of AS-scope LSA statistics";
      leaf lsa-type {
        type uint16;
        description "AS-Scope LSA type.";
      }
      leaf lsa-count {
        type yang:gauge32;
        description "The number of LSAs of the LSA type.";
      }
      leaf lsa-cksum-sum {
        type uint32;
        description
          "The module 2**32 sum of the LSA checksums
           for the LSAs of this type. The value should be
           treated as unsigned when comparing two sums of
           checksums. While differing checksums indicate a
           different combination of LSAs, equivalent checksums
           don't guarantee that the LSAs are the same given that
           multiple combinations of LSAs can result in the same
           checksum.";
      }
    }
  }
}
uses instance-fast-reroute-state;
```

```
}  
  
grouping area-stat {  
  description "Per-area statistics.";   
  leaf discontinuity-time {  
    type yang:date-and-time;  
    description  
      "The time on the most recent occasion at which any one or  
      more of this OSPF area's counters suffered a  
      discontinuity. If no such discontinuities have occurred  
      since the OSPF area was last re-initialized, then  
      this node contains the time the OSPF area was  
      re-initialized which normally occurs when it was  
      created.";   
  }  
  leaf spf-runs-count {  
    type yang:counter32;  
    description  
      "The number of times the intra-area SPF has run.  
      Discontinuities in the value of this counter can occur  
      when the OSPF area is re-initialized.";   
  }  
  leaf abr-count {  
    type yang:gauge32;  
    description  
      "The total number of Area Border Routers (ABRs)  
      reachable within this area.";   
  }  
  leaf asbr-count {  
    type yang:gauge32;  
    description  
      "The total number of AS Boundary Routers (ASBRs).";   
  }  
  leaf ar-nssa-translator-event-count {  
    type yang:counter32;  
    description  
      "The number of NSSA translator-state changes.  
      Discontinuities in the value of this counter can occur  
      when the OSPF area is re-initialized.";   
  }  
  leaf area-scope-lsa-count {  
    type yang:gauge32;  
    description  
      "The number of area-scope LSAs in the area.";   
  }  
  leaf area-scope-lsa-cksum-sum {  
    type uint32;  
    description
```

```
    "The module 2**32 sum of the LSA checksums
    for area-scope LSAs. The value should be treated as
    unsigned when comparing two sums of checksums. While
    differing checksums indicate a different combination
    of LSAs, equivalent checksums don't guarantee that the
    LSAs are the same given that multiple combinations of
    LSAs can result in the same checksum.";
  }
  container database {
    description "Container for area-scope LSA type statistics.";
    list area-scope-lsa-type {
      description "List of area-scope LSA statistics";
      leaf lsa-type {
        type uint16;
        description "Area-scope LSA type.";
      }
      leaf lsa-count {
        type yang:gauge32;
        description "The number of LSAs of the LSA type.";
      }
      leaf lsa-cksum-sum {
        type uint32;
        description
          "The module 2**32 sum of the LSA checksums
          for the LSAs of this type. The value should be
          treated as unsigned when comparing two sums of
          checksums. While differing checksums indicate a
          different combination of LSAs, equivalent checksums
          don't guarantee that the LSAs are the same given that
          multiple combinations of LSAs can result in the same
          checksum.";
      }
    }
  }
}

grouping interface-stat {
  description "Per-interface statistics";
  leaf discontinuity-time {
    type yang:date-and-time;
    description
      "The time on the most recent occasion at which any one or
      more of this OSPF interface's counters suffered a
      discontinuity. If no such discontinuities have occurred
      since the OSPF interface was last re-initialized, then
      this node contains the time the OSPF interface was
      re-initialized which normally occurs when it was
      created.";
```

```
}
leaf if-event-count {
  type yang:counter32;
  description
    "The number of times this interface has changed its
    state or an error has occurred. Discontinuities in the
    value of this counter can occur when the OSPF interface
    is re-initialized.";
}
leaf link-scope-lsa-count {
  type yang:gauge32;
  description "The number of link-scope LSAs.";
}
leaf link-scope-lsa-cksum-sum {
  type uint32;
  description
    "The module 2**32 sum of the LSA checksums
    for link-scope LSAs. The value should be treated as
    unsigned when comparing two sums of checksums. While
    differing checksums indicate a different combination
    of LSAs, equivalent checksums don't guarantee that the
    LSAs are the same given that multiple combinations of
    LSAs can result in the same checksum.";
}
container database {
  description "Container for link-scope LSA type statistics.";
  list link-scope-lsa-type {
    description "List of link-scope LSA statistics";
    leaf lsa-type {
      type uint16;
      description "Link scope LSA type.";
    }
    leaf lsa-count {
      type yang:gauge32;
      description "The number of LSAs of the LSA type.";
    }
  }
  leaf lsa-cksum-sum {
    type uint32;
    description
      "The module 2**32 sum of the LSA checksums
      for the LSAs of this type. The value should be
      treated as unsigned when comparing two sums of
      checksums. While differing checksums indicate a
      different combination of LSAs, equivalent checksums
      don't guarantee that the LSAs are the same given that
      multiple combinations of LSAs can result in the same
      checksum.";
  }
}
```

```
    }
  }
}

grouping neighbor-stat {
  description "Per-neighbor statistics.";
  leaf discontinuity-time {
    type yang:date-and-time;
    description
      "The time on the most recent occasion at which any one or
      more of this OSPF neighbor's counters suffered a
      discontinuity. If no such discontinuities have occurred
      since the OSPF neighbor was last re-initialized, then
      this node contains the time the OSPF neighbor was
      re-initialized which normally occurs when the neighbor
      is dynamically discovered and created.";
  }
  leaf nbr-event-count {
    type yang:counter32;
    description
      "The number of times this neighbor has changed
      state or an error has occurred. Discontinuities in the
      value of this counter can occur when the OSPF neighbor
      is re-initialized.";
  }
  leaf nbr-retrans-qlen {
    type yang:gauge32;
    description
      "The current length of the retransmission queue.";
  }
}

grouping instance-fast-reroute-config {
  description
    "This group defines global configuration of IP
    Fast ReRoute (FRR).";
  container fast-reroute {
    if-feature fast-reroute;
    description
      "This container may be augmented with global
      parameters for IP-FRR.";
    container lfa {
      if-feature lfa;
      description
        "This container may be augmented with
        global parameters for Loop-Free Alternatives (LFA).
        Container creation has no effect on LFA activation.";
    }
  }
}
```



```
    }  
  }  
  
  grouping instance-fast-reroute-state {  
    description "IP-FRR state data grouping";  
  
    container protected-routes {  
      if-feature fast-reroute;  
      config false;  
      description "Instance protection statistics";  
  
      list address-family-stats {  
        key "address-family prefix alternate";  
        description  
          "Per Address Family protected prefix information";  
  
        leaf address-family {  
          type iana-rt-types:address-family;  
          description  
            "Address-family";  
        }  
        leaf prefix {  
          type inet:ip-prefix;  
          description  
            "Protected prefix.";  
        }  
        leaf alternate {  
          type inet:ip-address;  
          description  
            "Alternate next hop for the prefix.";  
        }  
        leaf alternate-type {  
          type enumeration {  
            enum equal-cost {  
              description  
                "ECMP alternate.";  
            }  
            enum lfa {  
              description  
                "LFA alternate.";  
            }  
            enum remote-lfa {  
              description  
                "Remote LFA alternate.";  
            }  
            enum tunnel {  
              description  
                "Tunnel based alternate
```

```
        (like RSVP-TE or GRE).";
    }
    enum ti-lfa {
        description
            "TI-LFA alternate.";
    }
    enum mrt {
        description
            "MRT alternate.";
    }
    enum other {
        description
            "Unknown alternate type.";
    }
}
description
    "Type of alternate.";
}
leaf best {
    type boolean;
    description
        "Indicates that this alternate is preferred.";
}
leaf non-best-reason {
    type string {
        length "1..255";
    }
    description
        "Information field to describe why the alternate
        is not best.";
}
leaf protection-available {
    type bits {
        bit node-protect {
            position 0;
            description
                "Node protection available.";
        }
        bit link-protect {
            position 1;
            description
                "Link protection available.";
        }
        bit srlg-protect {
            position 2;
            description
                "SRLG protection available.";
        }
    }
}
```

```
        bit downstream-protect {
            position 3;
            description
                "Downstream protection available.";
        }
        bit other {
            position 4;
            description
                "Other protection available.";
        }
    }
    description "Protection provided by the alternate.";
}
leaf alternate-metric1 {
    type uint32;
    description
        "Metric from Point of Local Repair (PLR) to
        destination through the alternate path.";
}
leaf alternate-metric2 {
    type uint32;
    description
        "Metric from PLR to the alternate node";
}
leaf alternate-metric3 {
    type uint32;
    description
        "Metric from alternate node to the destination";
}
}
}

container unprotected-routes {
    if-feature fast-reroute;
    config false;
    description "List of prefixes that are not protected";

    list address-family-stats {
        key "address-family prefix";
        description
            "Per Address Family (AF) unprotected prefix statistics.";

        leaf address-family {
            type iana-rt-types:address-family;
            description "Address-family";
        }
        leaf prefix {
            type inet:ip-prefix;
        }
    }
}
```

```
        description "Unprotected prefix.";
    }
}

list protection-statistics {
    key frr-protection-method;
    config false;
    description "List protection method statistics";

    leaf frr-protection-method {
        type string;
        description "Protection method used.";
    }
    list address-family-stats {
        key address-family;
        description "Per Address Family protection statistics.";

        leaf address-family {
            type iana-rt-types:address-family;
            description "Address-family";
        }
        leaf total-routes {
            type uint32;
            description "Total prefixes.";
        }
        leaf unprotected-routes {
            type uint32;
            description
                "Total prefixes that are not protected.";
        }
        leaf protected-routes {
            type uint32;
            description
                "Total prefixes that are protected.";
        }
        leaf linkprotected-routes {
            type uint32;
            description
                "Total prefixes that are link protected.";
        }
        leaf nodeprotected-routes {
            type uint32;
            description
                "Total prefixes that are node protected.";
        }
    }
}
```

```
}

grouping interface-fast-reroute-config {
  description
    "This group defines interface configuration of IP-FRR.";
  container fast-reroute {
    if-feature fast-reroute;
    container lfa {
      if-feature lfa;
      leaf candidate-enable {
        type boolean;
        default true;
        description
          "Enable the interface to be used as backup.";
      }
      leaf enable {
        type boolean;
        default false;
        description
          "Activates LFA - Per-prefix LFA computation
           is assumed.";
      }
      container remote-lfa {
        if-feature remote-lfa;
        leaf enable {
          type boolean;
          default false;
          description
            "Activates Remote LFA (R-LFA).";
        }
      }
      description
        "Remote LFA configuration.";
    }
    description
      "LFA configuration.";
  }
  description
    "Interface IP Fast-reroute configuration.";
}

grouping interface-physical-link-config {
  description
    "Interface cost configuration that only applies to
     physical interfaces (non-virtual) and sham links.";
  leaf cost {
    type ospf-link-metric;
    description
```

```
        "Interface cost.";
    }
    leaf mtu-ignore {
        if-feature mtu-ignore;
        type boolean;
        description
            "Enable/Disable bypassing the MTU mismatch check in
            Database Description packets specified in RFC 2328,
            section 10.6.";
    }
    leaf prefix-suppression {
        if-feature prefix-suppression;
        type boolean;
        description
            "Suppress advertisement of the prefixes associated
            with the interface.";
    }
}

grouping interface-common-config {
    description
        "Common configuration for all types of interfaces,
        including virtual links and sham links.";

    leaf hello-interval {
        type uint16;
        units seconds;
        description
            "Interval between hello packets (seconds). It must
            be the same for all routers on the same network.
            Different networks, implementations, and deployments
            will use different hello-intervals. A sample value
            for a LAN network would be 10 seconds.";
        reference "RFC 2328: OSPF Version 2, Appendix C.3";
    }

    leaf dead-interval {
        type uint16;
        units seconds;
        must "../dead-interval > ../hello-interval" {
            error-message "The dead interval must be "
                + "larger than the hello interval";
        }
        description
            "The value must be greater than the 'hello-interval'.";
    }
    description
        "Interval after which a neighbor is declared down
        (seconds) if hello packets are not received. It is
```

```
        typically 3 or 4 times the hello-interval. A typical
        value for LAN networks is 40 seconds.";
        reference "RFC 2328: OSPF Version 2, Appendix C.3";
    }

    leaf retransmit-interval {
        type uint16 {
            range "1..3600";
        }
        units seconds;
        description
            "Interval between retransmitting unacknowledged Link
            State Advertisements (LSAs) (seconds). This should
            be well over the round-trip transmit delay for
            any two routers on the network. A sample value
            would be 5 seconds.";
        reference "RFC 2328: OSPF Version 2, Appendix C.3";
    }

    leaf transmit-delay {
        type uint16;
        units seconds;
        description
            "Estimated time needed to transmit Link State Update
            (LSU) packets on the interface (seconds). LSAs have
            their age incremented by this amount when advertised
            on the interface. A sample value would be 1 second.";
        reference "RFC 2328: OSPF Version 2, Appendix C.3";
    }

    leaf lls {
        if-feature lls;
        type boolean;
        description
            "Enable/Disable link-local signaling (LLS) support.";
    }

    container ttl-security {
        if-feature ttl-security;
        description "Time to Live (TTL) security check.";
        leaf enable {
            type boolean;
            description
                "Enable/Disable TTL security check.";
        }
        leaf hops {
            type uint8 {
                range "1..254";
            }
        }
    }
}
```

```
    }
    default 1;
    description
        "Maximum number of hops that an OSPF packet may
        have traversed before reception.";
    }
}
leaf enable {
    type boolean;
    default true;
    description
        "Enable/disable OSPF protocol on the interface.";
}

container authentication {
    description "Authentication configuration.";
    choice auth-type-selection {
        description
            "Options for OSPFv2/OSPFv3 authentication
            configuration.";
        case ospfv2-auth {
            when "derived-from-or-self ../../../../rt:type, "
                + "'ospfv2'" {
                description "Applied to OSPFv2 only.";
            }
            leaf ospfv2-auth-trailer-rfc {
                if-feature ospfv2-authentication-trailer;
                type ospfv2-auth-trailer-rfc-version;
                description
                    "Version of OSPFv2 authentication trailer support -
                    RFC 5709 or RFC 7474";
            }
        }
        choice ospfv2-auth-specification {
            description
                "Key chain or explicit key parameter specification";
            case auth-key-chain {
                if-feature key-chain;
                leaf ospfv2-key-chain {
                    type key-chain:key-chain-ref;
                    description
                        "key-chain name.";
                }
            }
            case auth-key-explicit {
                leaf ospfv2-key-id {
                    type uint32;
                    description
                        "Key Identifier";
                }
            }
        }
    }
}
```



```
    }
    leaf ospfv2-key {
      type string;
      description
        "OSPFv2 authentication key. The
         length of the key may be dependent on the
         cryptographic algorithm.";
    }
    leaf ospfv2-crypto-algorithm {
      type identityref {
        base key-chain:crypto-algorithm;
      }
      description
        "Cryptographic algorithm associated with key.";
    }
  }
}
case ospfv3-auth-ipsec {
  when "derived-from-or-self(.../.../.../.../rt:type, "
    + "'ospfv3')" {
    description "Applied to OSPFv3 only.";
  }
  if-feature ospfv3-authentication-ipsec;
  leaf sa {
    type string;
    description
      "Security Association (SA) name.";
  }
}
case ospfv3-auth-trailer {
  when "derived-from-or-self(.../.../.../.../rt:type, "
    + "'ospfv3')" {
    description "Applied to OSPFv3 only.";
  }
  if-feature ospfv3-authentication-trailer;
  choice ospfv3-auth-specification {
    description
      "Key chain or explicit key parameter specification";
    case auth-key-chain {
      if-feature key-chain;
      leaf ospfv3-key-chain {
        type key-chain:key-chain-ref;
        description
          "key-chain name.";
      }
    }
    case auth-key-explicit {
```

```
    leaf ospfv3-sa-id {
      type uint16;
      description
        "Security Association (SA) Identifier";
    }
    leaf ospfv3-key {
      type string;
      description
        "OSPFv3 authentication key. The
        length of the key may be dependent on the
        cryptographic algorithm.";
    }
    leaf ospfv3-crypto-algorithm {
      type identityref {
        base key-chain:crypto-algorithm;
      }
      description
        "Cryptographic algorithm associated with key.";
    }
  }
}
}
}
}
}

grouping interface-config {
  description "Configuration for real interfaces.";

  leaf interface-type {
    type enumeration {
      enum "broadcast" {
        description
          "Specify OSPF broadcast multi-access network.";
      }
      enum "non-broadcast" {
        description
          "Specify OSPF Non-Broadcast Multi-Access
          (NBMA) network.";
      }
      enum "point-to-multipoint" {
        description
          "Specify OSPF point-to-multipoint network.";
      }
      enum "point-to-point" {
        description
          "Specify OSPF point-to-point network.";
      }
    }
  }
}
```

```
    enum "hybrid" {
        if-feature hybrid-interface;
        description
            "Specify OSPF hybrid broadcast/P2MP network.";
    }
}
description
    "Interface type.";
}

leaf passive {
    type boolean;
    description
        "Enable/Disable passive interface - a passive interface's
        prefix will be advertised but no neighbor adjacencies
        will be formed on the interface.";
}

leaf demand-circuit {
    if-feature demand-circuit;
    type boolean;
    description
        "Enable/Disable demand circuit.";
}

leaf priority {
    type uint8;
    description
        "Configure OSPF router priority. On multi-access network
        this value is for Designated Router (DR) election. The
        priority is ignored on other interface types. A router
        with a higher priority will be preferred in the election
        and a value of 0 indicates the router is not eligible to
        become Designated Router or Backup Designated Router
        (BDR).";
}

container multi-areas {
    if-feature multi-area-adj;
    description "Container for multi-area config.";
    list multi-area {
        key multi-area-id;
        description
            "Configure OSPF multi-area adjacency.";
        leaf multi-area-id {
            type area-id-type;
            description
                "Multi-area adjacency area ID.";
        }
    }
}
```

```
    }
    leaf cost {
      type ospf-link-metric;
      description
        "Interface cost for multi-area adjacency.";
    }
  }
}

container static-neighbors {
  description "Statically configured neighbors.";

  list neighbor {
    key "identifier";
    description
      "Specify a static OSPF neighbor.";

    leaf identifier {
      type inet:ip-address;
      description
        "Neighbor Router ID, IPv4 address, or IPv6 address.";
    }

    leaf cost {
      type ospf-link-metric;
      description
        "Neighbor cost. Different implementations have different
        default costs with some defaulting to a cost inversely
        proportional to the interface speed. Others will
        default to 1 equating the cost to a hop count." ;
    }

    leaf poll-interval {
      type uint16;
      units seconds;
      description
        "Neighbor poll interval (seconds) for sending OSPF
        hello packets to discover the neighbor on NBMA
        networks. This interval dictates the granularity for
        discovery of new neighbors. A sample would be
        120 seconds (2 minutes) for a legacy Packet Data
        Network (PDN) X.25 network.";
      reference "RFC 2328: OSPF Version 2, Appendix C.5";
    }

    leaf priority {
      type uint8;
      description
        "Neighbor priority for DR election. A router with a
        higher priority will be preferred in the election

```

```
        and a value of 0 indicates the router is not
        eligible to become Designated Router or Backup
        Designated Router (BDR).";
    }
}

leaf node-flag {
    if-feature node-flag;
    type boolean;
    default false;
    description
        "Set prefix as identifying the advertising router.";
    reference "RFC 7684: OSPFv2 Prefix/Link Attribute
        Advertisement";
}

container bfd {
    if-feature bfd;
    description "BFD Client Configuration.";
    uses bfd-types:client-cfg-parms;
    reference "RFC YYYY: YANG Data Model for Bidirectional
        Forwarding Detection (BFD). Please replace YYYY with
        published RFC number for draft-ietf-bfd-yang.";
}

uses interface-fast-reroute-config;
uses interface-common-config;
uses interface-physical-link-config;
}

grouping neighbor-state {
    description
        "OSPF neighbor operational state.";

    leaf address {
        type inet:ip-address;
        config false;
        description
            "Neighbor address.";
    }

    leaf dr-router-id {
        type rt-types:router-id;
        config false;
        description "Neighbor's Designated Router (DR) Router ID.";
    }

    leaf dr-ip-addr {
```

```
    type inet:ip-address;
    config false;
    description "Neighbor's Designated Router (DR) IP address.";
}

leaf bdr-router-id {
    type rt-types:router-id;
    config false;
    description
        "Neighbor's Backup Designated Router (BDR) Router ID.";
}

leaf bdr-ip-addr {
    type inet:ip-address;
    config false;
    description
        "Neighbor's Backup Designated Router (BDR) IP Address.";
}

leaf state {
    type nbr-state-type;
    config false;
    description
        "OSPF neighbor state.";
}

leaf cost {
    type ospf-link-metric;
    config false;
    description "Cost to reach neighbor for Point-to-Multipoint
        and Hybrid networks";
}

leaf dead-timer {
    type rt-types:timer-value-seconds16;
    config false;
    description "This timer tracks the remaining time before
        the neighbor is declared dead.";
}

container statistics {
    config false;
    description "Per-neighbor statistics";
    uses neighbor-stat;
}

}

grouping interface-common-state {
    description
        "OSPF interface common operational state.";
    reference "RFC2328 Section 9: OSPF Version2 -
        The Interface Data Structure";
}
```

```
leaf state {
  type if-state-type;
  config false;
  description "Interface state.";
}

leaf hello-timer {
  type rt-types:timer-value-seconds16;
  config false;
  description "This timer tracks the remaining time before
               the next hello packet is sent on the
               interface.";
}

leaf wait-timer {
  type rt-types:timer-value-seconds16;
  config false;
  description "This timer tracks the remaining time before
               the interface exits the Waiting state.";
}

leaf dr-router-id {
  type rt-types:router-id;
  config false;
  description "Designated Router (DR) Router ID.";
}

leaf dr-ip-addr {
  type inet:ip-address;
  config false;
  description "Designated Router (DR) IP address.";
}

leaf bdr-router-id {
  type rt-types:router-id;
  config false;
  description "Backup Designated Router (BDR) Router ID.";
}

leaf bdr-ip-addr {
  type inet:ip-address;
  config false;
  description "Backup Designated Router (BDR) IP Address.";
}

container statistics {
  config false;
  description "Per-interface statistics";
}
```

```

    uses interface-stat;
}

container neighbors {
    config false;
    description "All neighbors for the interface.";
    list neighbor {
        key "neighbor-router-id";
        description
            "List of interface OSPF neighbors.";
        leaf neighbor-router-id {
            type rt-types:router-id;
            description
                "Neighbor Router ID.";
        }
        uses neighbor-state;
    }
}

container database {
    config false;
    description "Link-scope Link State Database.";
    list link-scope-lsa-type {
        key "lsa-type";
        description
            "List OSPF link-scope LSAs.";
        leaf lsa-type {
            type uint16;
            description "OSPF link-scope LSA type.";
        }
    }
    container link-scope-lsas {
        description
            "All link-scope LSAs of this LSA type.";
        list link-scope-lsa {
            key "lsa-id adv-router";
            description "List of OSPF link-scope LSAs";
            uses lsa-key;
            uses lsa {
                refine "version/ospfv2/ospfv2" {
                    must "derived-from-or-self( "
                        + ".../.../.../.../.../.../.../.../.../.../..."
                        + "rt:type, 'ospfv2') " {
                        description "OSPFv2 LSA.";
                    }
                }
                refine "version/ospfv3/ospfv3" {
                    must "derived-from-or-self( "
                        + ".../.../.../.../.../.../.../.../.../.../..."
                        + "rt:type, 'ospfv3') " {

```



```
        description "OSPFv3 LSA.";
    }
}
}
}
}
}
}

grouping interface-state {
    description
        "OSPF interface operational state.";
    reference "RFC2328 Section 9: OSPF Version2 -
        The Interface Data Structure";

    uses interface-common-state;
}

grouping virtual-link-config {
    description
        "OSPF virtual link configuration state.";

    uses interface-common-config;
}

grouping virtual-link-state {
    description
        "OSPF virtual link operational state.";

    leaf cost {
        type ospf-link-metric;
        config false;
        description
            "Virtual link interface cost.";
    }
    uses interface-common-state;
}

grouping sham-link-config {
    description
        "OSPF sham link configuration state.";

    uses interface-common-config;
    uses interface-physical-link-config;
}

grouping sham-link-state {
```

```
    description
      "OSPF sham link operational state.";
    uses interface-common-state;
  }

  grouping address-family-area-config {
    description
      "OSPF address-family specific area config state.";

    container ranges {
      description "Container for summary ranges";

      list range {
        key "prefix";
        description
          "Summarize routes matching address/mask -
           Applicable to Area Border Routers (ABRs) only.";
        leaf prefix {
          type inet:ip-prefix;
          description
            "IPv4 or IPv6 prefix";
        }
        leaf advertise {
          type boolean;
          description
            "Advertise or hide.";
        }
        leaf cost {
          type ospf-metric;
          description
            "Advertised cost of summary route.";
        }
      }
    }
  }

  grouping area-common-config {
    description
      "OSPF area common configuration state.";

    leaf summary {
      when "derived-from(..../area-type,'stub-nssa-area') " {
        description
          "Summary advertisement into the stub/NSSA area.";
      }
      type boolean;
      description
        "Enable/Disable summary advertisement into the stub or
```

```
        NSSA area.";
    }
    leaf default-cost {
        when "derived-from(..../area-type,'stub-nssa-area') " {
            description
                "Cost for LSA default route advertised into the
                stub or NSSA area.";
        }
        type ospf-metric;
        description
            "Set the summary default route cost for a
            stub or NSSA area.";
    }
}

grouping area-config {
    description
        "OSPF area configuration state.";

    leaf area-type {
        type identityref {
            base area-type;
        }
        default normal-area;
        description
            "Area type.";
    }

    uses area-common-config;
    uses address-family-area-config;
}

grouping area-state {
    description
        "OSPF area operational state.";

    container statistics {
        config false;
        description "Per-area statistics";
        uses area-stat;
    }

    container database {
        config false;
        description "Area-scope Link State Database.";
        list area-scope-lsa-type {
            key "lsa-type";
            description "List OSPF area-scope LSAs.";
        }
    }
}
```

```

    leaf lsa-type {
        type uint16;
        description "OSPF area-scope LSA type.";
    }
    container area-scope-lsas {
        description
            "All area-scope LSAs of an area-scope
            LSA type.";
        list area-scope-lsa {
            key "lsa-id adv-router";
            description "List of OSPF area-scope LSAs";
            uses lsa-key;
            uses lsa {
                refine "version/ospfv2/ospfv2" {
                    must "derived-from-or-self( "
                        + "../..../..../..../..../..../"
                        + "rt:type, 'ospfv2') " {
                        description "OSPFv2 LSA.";
                    }
                }
                refine "version/ospfv3/ospfv3" {
                    must "derived-from-or-self( "
                        + "../..../..../..../..../..../"
                        + "rt:type, 'ospfv3') " {
                        description "OSPFv3 LSA.";
                    }
                }
            }
        }
    }
}

grouping local-rib {
    description "Local-rib - RIB for Routes computed by the local
        OSPF routing instance.";
    container local-rib {
        config false;
        description "Local-rib.";
        list route {
            key "prefix";
            description "Routes";
            leaf prefix {
                type inet:ip-prefix;
                description "Destination prefix.";
            }
            container next-hops {

```

```
        description "Next hops for the route.";
        list next-hop {
            key "next-hop";
            description "List of next hops for the route";
            leaf outgoing-interface {
                type if:interface-ref;
                description
                    "Name of the outgoing interface.";
            }
            leaf next-hop {
                type inet:ip-address;
                description "Next hop address.";
            }
        }
    }
    leaf metric {
        type uint32;
        description "Metric for this route.";
    }
    leaf route-type {
        type route-type;
        description "Route type for this route.";
    }
    leaf route-tag {
        type uint32;
        description "Route tag for this route.";
    }
}

grouping ietf-spf-delay {
    leaf initial-delay {
        type uint32;
        units milliseconds;
        description
            "Delay used while in QUIET state (milliseconds).";
    }
    leaf short-delay {
        type uint32;
        units milliseconds;
        description
            "Delay used while in SHORT_WAIT state (milliseconds).";
    }
    leaf long-delay {
        type uint32;
        units milliseconds;
        description
```

```
        "Delay used while in LONG_WAIT state (milliseconds).";
    }
    leaf hold-down {
        type uint32;
        units milliseconds;
        description
            "Timer used to consider an IGP stability period
             (milliseconds).";
    }
    leaf time-to-learn {
        type uint32;
        units milliseconds;
        description
            "Duration used to learn all the IGP events
             related to a single component failure (milliseconds).";
    }
    leaf current-state {
        type enumeration {
            enum "quiet" {
                description "QUIET state";
            }
            enum "short-wait" {
                description "SHORT_WAIT state";
            }
            enum "long-wait" {
                description "LONG_WAIT state";
            }
        }
        config false;
        description
            "Current SPF back-off algorithm state.";
    }
    leaf remaining-time-to-learn {
        type rt-types:timer-value-milliseconds;
        config false;
        description
            "Remaining time until time-to-learn timer fires.";
    }
    leaf remaining-hold-down {
        type rt-types:timer-value-milliseconds;
        config false;
        description
            "Remaining time until hold-down timer fires.";
    }
    leaf last-event-received {
        type yang:timestamp;
        config false;
        description
```

```
        "Time of last SPF triggering event.";
    }
    leaf next-spf-time {
        type yang:timestamp;
        config false;
        description
            "Time when next SPF has been scheduled.";
    }
    leaf last-spf-time {
        type yang:timestamp;
        config false;
        description
            "Time of last SPF computation.";
    }
    description
        "Grouping for IETF SPF delay configuration and state";
}

grouping node-tag-config {
    description
        "OSPF node tag config state.";
    container node-tags {
        if-feature node-tag;
        list node-tag {
            key tag;
            leaf tag {
                type uint32;
                description
                    "Node tag value.";
            }
            description
                "List of tags.";
        }
        description
            "Container for node admin tags.";
    }
}

grouping instance-config {
    description
        "OSPF instance config state.";

    leaf enable {
        type boolean;
        default true;
        description
            "Enable/Disable the protocol.";
    }
}
```

```
leaf explicit-router-id {
  if-feature explicit-router-id;
  type rt-types:router-id;
  description
    "Defined in RFC 2328. A 32-bit number
     that uniquely identifies the router.";
}

container preference {
  description
    "Route preference configuration. In many
     implementations, preference is referred to as
     administrative distance.";
  reference
    "RFC 8349: A YANG Data Model for Routing Management
     (NMDA Version)";
  choice scope {
    description
      "Options for expressing preference
       as single or multiple values.";
    case single-value {
      leaf all {
        type uint8;
        description
          "Preference for intra-area, inter-area, and
           external routes.";
      }
    }
    case multi-values {
      choice granularity {
        description
          "Options for expressing preference
           for intra-area and inter-area routes.";
        case detail {
          leaf intra-area {
            type uint8;
            description
              "Preference for intra-area routes.";
          }
          leaf inter-area {
            type uint8;
            description
              "Preference for inter-area routes.";
          }
        }
        case coarse {
          leaf internal {
            type uint8;
          }
        }
      }
    }
  }
}
```



```
        description
            "Preference for both intra-area and
            inter-area routes.";
    }
}
leaf external {
    type uint8;
    description
        "Preference for AS external routes.";
}
}
}

container nsr {
    if-feature nsr;
    description
        "Non-Stop Routing (NSR) config state.";
    leaf enable {
        type boolean;
        description
            "Enable/Disable NSR.";
    }
}

container graceful-restart {
    if-feature graceful-restart;
    description
        "Graceful restart config state.";
    reference "RFC 3623: OSPF Graceful Restart
        RFC 5187: OSPFv3 Graceful Restart";
    leaf enable {
        type boolean;
        description
            "Enable/Disable graceful restart as defined in RFC 3623
            for OSPFv2 and RFC 5187 for OSPFv3.";
    }
    leaf helper-enable {
        type boolean;
        description
            "Enable graceful restart helper support for restarting
            routers (RFC 3623 Section 3).";
    }
    leaf restart-interval {
        type uint16 {
            range "1..1800";
        }
    }
}
```

```
        units seconds;
        default "120";
        description
            "Interval to attempt graceful restart prior
             to failing (RFC 3623 Section B.1) (seconds)";
    }
    leaf helper-strict-lsa-checking {
        type boolean;
        description
            "Terminate graceful restart when an LSA topology change
             is detected (RFC 3623 Section B.2).";
    }
}

container auto-cost {
    if-feature auto-cost;
    description
        "Interface Auto-cost configuration state.";
    leaf enable {
        type boolean;
        description
            "Enable/Disable interface auto-cost.";
    }
    leaf reference-bandwidth {
        when "../enable = 'true'" {
            description "Only when auto cost is enabled";
        }
        type uint32 {
            range "1..4294967";
        }
        units Mbits;
        description
            "Configure reference bandwidth used to automatically
             determine interface cost (Mbits). The cost is the
             reference bandwidth divided by the interface speed
             with 1 being the minimum cost.";
    }
}

container spf-control {
    leaf paths {
        if-feature max-ecmp;
        type uint16 {
            range "1..65535";
        }
        description
            "Maximum number of Equal-Cost Multi-Path (ECMP) paths.";
    }
}
```

```
    container ietf-spf-delay {
      if-feature ietf-spf-delay;
      uses ietf-spf-delay;
      description
        "IETF SPF delay algorithm configuration.";
    }
    description "SPF calculation control.";
  }

  container database-control {
    leaf max-lsa {
      if-feature max-lsa;
      type uint32 {
        range "1..4294967294";
      }
      description
        "Maximum number of LSAs OSPF the router will accept.";
    }
    description "Database maintenance control.";
  }

  container stub-router {
    if-feature stub-router;
    description "Set maximum metric configuration";

    choice trigger {
      description
        "Specific triggers which will enable stub
        router state.";
      container always {
        presence
          "Enables unconditional stub router support";
        description
          "Unconditional stub router state (advertise
          transit links with MaxLinkMetric";
        reference "RFC 6987: OSPF Stub Router
          Advertisement";
      }
    }
  }

  container mpls {
    description
      "OSPF MPLS config state.";
    container te-rid {
      if-feature te-rid;
      description
        "Stable OSPF Router IP Address used for Traffic
```

```
        Engineering (TE)";
    leaf ipv4-router-id {
        type inet:ipv4-address;
        description
            "Explicitly configure the TE IPv4 Router ID.";
    }
    leaf ipv6-router-id {
        type inet:ipv6-address;
        description
            "Explicitly configure the TE IPv6 Router ID.";
    }
}
container ldp {
    description
        "OSPF MPLS LDP config state.";
    leaf igp-sync {
        if-feature ldp-igp-sync;
        type boolean;
        description
            "Enable LDP IGP synchronization.";
    }
}
}
uses instance-fast-reroute-config;
uses node-tag-config;
}

grouping instance-state {
    description
        "OSPF instance operational state.";

    leaf router-id {
        type rt-types:router-id;
        config false;
        description
            "Defined in RFC 2328. A 32-bit number
             that uniquely identifies the router.";
    }

    uses local-rib;

    container statistics {
        config false;
        description "Per-instance statistics";
        uses instance-stat;
    }

    container database {
```

```

    config false;
    description "AS-scope Link State Database.";
    list as-scope-lsa-type {
        key "lsa-type";
        description "List OSPF AS-scope LSAs.";
        leaf lsa-type {
            type uint16;
            description "OSPF AS scope LSA type.";
        }
        container as-scope-lsas {
            description "All AS-scope of LSA of this LSA type.";
            list as-scope-lsa {
                key "lsa-id adv-router";
                description "List of OSPF AS-scope LSAs";
                uses lsa-key;
                uses lsa {
                    refine "version/ospfv2/ospfv2" {
                        must "derived-from-or-self( "
                            + "../.../.../.../.../"
                            + "rt:type, 'ospfv2') " {
                            description "OSPFv2 LSA.";
                        }
                    }
                    refine "version/ospfv3/ospfv3" {
                        must "derived-from-or-self( "
                            + "../.../.../.../.../"
                            + "rt:type, 'ospfv3') " {
                            description "OSPFv3 LSA.";
                        }
                    }
                }
            }
        }
    }
    uses spf-log;
    uses lsa-log;
}

grouping multi-topology-area-common-config {
    description
        "OSPF multi-topology area common configuration state.";
    leaf summary {
        when "derived-from(.../.../.../area-type, 'stub-nssa-area') " {
            description
                "Summary advertisement into the stub/NSSA area.";
        }
        type boolean;
    }
}

```

```
        description
            "Enable/Disable summary advertisement into the
            topology in the stub or NSSA area.";
    }
    leaf default-cost {
        when "derived-from ../../../../area-type, 'stub-nssa-area'" {
            description
                "Cost for LSA default route advertised into the
                topology into the stub or NSSA area.";
        }
        type ospf-metric;
        description
            "Set the summary default route cost for a
            stub or NSSA area.";
    }
}

grouping multi-topology-area-config {
    description
        "OSPF multi-topology area configuration state.";

    uses multi-topology-area-common-config;
    uses address-family-area-config;
}

grouping multi-topology-state {
    description
        "OSPF multi-topology operational state.";

    uses local-rib;
}

grouping multi-topology-interface-config {
    description
        "OSPF multi-topology configuration state.";

    leaf cost {
        type ospf-link-metric;
        description
            "Interface cost for this topology.";
    }
}

grouping ospfv3-interface-config {
    description
        "OSPFv3 interface specific configuration state.";

    leaf instance-id {
```

```
        type uint8 {
            range "0 .. 31";
        }
        description
            "OSPFv3 instance ID.";
    }
}

grouping ospfv3-interface-state {
    description
        "OSPFv3 interface specific operational state.";

    leaf interface-id {
        type uint16;
        config false;
        description
            "OSPFv3 interface ID.";
    }
}

grouping lsa-identifiers {
    description
        "The parameters that uniquely identify an LSA.";
    leaf area-id {
        type area-id-type;
        description
            "Area ID";
    }
    leaf type {
        type uint16;
        description
            "LSA type.";
    }
    leaf lsa-id {
        type union {
            type inet:ipv4-address;
            type yang:dotted-quad;
        }
        description "Link-State ID.";
    }
    leaf adv-router {
        type rt-types:router-id;
        description
            "LSA advertising router.";
    }
    leaf seq-num {
        type uint32;
        description
```

```
        "LSA sequence number.";
    }
}

grouping spf-log {
    description
        "Grouping for SPF log.";
    container spf-log {
        config false;
        description
            "This container lists the SPF log.";
        list event {
            key id;
            description
                "List of SPF log entries represented
                 as a wrapping buffer in chronological
                 order with the oldest entry returned
                 first.";
            leaf id {
                type uint32;
                description
                    "Event identifier - Purely internal value.";
            }
            leaf spf-type {
                type enumeration {
                    enum full {
                        description
                            "SPF computation was a Full SPF.";
                    }
                    enum intra {
                        description
                            "SPF computation was only for intra-area routes.";
                    }
                    enum inter {
                        description
                            "SPF computation was only for inter-area
                             summary routes.";
                    }
                    enum external {
                        description
                            "SPF computation was only for AS external routes.";
                    }
                }
            }
            description
                "The SPF computation type for the SPF log entry.";
        }
        leaf schedule-timestamp {
            type yang:timestamp;
        }
    }
}
```



```
        description
            "This is the timestamp when the computation was
            scheduled.";
    }
    leaf start-timestamp {
        type yang:timestamp;
        description
            "This is the timestamp when the computation was
            started.";
    }
    leaf end-timestamp {
        type yang:timestamp;
        description
            "This the timestamp when the computation was
            completed.";
    }
    list trigger-lsa {
        description
            "The list of LSAs that triggered the computation.";
        uses lsa-identifiers;
    }
}

}

}

grouping lsa-log {
    description
        "Grouping for the LSA log.";
    container lsa-log {
        config false;
        description
            "This container lists the LSA log.
            Local LSA modifications are also included
            in the list.";
        list event {
            key id;
            description
                "List of LSA log entries represented
                as a wrapping buffer in chronological order
                with the oldest entries returned first.";
            leaf id {
                type uint32;
                description
                    "Event identifier - purely internal value.";
            }
        }
        container lsa {
            description
                "This container describes the logged LSA.";
        }
    }
}
```

```
        uses lsa-identifiers;
    }
    leaf received-timestamp {
        type yang:timestamp;
        description
            "This is the timestamp when the LSA was received.
            In case of local LSA update, the timestamp refers
            to the LSA origination time.";
    }
    leaf reason {
        type identityref {
            base lsa-log-reason;
        }
        description
            "This reason for the LSA log entry.";
    }
}
}
}

augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol" {
    when "derived-from(rt:type, 'ospf')" {
        description
            "This augmentation is only valid for a routing protocol
            instance of OSPF (type 'ospfv2' or 'ospfv3').";
    }
    description "OSPF protocol ietf-routing module
        control-plane-protocol augmentation.";

    container ospf {
        description
            "OSPF protocol Instance";

        leaf address-family {
            type iana-rt-types:address-family;
            description
                "Address-family of the instance.";
        }

        uses instance-config;
        uses instance-state;

        container areas {
            description "All areas.";
            list area {
                key "area-id";
                description

```

```
    "List of OSPF areas";
  leaf area-id {
    type area-id-type;
    description
      "Area ID";
  }

  uses area-config;
  uses area-state;

  container virtual-links {
    when "derived-from-or-self(..../area-type, 'normal-area') "
      + "and ..../area-id = '0.0.0.0'" {
      description
        "Virtual links must be in backbone area.";
    }
    description "All virtual links.";
    list virtual-link {
      key "transit-area-id router-id";
      description
        "OSPF virtual link";
      leaf transit-area-id {
        type leafref {
          path "../..../..../area/area-id";
        }
        must "derived-from-or-self("
          + "../..../..../area[area-id=current()]/area-type, "
          + "'normal-area') and "
          + "../..../..../area[area-id=current()]/area-id != "
          + "'0.0.0.0'" {
          error-message "Virtual link transit area must "
            + "be non-zero.";
          description
            "Virtual-link transit area must be
              non-zero area.";
        }
        description
          "Virtual link transit area ID.";
      }
      leaf router-id {
        type rt-types:router-id;
        description
          "Virtual Link remote endpoint Router ID.";
      }
    }

    uses virtual-link-config;
    uses virtual-link-state;
  }
```

```

    }
    container sham-links {
      if-feature pe-ce-protocol;
      description "All sham links.";
      list sham-link {
        key "local-id remote-id";
        description
          "OSPF sham link";
        leaf local-id {
          type inet:ip-address;
          description
            "Address of the local sham Link endpoint.";
        }
        leaf remote-id {
          type inet:ip-address;
          description
            "Address of the remote sham Link endpoint.";
        }
        uses sham-link-config;
        uses sham-link-state;
      }
    }
    container interfaces {
      description "All interfaces.";
      list interface {
        key "name";
        description
          "List of OSPF interfaces.";
        leaf name {
          type if:interface-ref;
          description
            "Interface name reference.";
        }
        uses interface-config;
        uses interface-state;
      }
    }
  }
}

augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol/ospf" {
  when "derived-from(../rt:type, 'ospf')" {
    description
      "This augmentation is only valid for OSPF
      (type 'ospfv2' or 'ospfv3').";
  }
}

```

```

    }
    if-feature multi-topology;
    description
        "OSPF multi-topology instance configuration
        state augmentation.";
    container topologies {
        description "All topologies.";
        list topology {
            key "name";
            description
                "OSPF topology - The OSPF topology address-family
                must coincide with the routing-instance
                address-family.";
            leaf name {
                type leafref {
                    path "../.../.../.../rt:ribs/rt:rib/rt:name";
                }
                description "RIB name corresponding to the OSPF
                    topology.";
            }

            uses multi-topology-state;
        }
    }
}

augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol/ospf/"
+ "areas/area" {
    when "derived-from-or-self(.../.../.../rt:type, "
    + "'ospfv2') " {
        description
            "This augmentation is only valid for OSPFv2.";
    }
    if-feature multi-topology;
    description
        "OSPF multi-topology area configuration state
        augmentation.";
    container topologies {
        description "All topologies for the area.";
        list topology {
            key "name";
            description "OSPF area topology.";
            leaf name {
                type leafref {
                    path "../.../.../.../.../.../.../..."
                    + "rt:ribs/rt:rib/rt:name";
                }
            }
        }
    }
}

```

```
        description
            "Single topology enabled for this area.";
    }

    uses multi-topology-area-config;
}

}

augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol/ospf/"
+ "areas/area/interfaces/interface" {
    when "derived-from-or-self ../../../../rt:type, "
    + "'ospfv2'" {
        description
            "This augmentation is only valid for OSPFv2.";
    }
    if-feature multi-topology;
    description
        "OSPF multi-topology interface configuration state
        augmentation.";
    container topologies {
        description "All topologies for the interface.";
        list topology {
            key "name";
            description "OSPF interface topology.";
            leaf name {
                type leafref {
                    path "../../../../../rt:ribs/rt:rib/rt:name";
                }
            }
            description
                "Single topology enabled on this interface.";
        }

        uses multi-topology-interface-config;
    }
}

}

augment "/rt:routing/rt:control-plane-protocols/"
+ "rt:control-plane-protocol/ospf/"
+ "areas/area/interfaces/interface" {
    when "derived-from-or-self ../../../../rt:type, "
    + "'ospfv3'" {
        description
            "This augmentation is only valid for OSPFv3.";
    }
}
```

```
    description
      "OSPFv3 interface specific configuration state
      augmentation.";
    uses ospfv3-interface-config;
    uses ospfv3-interface-state;
  }

  grouping route-content {
    description
      "This grouping defines OSPF-specific route attributes.";
    leaf metric {
      type uint32;
      description "OSPF route metric.";
    }
    leaf tag {
      type uint32;
      default "0";
      description "OSPF route tag.";
    }
    leaf route-type {
      type route-type;
      description "OSPF route type";
    }
  }

  augment "/rt:routing/rt:ribs/rt:rib/rt:routes/rt:route" {
    when "derived-from(rt:source-protocol, 'ospf')";
    description
      "This augmentation is only valid for routes whose
      source protocol is OSPF.";
  }
  description
    "OSPF-specific route attributes.";
  uses route-content;
}

/*
 * RPCs
 */

rpc clear-neighbor {
  description
    "This RPC request clears a particular set of OSPF neighbors.
    If the operation fails for OSPF internal reason, then
    error-tag and error-app-tag should be set to a meaningful
    value.";
  input {
    leaf routing-protocol-name {
```

```
    type leafref {
      path "/rt:routing/rt:control-plane-protocols/"
        + "rt:control-plane-protocol/rt:name";
    }
    mandatory "true";
    description
      "OSPF protocol instance which information for neighbors
      are to be cleared.

      If the referenced OSPF instance doesn't exist, then
      this operation SHALL fail with error-tag 'data-missing'
      and error-app-tag
      'routing-protocol-instance-not-found'.";
  }

  leaf interface {
    type if:interface-ref;
    description
      "Name of the OSPF interface for which neighbors are to
      be cleared.

      If the referenced OSPF interface doesn't exist, then
      this operation SHALL fail with error-tag
      'data-missing' and error-app-tag
      'ospf-interface-not-found'.";
  }
}

rpc clear-database {
  description
    "This RPC request clears a particular OSPF Link State
    Database. If the operation fails for OSPF internal reason,
    then error-tag and error-app-tag should be set to a
    meaningful value.";
  input {
    leaf routing-protocol-name {
      type leafref {
        path "/rt:routing/rt:control-plane-protocols/"
          + "rt:control-plane-protocol/rt:name";
      }
      mandatory "true";
      description
        "OSPF protocol instance whose Link State Database is to
        be cleared.

        If the referenced OSPF instance doesn't exist, then
        this operation SHALL fail with error-tag 'data-missing'";
    }
  }
}
```



```
        and error-app-tag
        'routing-protocol-instance-not-found'. ";
    }
}

/*
 * Notifications
 */

grouping notification-instance-hdr {
  description
    "This grouping describes common instance specific
    data for OSPF notifications.";

  leaf routing-protocol-name {
    type leafref {
      path "/rt:routing/rt:control-plane-protocols/"
        + "rt:control-plane-protocol/rt:name";
    }
    must "derived-from( "
      + "/rt:routing/rt:control-plane-protocols/"
      + "rt:control-plane-protocol[rt:name=current()]/"
      + "rt:type, 'ospf')";
    description
      "OSPF routing protocol instance name.";
  }

  leaf address-family {
    type leafref {
      path "/rt:routing/"
        + "rt:control-plane-protocols/rt:control-plane-protocol"
        + "[rt:name=current()]/../routing-protocol-name]/"
        + "ospf/address-family";
    }
    description
      "Address family of the OSPF instance.";
  }
}

grouping notification-interface {
  description
    "This grouping provides interface information
    for the OSPF interface specific notification.";

  choice if-link-type-selection {
    description
      "Options for link type.";
  }
}
```

```
    container interface {
      description "Normal interface.";
      leaf interface {
        type if:interface-ref;
        description "Interface.";
      }
    }
    container virtual-link {
      description "virtual-link.";
      leaf transit-area-id {
        type area-id-type;
        description "Area ID.";
      }
      leaf neighbor-router-id {
        type rt-types:router-id;
        description "Neighbor Router ID.";
      }
    }
    container sham-link {
      description "sham link.";
      leaf area-id {
        type area-id-type;
        description "Area ID.";
      }
      leaf local-ip-addr {
        type inet:ip-address;
        description "Sham link local address.";
      }
      leaf remote-ip-addr {
        type inet:ip-address;
        description "Sham link remote address.";
      }
    }
  }
}

grouping notification-neighbor {
  description
    "This grouping provides the neighbor information
    for neighbor specific notifications.";

  leaf neighbor-router-id {
    type rt-types:router-id;
    description "Neighbor Router ID.";
  }

  leaf neighbor-ip-addr {
    type inet:ip-address;
  }
}
```

```
        description "Neighbor address.";
    }
}

notification if-state-change {
    uses notification-instance-hdr;
    uses notification-interface;

    leaf state {
        type if-state-type;
        description "Interface state.";
    }
    description
        "This notification is sent when an interface
        state change is detected.";
}

notification if-config-error {
    uses notification-instance-hdr;
    uses notification-interface;

    leaf packet-source {
        type inet:ip-address;
        description "Source address.";
    }

    leaf packet-type {
        type packet-type;
        description "OSPF packet type.";
    }

    leaf error {
        type enumeration {
            enum "bad-version" {
                description "Bad version.";
            }
            enum "area-mismatch" {
                description "Area mismatch.";
            }
            enum "unknown-nbma-nbr" {
                description "Unknown NBMA neighbor.";
            }
            enum "unknown-virtual-nbr" {
                description "Unknown virtual link neighbor.";
            }
            enum "auth-type-mismatch" {
                description "Auth type mismatch.";
            }
        }
    }
}
```

```
    enum "auth-failure" {
      description "Auth failure.";
    }
    enum "net-mask-mismatch" {
      description "Network mask mismatch.";
    }
    enum "hello-interval-mismatch" {
      description "Hello interval mismatch.";
    }
    enum "dead-interval-mismatch" {
      description "Dead interval mismatch.";
    }
    enum "option-mismatch" {
      description "Option mismatch.";
    }
    enum "mtu-mismatch" {
      description "MTU mismatch.";
    }
    enum "duplicate-router-id" {
      description "Duplicate Router ID.";
    }
    enum "no-error" {
      description "No error.";
    }
  }
  description "Error code.";
}
description
  "This notification is sent when an interface
  config error is detected.";
}

notification nbr-state-change {
  uses notification-instance-hdr;
  uses notification-interface;
  uses notification-neighbor;

  leaf state {
    type nbr-state-type;
    description "Neighbor state.";
  }

  description
    "This notification is sent when a neighbor
    state change is detected.";
}

notification nbr-restart-helper-status-change {
```

```
    uses notification-instance-hdr;
    uses notification-interface;
    uses notification-neighbor;

    leaf status {
        type restart-helper-status-type;
        description "Restart helper status.";
    }

    leaf age {
        type rt-types:timer-value-seconds16;
        description
            "Remaining time in current OSPF graceful restart
            interval when the router is acting as a restart
            helper for the neighbor.";
    }

    leaf exit-reason {
        type restart-exit-reason-type;
        description
            "Restart helper exit reason.";
    }
    description
        "This notification is sent when a neighbor restart
        helper status change is detected.";
}

notification if-rx-bad-packet {
    uses notification-instance-hdr;
    uses notification-interface;

    leaf packet-source {
        type inet:ip-address;
        description "Source address.";
    }

    leaf packet-type {
        type packet-type;
        description "OSPF packet type.";
    }

    description
        "This notification is sent when an OSPF packet that
        cannot be parsed is received on an OSPF interface.";
}

notification lsdb-approaching-overflow {
    uses notification-instance-hdr;
```

```
    leaf ext-lsdb-limit {
      type uint32;
      description
        "The maximum number of non-default AS-external LSAs
        entries that can be stored in the Link State Database.";
    }

    description
      "This notification is sent when the number of LSAs
      in the router's Link State Database has exceeded
      ninety percent of the AS-external limit (ext-lsdb-limit).";
  }

  notification lsdb-overflow {
    uses notification-instance-hdr;

    leaf ext-lsdb-limit {
      type uint32;
      description
        "The maximum number of non-default AS-external LSAs
        entries that can be stored in the Link State Database.";
    }

    description
      "This notification is sent when the number of LSAs
      in the router's Link State Database has exceeded the
      AS-external limit (ext-lsdb-limit).";
  }

  notification nssa-translator-status-change {
    uses notification-instance-hdr;

    leaf area-id {
      type area-id-type;
      description "Area ID.";
    }

    leaf status {
      type nssa-translator-state-type;
      description
        "NSSA translator status.";
    }

    description
      "This notification is sent when there is a change
      in the router's role in translating OSPF NSSA LSAs
      to OSPF AS-External LSAs.";
  }
```

```
notification restart-status-change {
  uses notification-instance-hdr;

  leaf status {
    type restart-status-type;
    description
      "Restart status.";
  }

  leaf restart-interval {
    type uint16 {
      range 1..1800;
    }
    units seconds;
    default "120";
    description
      "Restart interval.";
  }

  leaf exit-reason {
    type restart-exit-reason-type;
    description
      "Restart exit reason.";
  }

  description
    "This notification is sent when the graceful restart
     state for the router has changed.";
}
}
<CODE ENDS>
```

4. Security Considerations

The YANG modules specified in this document define a schema for data that is designed to be accessed via network management protocols such as NETCONF [RFC6241] or RESTCONF [RFC8040]. The lowest NETCONF layer is the secure transport layer, and the mandatory-to-implement secure transport is Secure Shell (SSH) [RFC6242]. The lowest RESTCONF layer is HTTPS, and the mandatory-to-implement secure transport is TLS [RFC8446].

The NETCONF Access Control Model (NACM) [RFC8341] provides the means to restrict access for particular NETCONF or RESTCONF users to a pre-configured subset of all available NETCONF or RESTCONF protocol operations and content.

There are a number of data nodes defined in `ietf-ospf.yang` module that are writable/creatable/deletable (i.e., `config true`, which is the default). These data nodes may be considered sensitive or vulnerable in some network environments. Write operations (e.g., `edit-config`) to these data nodes without proper protection can have a negative effect on network operations. Writable data node represent configuration of each instance, area, virtual link, sham-link, and interface. These correspond to the following schema nodes:

```
/ospf
/ospf/areas/
/ospf/areas/area[area-id]
/ospf/virtual-links/
/ospf/virtual-links/virtual-link[transit-area-id router-id]
/ospf/areas/area[area-id]/interfaces
/ospf/areas/area[area-id]/interfaces/interface[name]
/ospf/area/area[area-id]/sham-links
/ospf/area/area[area-id]/sham-links/sham-link[local-id remote-id]
```

For OSPF, the ability to modify OSPF configuration will allow the entire OSPF domain to be compromised including peering with unauthorized routers to misroute traffic or mount a massive Denial-of-Service (DoS) attack. For example, adding OSPF on any unprotected interface could allow an OSPF adjacency to be formed with an unauthorized and malicious neighbor. Once an adjacency is formed, traffic could be hijacked. As a simpler example, a Denial-of-Service attack could be mounted by changing the cost of an OSPF interface to be asymmetric such that a hard routing loop ensues. In general, unauthorized modification of most OSPF features will pose there own set of security risks and the "Security Considerations" in the respective reference RFCs should be consulted.

Some of the readable data nodes in the `ietf-ospf.yang` module may be considered sensitive or vulnerable in some network environments. It is thus important to control read access (e.g., via `get`, `get-config`, or `notification`) to these data nodes. The exposure of the Link State Database (LSDB) will expose the detailed topology of the network. There is a separate Link State Database for each instance, area, virtual link, sham-link, and interface. These correspond to the following schema nodes:


```
/ospf/database  
  
/ospf/areas/area[area-id]/database  
  
/ospf/virtual-links/virtual-link[transit-area-id router-  
id]/database  
  
/ospf/areas/area[area-id]/interfaces/interface[name]/database  
  
/ospf/area/area[area-id]/sham-links/sham-link[local-id remote-  
id]/database
```

Exposure of the Link State Database includes information beyond the scope of the OSPF router and this may be undesirable since exposure may facilitate other attacks. Additionally, in the case of an area LSDB, the complete IP network topology and, if deployed, the traffic engineering topology of the OSPF area can be reconstructed. Network operators may consider their topologies to be sensitive confidential data.

For OSPF authentication, configuration is supported via the specification of key-chains [RFC8177] or the direct specification of key and authentication algorithm. Hence, authentication configuration using the "auth-table-trailer" case in the "authentication" container inherits the security considerations of [RFC8177]. This includes the considerations with respect to the local storage and handling of authentication keys.

Additionally, local specification of OSPF authentication keys and the associated authentication algorithm is supported for legacy implementations that do not support key-chains [RFC8177]. It is RECOMMENDED that implementations migrate to key-chains due the seamless support of key and algorithm rollover, as well as, the hexadecimal key specification affording more key entropy, and encryption of keys using the Advanced Encryption Standard (AES) Key Wrap Padding Algorithm [RFC5649].

Some of the RPC operations in this YANG module may be considered sensitive or vulnerable in some network environments. It is thus important to control access to these operations. The OSPF YANG module supports the "clear-neighbor" and "clear-database" RPCs. If access to either of these is compromised, they can result in temporary network outages be employed to mount DoS attacks.

The actual authentication key data (whether locally specified or part of a key-chain) is sensitive and needs to be kept secret from unauthorized parties; compromise of the key data would allow an

attacker to forge OSPF traffic that would be accepted as authentic, potentially compromising the entirety OSPF domain.

5. IANA Considerations

This document registers a URI in the IETF XML registry [RFC3688]. Following the format in [RFC3688], the following registration is requested to be made:

```
URI: urn:ietf:params:xml:ns:yang:ietf-ospf
Registrant Contact: The IESG.
XML: N/A, the requested URI is an XML namespace.
```

This document registers a YANG module in the YANG Module Names registry [RFC6020].

```
name: ietf-ospf
namespace: urn:ietf:params:xml:ns:yang:ietf-ospf
prefix: ospf
reference: RFC XXXX
```

6. Acknowledgements

The authors wish to thank Yi Yang, Alexander Clemm, Gaurav Gupta, Ladislav Lhotka, Stephane Litkowski, Greg Hankins, Manish Gupta, Michael Darwish, and Alan Davey for their thorough reviews and helpful comments.

Thanks to Tom Petch for last call review and improvement of the document organization.

Thanks to Alvaro Retana for AD comments.

Thanks to Benjamin Kaduk, Suresh Krishnan, and Roman Danyliw for IESG review comments.

This document was produced using Marshall Rose's xml2rfc tool.

Author affiliation with The MITRE Corporation is provided for identification purposes only, and is not intended to convey or imply MITRE's concurrence with, or support for, the positions, opinions or viewpoints expressed. MITRE has approved this document for Public Release, Distribution Unlimited, with Public Release Case Number 18-3194.

7. References

7.1. Normative References

- [I-D.ietf-bfd-yang]
Rahman, R., Zheng, L., Jethanandani, M., Networks, J., and G. Mirsky, "YANG Data Model for Bidirectional Forwarding Detection (BFD)", draft-ietf-bfd-yang-17 (work in progress), August 2018.
- [RFC1765] Moy, J., "OSPF Database Overflow", RFC 1765, DOI 10.17487/RFC1765, March 1995, <<https://www.rfc-editor.org/info/rfc1765>>.
- [RFC1793] Moy, J., "Extending OSPF to Support Demand Circuits", RFC 1793, DOI 10.17487/RFC1793, April 1995, <<https://www.rfc-editor.org/info/rfc1793>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC3101] Murphy, P., "The OSPF Not-So-Stubby Area (NSSA) Option", RFC 3101, DOI 10.17487/RFC3101, January 2003, <<https://www.rfc-editor.org/info/rfc3101>>.
- [RFC3623] Moy, J., Pillay-Esnault, P., and A. Lindem, "Graceful OSPF Restart", RFC 3623, DOI 10.17487/RFC3623, November 2003, <<https://www.rfc-editor.org/info/rfc3623>>.
- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", RFC 3630, DOI 10.17487/RFC3630, September 2003, <<https://www.rfc-editor.org/info/rfc3630>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC4552] Gupta, M. and N. Melam, "Authentication/Confidentiality for OSPFv3", RFC 4552, DOI 10.17487/RFC4552, June 2006, <<https://www.rfc-editor.org/info/rfc4552>>.

- [RFC4576] Rosen, E., Psenak, P., and P. Pillay-Esnault, "Using a Link State Advertisement (LSA) Options Bit to Prevent Looping in BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4576, DOI 10.17487/RFC4576, June 2006, <<https://www.rfc-editor.org/info/rfc4576>>.
- [RFC4577] Rosen, E., Psenak, P., and P. Pillay-Esnault, "OSPF as the Provider/Customer Edge Protocol for BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4577, DOI 10.17487/RFC4577, June 2006, <<https://www.rfc-editor.org/info/rfc4577>>.
- [RFC4915] Psenak, P., Mirtorabi, S., Roy, A., Nguyen, L., and P. Pillay-Esnault, "Multi-Topology (MT) Routing in OSPF", RFC 4915, DOI 10.17487/RFC4915, June 2007, <<https://www.rfc-editor.org/info/rfc4915>>.
- [RFC4973] Srisuresh, P. and P. Joseph, "OSPF-xTE: Experimental Extension to OSPF for Traffic Engineering", RFC 4973, DOI 10.17487/RFC4973, July 2007, <<https://www.rfc-editor.org/info/rfc4973>>.
- [RFC5082] Gill, V., Heasley, J., Meyer, D., Savola, P., Ed., and C. Pignataro, "The Generalized TTL Security Mechanism (GTSM)", RFC 5082, DOI 10.17487/RFC5082, October 2007, <<https://www.rfc-editor.org/info/rfc5082>>.
- [RFC5185] Mirtorabi, S., Psenak, P., Lindem, A., Ed., and A. Oswal, "OSPF Multi-Area Adjacency", RFC 5185, DOI 10.17487/RFC5185, May 2008, <<https://www.rfc-editor.org/info/rfc5185>>.
- [RFC5187] Pillay-Esnault, P. and A. Lindem, "OSPFv3 Graceful Restart", RFC 5187, DOI 10.17487/RFC5187, June 2008, <<https://www.rfc-editor.org/info/rfc5187>>.
- [RFC5250] Berger, L., Bryskin, I., Zinin, A., and R. Coltun, "The OSPF Opaque LSA Option", RFC 5250, DOI 10.17487/RFC5250, July 2008, <<https://www.rfc-editor.org/info/rfc5250>>.
- [RFC5286] Atlas, A., Ed. and A. Zinin, Ed., "Basic Specification for IP Fast Reroute: Loop-Free Alternates", RFC 5286, DOI 10.17487/RFC5286, September 2008, <<https://www.rfc-editor.org/info/rfc5286>>.
- [RFC5309] Shen, N., Ed. and A. Zinin, Ed., "Point-to-Point Operation over LAN in Link State Routing Protocols", RFC 5309, DOI 10.17487/RFC5309, October 2008, <<https://www.rfc-editor.org/info/rfc5309>>.

- [RFC5329] Ishiguro, K., Manral, V., Davey, A., and A. Lindem, Ed., "Traffic Engineering Extensions to OSPF Version 3", RFC 5329, DOI 10.17487/RFC5329, September 2008, <<https://www.rfc-editor.org/info/rfc5329>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC5613] Zinin, A., Roy, A., Nguyen, L., Friedman, B., and D. Yeung, "OSPF Link-Local Signaling", RFC 5613, DOI 10.17487/RFC5613, August 2009, <<https://www.rfc-editor.org/info/rfc5613>>.
- [RFC5642] Venkata, S., Harwani, S., Pignataro, C., and D. McPherson, "Dynamic Hostname Exchange Mechanism for OSPF", RFC 5642, DOI 10.17487/RFC5642, August 2009, <<https://www.rfc-editor.org/info/rfc5642>>.
- [RFC5709] Bhatia, M., Manral, V., Fanto, M., White, R., Barnes, M., Li, T., and R. Atkinson, "OSPFv2 HMAC-SHA Cryptographic Authentication", RFC 5709, DOI 10.17487/RFC5709, October 2009, <<https://www.rfc-editor.org/info/rfc5709>>.
- [RFC5714] Shand, M. and S. Bryant, "IP Fast Reroute Framework", RFC 5714, DOI 10.17487/RFC5714, January 2010, <<https://www.rfc-editor.org/info/rfc5714>>.
- [RFC5838] Lindem, A., Ed., Mirtorabi, S., Roy, A., Barnes, M., and R. Aggarwal, "Support of Address Families in OSPFv3", RFC 5838, DOI 10.17487/RFC5838, April 2010, <<https://www.rfc-editor.org/info/rfc5838>>.
- [RFC6020] Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.
- [RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/info/rfc6241>>.
- [RFC6242] Wasserman, M., "Using the NETCONF Protocol over Secure Shell (SSH)", RFC 6242, DOI 10.17487/RFC6242, June 2011, <<https://www.rfc-editor.org/info/rfc6242>>.

- [RFC6565] Pillay-Esnault, P., Moyer, P., Doyle, J., Ertekin, E., and M. Lundberg, "OSPFv3 as a Provider Edge to Customer Edge (PE-CE) Routing Protocol", RFC 6565, DOI 10.17487/RFC6565, June 2012, <<https://www.rfc-editor.org/info/rfc6565>>.
- [RFC6845] Sheth, N., Wang, L., and J. Zhang, "OSPF Hybrid Broadcast and Point-to-Multipoint Interface Type", RFC 6845, DOI 10.17487/RFC6845, January 2013, <<https://www.rfc-editor.org/info/rfc6845>>.
- [RFC6860] Yang, Y., Retana, A., and A. Roy, "Hiding Transit-Only Networks in OSPF", RFC 6860, DOI 10.17487/RFC6860, January 2013, <<https://www.rfc-editor.org/info/rfc6860>>.
- [RFC6987] Retana, A., Nguyen, L., Zinin, A., White, R., and D. McPherson, "OSPF Stub Router Advertisement", RFC 6987, DOI 10.17487/RFC6987, September 2013, <<https://www.rfc-editor.org/info/rfc6987>>.
- [RFC6991] Schoenwaelder, J., Ed., "Common YANG Data Types", RFC 6991, DOI 10.17487/RFC6991, July 2013, <<https://www.rfc-editor.org/info/rfc6991>>.
- [RFC7166] Bhatia, M., Manral, V., and A. Lindem, "Supporting Authentication Trailer for OSPFv3", RFC 7166, DOI 10.17487/RFC7166, March 2014, <<https://www.rfc-editor.org/info/rfc7166>>.
- [RFC7474] Bhatia, M., Hartman, S., Zhang, D., and A. Lindem, Ed., "Security Extension for OSPFv2 When Using Manual Key Management", RFC 7474, DOI 10.17487/RFC7474, April 2015, <<https://www.rfc-editor.org/info/rfc7474>>.
- [RFC7490] Bryant, S., Filsfils, C., Previdi, S., Shand, M., and N. So, "Remote Loop-Free Alternate (LFA) Fast Reroute (FRR)", RFC 7490, DOI 10.17487/RFC7490, April 2015, <<https://www.rfc-editor.org/info/rfc7490>>.
- [RFC7684] Psenak, P., Gredler, H., Shakir, R., Henderickx, W., Tantsura, J., and A. Lindem, "OSPFv2 Prefix/Link Attribute Advertisement", RFC 7684, DOI 10.17487/RFC7684, November 2015, <<https://www.rfc-editor.org/info/rfc7684>>.
- [RFC7770] Lindem, A., Ed., Shen, N., Vasseur, JP., Aggarwal, R., and S. Shaffer, "Extensions to OSPF for Advertising Optional Router Capabilities", RFC 7770, DOI 10.17487/RFC7770, February 2016, <<https://www.rfc-editor.org/info/rfc7770>>.

- [RFC7777] Hegde, S., Shakir, R., Smirnov, A., Li, Z., and B. Decraene, "Advertising Node Administrative Tags in OSPF", RFC 7777, DOI 10.17487/RFC7777, March 2016, <<https://www.rfc-editor.org/info/rfc7777>>.
- [RFC7884] Pignataro, C., Bhatia, M., Aldrin, S., and T. Ranganath, "OSPF Extensions to Advertise Seamless Bidirectional Forwarding Detection (S-BFD) Target Discriminators", RFC 7884, DOI 10.17487/RFC7884, July 2016, <<https://www.rfc-editor.org/info/rfc7884>>.
- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC8040] Bierman, A., Bjorklund, M., and K. Watsen, "RESTCONF Protocol", RFC 8040, DOI 10.17487/RFC8040, January 2017, <<https://www.rfc-editor.org/info/rfc8040>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8177] Lindem, A., Ed., Qu, Y., Yeung, D., Chen, I., and J. Zhang, "YANG Data Model for Key Chains", RFC 8177, DOI 10.17487/RFC8177, June 2017, <<https://www.rfc-editor.org/info/rfc8177>>.
- [RFC8294] Liu, X., Qu, Y., Lindem, A., Hopps, C., and L. Berger, "Common YANG Data Types for the Routing Area", RFC 8294, DOI 10.17487/RFC8294, December 2017, <<https://www.rfc-editor.org/info/rfc8294>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.
- [RFC8341] Bierman, A. and M. Bjorklund, "Network Configuration Access Control Model", STD 91, RFC 8341, DOI 10.17487/RFC8341, March 2018, <<https://www.rfc-editor.org/info/rfc8341>>.
- [RFC8342] Bjorklund, M., Schoenwaelder, J., Shafer, P., Watsen, K., and R. Wilton, "Network Management Datastore Architecture (NMDA)", RFC 8342, DOI 10.17487/RFC8342, March 2018, <<https://www.rfc-editor.org/info/rfc8342>>.

- [RFC8343] Bjorklund, M., "A YANG Data Model for Interface Management", RFC 8343, DOI 10.17487/RFC8343, March 2018, <<https://www.rfc-editor.org/info/rfc8343>>.
- [RFC8349] Lhotka, L., Lindem, A., and Y. Qu, "A YANG Data Model for Routing Management (NMDA Version)", RFC 8349, DOI 10.17487/RFC8349, March 2018, <<https://www.rfc-editor.org/info/rfc8349>>.
- [RFC8405] Decraene, B., Litkowski, S., Gredler, H., Lindem, A., Francois, P., and C. Bowers, "Shortest Path First (SPF) Back-Off Delay Algorithm for Link-State IGP", RFC 8405, DOI 10.17487/RFC8405, June 2018, <<https://www.rfc-editor.org/info/rfc8405>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.
- [RFC8476] Tantsura, J., Chunduri, U., Aldrin, S., and P. Psenak, "Signaling Maximum SID Depth (MSD) Using OSPF", RFC 8476, DOI 10.17487/RFC8476, December 2018, <<https://www.rfc-editor.org/info/rfc8476>>.

7.2. Informative References

- [RFC0905] "ISO Transport Protocol specification ISO DP 8073", RFC 905, DOI 10.17487/RFC0905, April 1984, <<https://www.rfc-editor.org/info/rfc905>>.
- [RFC4750] Joyal, D., Ed., Galecki, P., Ed., Giacalone, S., Ed., Coltun, R., and F. Baker, "OSPF Version 2 Management Information Base", RFC 4750, DOI 10.17487/RFC4750, December 2006, <<https://www.rfc-editor.org/info/rfc4750>>.
- [RFC5443] Jork, M., Atlas, A., and L. Fang, "LDP IGP Synchronization", RFC 5443, DOI 10.17487/RFC5443, March 2009, <<https://www.rfc-editor.org/info/rfc5443>>.
- [RFC5643] Joyal, D., Ed. and V. Manral, Ed., "Management Information Base for OSPFv3", RFC 5643, DOI 10.17487/RFC5643, August 2009, <<https://www.rfc-editor.org/info/rfc5643>>.
- [RFC5649] Housley, R. and M. Dworkin, "Advanced Encryption Standard (AES) Key Wrap with Padding Algorithm", RFC 5649, DOI 10.17487/RFC5649, September 2009, <<https://www.rfc-editor.org/info/rfc5649>>.

- [RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", RFC 5880, DOI 10.17487/RFC5880, June 2010, <<https://www.rfc-editor.org/info/rfc5880>>.
- [RFC5881] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD) for IPv4 and IPv6 (Single Hop)", RFC 5881, DOI 10.17487/RFC5881, June 2010, <<https://www.rfc-editor.org/info/rfc5881>>.

Appendix A. Contributors' Addresses

Dean Bogdanovic
Volta Networks, Inc.

EMail: dean@voltanet.io

Kiran Koushik Agrahara Sreenivasa
Verizon
500 W Dove Rd
Southlake, TX 76092
USA

EMail: kk@employees.org

Authors' Addresses

Derek Yeung
Arrcus

EMail: derek@arrcus.com

Yingzhen Qu
Futurewei
2330 Central Expressway
Santa Clara, CA 95050
USA

EMail: yingzhen.qu@futurewei.com

Jeffrey Zhang
Juniper Networks
10 Technology Park Drive
Westford, MA 01886
USA

EMail: zzhang@juniper.net

Ing-Wher Chen
The MITRE Corporation

EMail: ingwherchen@mitre.org

Acee Lindem
Cisco Systems
301 Midenhall Way
Cary, NC 27513

EMail: acee@cisco.com

Internet Engineering Task Force
Internet-Draft
Intended status: Informational
Expires: December 30, 2018

T. Li
Arista Networks
P. Psenak
Cisco Systems, Inc.
June 28, 2018

Dynamic Flooding on Dense Graphs
draft-li-dynamic-flooding-05

Abstract

Routing with link state protocols in dense network topologies can result in sub-optimal convergence times due to the overhead associated with flooding. This can be addressed by decreasing the flooding topology so that it is less dense.

This document discusses the problem in some depth and an architectural solution. Specific protocol changes for IS-IS, OSPFv2, and OSPFv3 are described in this document.

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 <https://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 December 30, 2018.

Copyright Notice

Copyright (c) 2018 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 (<https://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. Requirements Language	4
2. Problem Statement	4
3. Solution Requirements	4
4. Dynamic Flooding	5
4.1. Applicability	6
4.2. Leader election	7
4.3. Computing the Flooding Topology	7
4.4. Topologies on Complete Bipartite Graphs	8
4.4.1. A Minimal Flooding Topology	8
4.4.2. Xia Topologies	9
4.4.3. Optimization	10
4.5. Encoding the Flooding Topology	10
4.6. Analysis of Topology Changes	10
4.6.1. Link Addition	10
4.6.2. Node Addition	11
4.6.3. Link Failures Off the Flooding Topology	11
4.6.4. Failure of the Area Leader	11
4.6.5. Failures on the Flooding Topology	11
4.6.6. Recovery from Multiple Failures	12
5. Protocol Elements	12
5.1. IS-IS TLVs	12
5.1.1. IS-IS Area Leader Sub-TLV	13
5.1.2. IS-IS Area System IDs TLV	14
5.1.3. IS-IS Flooding Path TLV	15
5.2. OSPF LSAs and TLVs	16
5.2.1. OSPF Area Leader Sub-TLV	16
5.2.2. OSPFv2 Dynamic Flooding Opaque LSA	16
5.2.3. OSPFv3 Dynamic Flooding LSA	18
5.2.4. OSPF Area Router IDs TLV	18
5.2.5. OSPF Flooding Path TLV	19
6. Behavioral Specification	20
6.1. Leader Election	21
6.2. Area Leader Responsibilities	21
6.3. Distributed Flooding Topology Calculation	21
6.4. Flooding Behavior	22
7. IANA Considerations	22
7.1. IS-IS	22
7.2. OSPF	23
7.2.1. OSPF Dynamic Flooding LSA TLVs Registry	24
7.3. IGP	24

8. Security Considerations	25
9. Acknowledgements	25
10. References	25
10.1. Normative References	25
10.2. Informative References	27
Authors' Addresses	27

1. Introduction

In recent years, there has been increased focused on how to address the dynamic routing of networks that have a bipartite (a.k.a. spine-leaf or leaf-spine), Clos [Clos], or Fat Tree [Leiserson] topology. Conventional Interior Gateway Protocols (IGPs, i.e., IS-IS [ISO10589], OSPFv2 [RFC2328], and OSPFv3 [RFC5340]) under-perform, redundantly flooding information throughout the dense topology, leading to overloaded control plane inputs and thereby creating operational issues. For practical considerations, network architects have resorted to applying unconventional techniques to address the problem, applying BGP in the data center [RFC7938]. However it is very clear that using an Exterior Gateway Protocol as an IGP is sub-optimal, if only due to the configuration overhead.

The primary issue that is demonstrated when conventional mechanisms are applied is the poor reaction of the network to topology changes. Normal link state routing protocols rely on a flooding algorithm for state distribution. In a dense topology, this flooding algorithm is highly redundant, resulting in unnecessary overhead. Each node in the topology receives each link state update multiple times. Ultimately, all of the redundant copies will be discarded, but only after they have reached the control plane and been processed. This creates issues because significant link state database updates can become queued behind many redundant copies of another update. This delays convergence as the link state database does not stabilize promptly.

In a real world implementation, the packet queues leading to the control plane are necessarily of finite size, so if the flooding rate exceeds the update processing rate for long enough, the control plane will be obligated to drop incoming updates. If these lost updates are of significance, this will further delay stabilization of the link state database and the convergence of the network.

This is not a new problem. Historically, when routing protocols have been deployed in networks where the underlying topology is a complete graph, there have been similar issues. This was more common when the underlying link layer fabric presented the network layer with a full mesh of virtual connections. This was addressed by reducing the

flooding topology through IS-IS Mesh Groups [RFC2973], but this approach requires careful configuration of the flooding topology.

Thus, the root problem is not limited to massively scalable data centers. It exists with any dense topology at scale.

This problem is not entirely surprising. Link state routing protocols were conceived when links were very expensive and topologies were sparse. The fact that those same designs are sub-optimal in a dense topology should not come as a huge surprise. The fundamental premise that was addressed by the original designs was an environment of extreme cost and scarcity. Technology has progressed to the point where links are cheap and common. This represents a complete reversal in the economic fundamentals of network engineering. The original designs are to be commended for continuing to provide correct operation to this point, and optimizations for operation in today's environment are to be expected.

1.1. Requirements Language

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 [RFC2119].

2. Problem Statement

In a dense topology, the flooding algorithm that is the heart of conventional link state routing protocols causes a great deal of redundant messaging. This is exacerbated by scale. While the protocol can survive this combination, the redundant messaging is unnecessary overhead and delays convergence. Thus, the problem is to provide routing in dense, scalable topologies with rapid convergence.

3. Solution Requirements

A solution to this problem must then meet the following requirements:

Requirement 1 Provide a dynamic routing solution. Reachability must be restored after any topology change.

Requirement 2 Provide a significant improvement in convergence.

Requirement 3 The solution should address a variety of dense topologies. Just addressing a complete bipartite topology such as K5,8 is insufficient. Multi-stage Clos topologies must also be addressed, as well as topologies that are slight variants. Addressing complete graphs is a good demonstration of generality.

Requirement 4 There must be no single point of failure. The loss of any link or node should not unduly hinder convergence.

Requirement 5 Dense topologies are subgraphs of much larger topologies. Operational efficiency requires that the dense subgraph not operate in a radically different manner than the remainder of the topology. While some operational differences are permissible, they should be minimized. Changes to nodes outside of the dense subgraph are not acceptable. These situations occur when massively scaled data centers are part of an overall larger wide-area network. Having a second protocol operating just on this subgraph would add much more complexity at the edge of the subgraph where the two protocols would have to inter-operate.

4. Dynamic Flooding

We have observed that the combination of the dense topology and flooding on the physical topology in a scalable network is sub-optimal. However, if we decouple the flooding topology from the physical topology and only flood on a greatly reduced portion of that topology, we can have efficient flooding and retain all of the resilience of existing protocols.

In this idea, the flooding topology is computed either centrally on an elected node or in a distributed manner on all nodes that are supporting Dynamic Flooding. If the flooding topology is computed centrally, it is encoded into and distributed as part of the normal link state database. We call this the centralized mode of operation. If the flooding topology is computed in a distributed fashion, we call this the distributed mode of operation. Nodes within the dense topology would only flood on the flooding topology. On links outside of the normal flooding topology, normal database synchronization mechanisms (i.e., OSPF database exchange, IS-IS CSNPs) would apply, but flooding would not. New link state information that arrives from outside of the flooding topology suggests that the sender has a different or no flooding topology information and that the link state update should be flooded on the flooding topology as well.

Since the flooding topology is computed prior to topology changes, it does not factor into the convergence time and can be done when the topology is stable. The speed of the computation and its distribution, in the case of a centralized mode, is not a significant issue.

If a node does not have any flooding topology information when it receives new link state information, it should flood according to

legacy flooding rules. This situation will occur when the dense topology is first established, but is unlikely to recur.

When centralized mode is used and if, during a transient, there are multiple flooding topologies being advertised, then nodes should flood link state updates on all of the flooding topologies. Each node should locally evaluate the election of the lead node for the dense subgraph and first flood on the topology of the lead node. The rationale behind this is straightforward: if there is a transient and there has been a recent change in the elected node, then propagating topology information promptly along the most likely flooding topology should be the priority.

During transients, it is possible that loops will form in the flooding topology. This is not problematic, as the legacy flooding rules would cause duplicate updates to be ignored. Similarly, during transients, it is possible that the forwarding topology may become disconnected. To address this, nodes can perform a database synchronization check anytime a link is added to or removed from the flooding topology.

4.1. Applicability

In a complete graph, this approach is appealing because it drastically decreases the flooding topology without the manual configuration of mesh groups. By controlling the diameter of the flooding topology, as well as the maximum degree node in the flooding topology, convergence time goals can be met and the stability of the control plane can be assured.

Similarly, in a massively scaled data center, where there are many opportunities for redundant flooding, this mechanism ensures that flooding is redundant, with each leaf and spine well connected, while ensuring that no update need make too many hops and that no node shares an undue portion of the flooding effort.

In a network where only a portion of the nodes support Dynamic Flooding, the remaining nodes will continue to perform universal flooding. This is not an issue for correctness, as no node can become isolated.

Flooding that is initiated within the flooding topology will remain within that flooding topology until it reaches a legacy node, which will resume legacy flooding. Legacy flooding will be bounded by the flooding topology, which can help limit the propagation of unnecessary flooding. Whether or not the network can remain stable in this condition is unknown and may be very dependent on the number and location of the nodes that support Dynamic Flooding.

4.2. Leader election

A single node within the dense topology is elected as an Area Leader.

A generalization of the mechanisms used in existing Designated Router (OSPF) or Designated Intermediate-System (IS-IS) elections suffices. The elected node is known as the Area Leader.

In the case of centralized mode, the Area Leader is responsible for computing and distributing the flooding topology. When a new node is elected and has distributed new flooding topology information, then the old node should withdraw its flooding topology information from the link state database. If the old node does not return to the topology in a timely manner, the new node may remove the old node's information from the link state database.

In the case of distributed mode, the distributed algorithm advertised by the Area Leader **MUST** be used by all routers that participate in Dynamic Flooding.

Not every router needs to be a candidate to be Area Leader within an area, as a single candidate is sufficient for correct operation. For redundancy, however, it is strongly **RECOMMENDED** that there be multiple candidates.

4.3. Computing the Flooding Topology

There is a great deal of flexibility in how the flooding topology may be computed. For resilience, it needs to at least contain a cycle of all nodes in the dense subgraph. However, additional links could be added to decrease the convergence time. The trade-off between the density of the flooding topology and the convergence time is a matter for further study. The exact algorithm for computing the flooding topology in the case of the centralized computation need not be standardized, as it is not an interoperability issue. Only the encoding of the result needs to be documented. In the case of distributed mode, all nodes in the IGP area need to use the same algorithm to compute the flooding topology. It is possible to use private algorithms to compute flooding topology, so long as all nodes in the IGP area use the same algorithm.

While the flooding topology should be a covering cycle, it need not be a Hamiltonian cycle where each node appears only once. In fact, in many relevant topologies this will not be possible e.g., $K_{5,8}$. This is fortunate, as computing a Hamiltonian cycle is known to be NP-complete.

A simple algorithm to compute the topology for a complete bipartite graph is to simply select unvisited nodes on each side of the graph until both sides are completely visited. If the number of nodes on each side of the graph are unequal, then revisiting nodes on the less populated side of the graph will be inevitable. This algorithm can run in $O(N)$ time, so is quite efficient.

While a simple cycle is adequate for correctness and resiliency, it may not be optimal for convergence. At scale, a cycle may have a diameter that is half the number of nodes in the graph. This could cause an undue delay in link state update propagation. Therefore it may be useful to have a bound on the diameter of the flooding topology. Introducing more links into the flooding topology would reduce the diameter, but at the trade-off of possibly adding redundant messaging. The optimal trade-off between convergence time and graph diameter is for further study.

Similarly, if additional redundancy is added to the flooding topology, specific nodes in that topology may end up with a very high degree. This could result in overloading the control plane of those nodes, resulting in poor convergence. Thus, it may be optimal to have an upper bound on the degree of nodes in the flooding topology. Again, the optimal trade-off between graph diameter, node degree, and convergence time, and topology computation time is for further study.

If the leader chooses to include a multi-node broadcast LAN segment as part of the flooding topology, all of the connectivity to that LAN segment should be included as well. Once updates are flooded onto the LAN, they will be received by every attached node.

4.4. Topologies on Complete Bipartite Graphs

Complete bipartite graph topologies have become popular for data center applications and are commonly called leaf-spine or spine-leaf topologies. In this section, we discuss some flooding topologies that are of particular interest in these networks.

4.4.1. A Minimal Flooding Topology

We define a Minimal Flooding Topology on a complete bipartite graph as one in which the topology is connected and each node has at least degree two. This is of interest because it guarantees that the flooding topology has no single points of failure.

In practice, this implies that every leaf node in the flooding topology will have a degree of two. As there are usually more leaves than spines, the degree of the spines will be higher, but the load on the individual spines can be evenly distributed.

This type of flooding topology is also of interest because it scales well. As the number of leaves increases, we can construct flooding topologies that perform well. Specifically, for n spines and m leaves, if $m \geq n(n/2-1)$, then there is a flooding topology that has a diameter of four.

4.4.2. Xia Topologies

We define a Xia Topology on a complete bipartite graph as one in which all spine nodes are bi-connected through leaves with degree two, but the remaining leaves all have degree one and are evenly distributed across the spines.

Constructively, we can create a Xia topology by iterating through the spines. Each spine can be connected to the next spine by selecting any unused leaf. Since leaves are connected to all spines, all leaves will have a connection to both the first and second spine and we can therefore choose any leaf without loss of generality. Continuing this iteration across all of the spines, selecting a new leaf at each iteration, will result in a path that connects all spines. Adding one more leaf between the last and first spine will produce a cycle of n spines and n leaves.

At this point, $m-n$ leaves remain unconnected. These can be distributed evenly across the remaining spines, connected by a single link.

Xia topologies represent a compromise that trades off increased risk and decreased performance for lower flooding amplification. Xia topologies will have a larger diameter. For m spines, the diameter will be $m + 2$.

In a Xia topology, some leaves are singly connected. This represents a risk in that in some failures, convergence may be delayed. However, there may be some alternate behaviors that can be employed to mitigate these risks. If a leaf node sees that its single link on the flooding topology has failed, it can compensate by performing a database synchronization check with a different spine. Similarly, if a leaf determines that its connected spine on the flooding topology has failed, it can compensate by performing a database synchronization check with a different spine. In both of these cases, the synchronization check is intended to ameliorate any delays in link state propagation due to the fragmentation of the flooding topology.

The benefit of this topology is that flooding load is easily understood. Each node in the spine cycle will never receive an

update more than twice. For n leaves and m spines, a spine never transmits more than m/n updates.

4.4.3. Optimization

If two systems have multiple links between them, only one of the links should be part of the flooding topology. Moreover, symmetric selection of the link to use for flooding is not required.

4.5. Encoding the Flooding Topology

There are a variety of ways that the flooding topology could be encoded efficiently. If the topology was only a cycle, a simple list of the nodes in the topology would suffice. However, this is insufficiently flexible as it would require a slightly different encoding scheme as soon as a single additional link is added. Instead, we choose to encode the flooding topology as a set of intersecting paths, where each path is a set of connected edges.

Other encodings are certainly possible. We have attempted to make a useful trade off between simplicity, generality, and space.

4.6. Analysis of Topology Changes

In this section, we explicitly consider a variety of different topological failures in the network and how dynamic flooding should address them.

4.6.1. Link Addition

If a link is added to the topology, the protocol will form a normal adjacency on the link and update the appropriate link state advertisements for the routers on either end of the link. These link state updates will be flooded on the flooding topology.

In centralized mode, the Area Leader, upon receiving these updates, may choose to retain the existing flooding topology or may choose to modify the flooding topology. If it elects to change the flooding topology, it will update the flooding topology in the link state database and flood it using the new flooding topology.

In distributed mode, any change in the topology, including the link addition, should trigger the flooding topology recalculation. This is done to ensure that all nodes converge on the same flooding topology, regardless of the time of the calculation.

4.6.2. Node Addition

In centralized mode, if a node is added to the topology, then at least one link is also added to the topology. The paragraph above applies and the Area Leader will necessarily need to add the new node to the flooding topology.

In distributed mode, the addition of a node should trigger flooding topology recalculation.

Until the new node is incorporated into the flooding topology at least a single link towards the new node **MUST** be added to the flooding topology locally on all of its neighbors.

4.6.3. Link Failures Off the Flooding Topology

If a link that is not part of the flooding topology fails, then the adjoining routers will update their link state advertisements and flood them on the flooding topology. There is no need for changes to the flooding topology.

4.6.4. Failure of the Area Leader

The failure of the Area Leader can be detected by observing that it is disconnected from the area topology. In this case, the Area Leader election process is repeated and a new Area Leader is elected.

In the centralized mode, the new Area Leader will compute a new flooding topology and flood it using the new flooding topology.

As an optimization, applicable to centralized mode, the new Area Leader **MAY** compute a new flooding topology that has as much in common as possible with the old flooding topology. This will minimize the risk of over-flooding.

4.6.5. Failures on the Flooding Topology

If there is a failure on the flooding topology, the adjoining routers will update their link state advertisements and flood them. If the original flooding topology is bi-connected, the flooding topology should still be connected despite a single failure.

In centralized mode, the Area Leader will notice the change in the flooding topology, recompute the flooding topology, and flood it using the new flooding topology.

In distributed mode, all routers supporting dynamic flooding will notice the change in the flooding topology and recompute the new flooding topology.

4.6.6. Recovery from Multiple Failures

In the unlikely event of multiple failures on the flooding topology, it may become disconnected. The nodes that remain active on the edges of the flooding topology will recognize this, update their own link state advertisements and flood them on the remainder of the flooding topology. At this point, nodes will be able to compute that the flooding topology is partitioned.

Note that this is very different from partitioning the area itself. The area may remain connected and forwarding may still be effective.

When this condition is detected, the flooding topology can no longer be expected to deliver link state updates in a prompt manner. Nodes on the edges of the flooding topology should perform database synchronization on all links not on the flooding topology. Updates received from off of the flooding topology should be flooded on the remaining flooding topology. Any links that provide updates or require updates that are not part of the flooding topology should temporarily be added to the flooding topology. This should repair the current flooding topology, albeit in a sub-optimal manner.

In centralized mode, the Area Leader will also detect this condition, compute a new flooding topology, and flood it using the new flooding topology.

In distributed mode, all routers that actively participate in Dynamic Flooding will compute the new flooding topology.

5. Protocol Elements

5.1. IS-IS TLVs

The following TLVs are added to IS-IS:

1. A TLV that an IS may inject into its LSP to indicate its preference for becoming Area Leader.
2. A TLV to carry the list of system IDs that compromise the flooding topology for the area.
3. A TLV to carry the adjacency matrix for the flooding topology for the area.

5.1.1.1. IS-IS Area Leader Sub-TLV

The Area Leader Sub-TLV allows a system to:

1. Indicate its eligibility and priority for becoming Area Leader.
2. Indicate whether centralized or distributed mode is to be used to compute the flooding topology in the area.
3. Indicate the algorithm identifier for the algorithm that is used to compute the flooding topology in distributed mode.

Intermediate Systems (routers) that are not advertising this Sub-TLV are not eligible to become Area Leader.

The Area Leader is the router with the numerically highest Area Leader priority in the area. In the event of ties, the router with the numerically highest system ID is the Area Leader. Due to transients during database flooding, different routers may not agree on the Area Leader.

The Area Leader Sub-TLV is advertised as a Sub-TLV of the IS-IS Router Capability TLV-242 that is defined in [RFC7981] and has the following format:

0										1										2										3									
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	2	3	4	5	6	7	8	9
TLV Type										TLV Length										Priority										Algorithm									

TLV Type: TBD1

TLV Length: 2

Priority: 0-255, unsigned integer

Algorithm - a numeric identifier in the range 0-255 that identifies the algorithm used to calculate the flooding topology. The following values are defined:

0: Centralized computation by the Area Leader.

1-127: Standardized distributed algorithms. Individual values are assigned and managed by IANA. Before any assignments can be made, there MUST be an IETF specification that specifies IANA allocation for any value from this range (see Section 7.3).

128-254: Private distributed algorithms. Values from this range will not be registered with IANA and MUST NOT be mentioned by RFCs.

255: Reserved

5.1.2. IS-IS Area System IDs TLV

IS-IS Area System IDs TLV is only used in centralized mode.

The Area System IDs TLV is used by the Area Leader to enumerate the system IDs that it has used in computing the flooding topology. Conceptually, the Area Leader creates a list of system IDs for all routers in the area, assigning indices to each system, starting with index 0.

Because the space in a single TLV is small, more than one TLV may be required to encode all of the system IDs in the area. This TLV may be present in multiple LSPs.

The format of the Area System IDs TLV is:

```

      0               1               2               3
      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
+-----+-----+-----+-----+-----+-----+-----+-----+
| TLV Type           | TLV Length       | Starting Index           |
+-----+-----+-----+-----+-----+-----+-----+-----+
| L | Reserved       | System IDs ...
+-----+-----+-----+-----+-----+-----+-----+-----+
System IDs continued ....
+-----+-----+-----+-----+-----+-----+-----+-----+

```

TLV Type: TBD2

TLV Length: $3 + (\text{System ID length} * (\text{number of System IDs}))$

Starting index: The index of the first system ID that appears in this TLV.

L (Last): This bit is set if the index of the last system ID that appears in this TLV is equal to the last index in the full list of system IDs for the area.

System IDs: A concatenated list of system IDs for the area.

If there are multiple IS-IS Area System IDs TLVs with the L bit set advertised by the same router, the TLV which specifies the smaller maximum index is used and the other TLV(s) with L bit set are

ignored. TLVs which specify system IDs with indices greater than that specified by the TLV with the L bit set are also ignored.

5.1.3. IS-IS Flooding Path TLV

IS-IS Flooding Path TLV is only used in centralized mode.

The Flooding Path TLV is used to denote a path in the flooding topology. The goal is an efficient encoding of the links of the topology. A single link is a simple case of a path that only covers two nodes. A connected path may be described as a sequence of indices: (I1, I2, I3, ...), denoting a link from the system with index 1 to the system with index 2, a link from the system with index 2 to the system with index 3, and so on.

If a path exceeds the size that can be stored in a single TLV, then the path may be distributed across multiple TLVs by the replication of a single system index.

Complex topologies that are not a single path can be described using multiple TLVs.

The Flooding Path TLV contains a list of system indices relative to the systems advertised through the Area System IDs TLV. At least 2 indices must be included in the TLV. Due to the length restriction of TLVs, this TLV can contain at most 126 system indices.

The Flooding Path TLV has the format:

```

      0                               1                               2                               3
      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| TLV Type          | TLV Length      | Starting Index          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Index 2           | Additional indices ...
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

TLV Type: TBD3

TLV Length: 2 * (number of indices in the path)

Starting index: The index of the first system in the path.

Index 2: The index of the next system in the path.

Additional indices (optional): A sequence of additional indices to systems along the path.

5.2. OSPF LSAs and TLVs

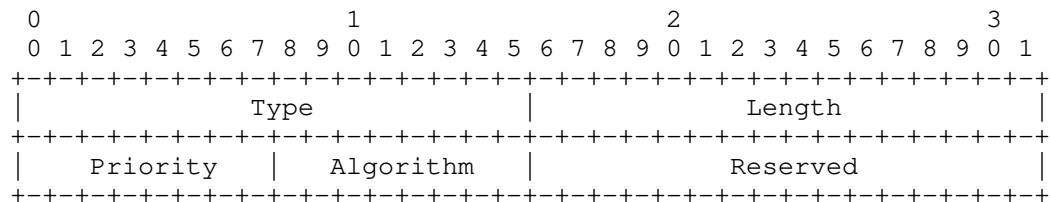
This section defines new LSAs and TLVs for both OSPFv2 and OSPFv3.

5.2.1. OSPF Area Leader Sub-TLV

The usage of the OSPF Area Leader Sub-TLV is identical to IS-IS and is described in Section 5.1.1.

The OSPF Area Leader Sub-TLV is used by both OSPFv2 and OSPFv3.

The OSPF Area Leader Sub-TLV is advertised as a top-level TLV of the RI LSA that is defined in [RFC7770] and has the following format:



Type: TBD4

Length: 4 octets

Priority: 0-255, unsigned integer

Algorithm: as defined in Section 5.1.1.

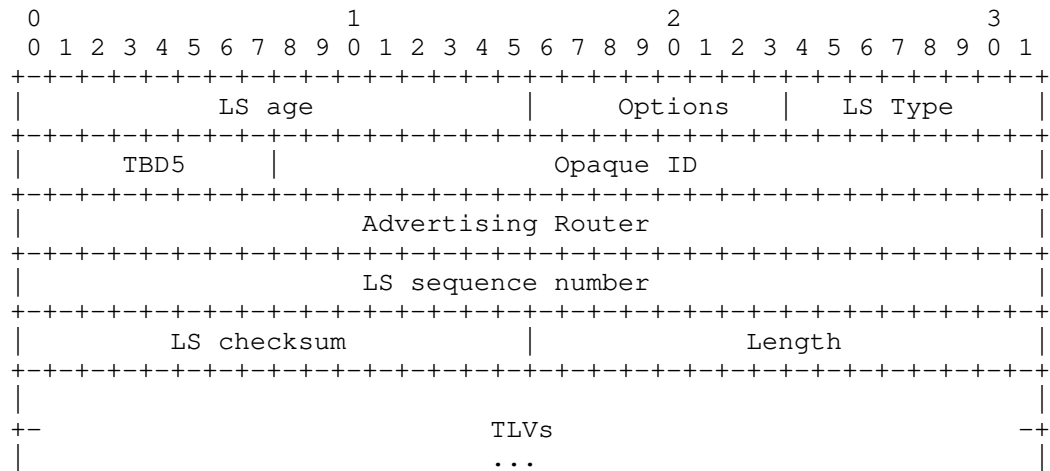
5.2.2. OSPFv2 Dynamic Flooding Opaque LSA

The OSPFv2 Dynamic Flooding Opaque LSA is only used in centralized mode.

The OSPFv2 Dynamic Flooding Opaque LSA is used to advertise additional data related to the dynamic flooding in OSPFv2. OSPFv2 Opaque LSAs are described in [RFC5250].

Multiple OSPFv2 Dynamic Flooding Opaque LSAs can be advertised by an OSPFv2 router. The flooding scope of the OSPFv2 Dynamic Flooding Opaque LSA is area-local.

The format of the OSPFv2 Dynamic Flooding Opaque LSA is as follows:



OSPFv2 Dynamic Flooding Opaque LSA

The opaque type used by OSPFv2 Dynamic Flooding Opaque LSA is TBD. The opaque type is used to differentiate the various type of OSPFv2 Opaque LSAs and is described in section 3 of [RFC5250]. The LS Type is 10. The LSA Length field [RFC2328] represents the total length (in octets) of the Opaque LSA including the LSA header and all TLVs (including padding).

The Opaque ID field is an arbitrary value used to maintain multiple Dynamic Flooding Opaque LSAs. For OSPFv2 Dynamic Flooding Opaque LSAs, the Opaque ID has no semantic significance other than to differentiate Dynamic Flooding Opaque LSAs originated by the same OSPFv2 router.

The format of the TLVs within the body of the OSPFv2 Dynamic Flooding Opaque LSA is the same as the format used by the Traffic Engineering Extensions to OSPF [RFC3630].

The Length field defines the length of the value portion in octets (thus a TLV with no value portion would have a length of 0). The TLV is padded to 4-octet alignment; padding is not included in the length field (so a 3-octet value would have a length of 3, but the total size of the TLV would be 8 octets). Nested TLVs are also 32-bit aligned. For example, a 1-octet value would have the length field set to 1, and 3 octets of padding would be added to the end of the value portion of the TLV. The padding is composed of zeros.

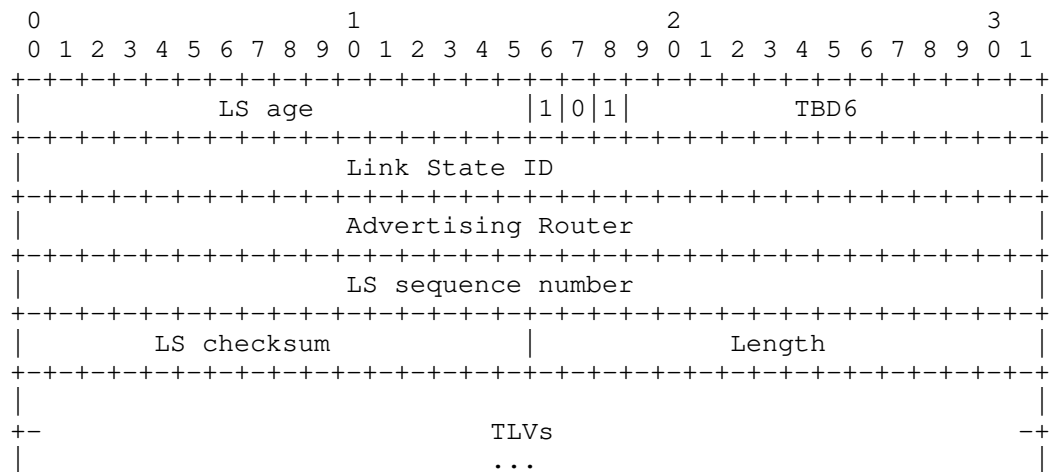
5.2.3. OSPFv3 Dynamic Flooding LSA

The OSPFv3 Dynamic Flooding Opaque LSA is only used in centralized mode.

The OSPFv3 Dynamic Flooding LSA is used to advertise additional data related to the dynamic flooding in OSPFv3.

The OSPFv3 Dynamic Flooding LSA has a function code of TBD. The flooding scope of the OSPFv3 Dynamic Flooding LSA is area-local. The U bit will be set indicating that the OSPFv3 Dynamic Flooding LSA should be flooded even if it is not understood. The Link State ID (LSID) value for this LSA is the Instance ID. OSPFv3 routers MAY advertise multiple Dynamic Flooding Opaque LSAs in each area.

The format of the OSPFv3 Dynamic Flooding LSA is as follows:



OSPFv3 Dynamic Flooding LSA

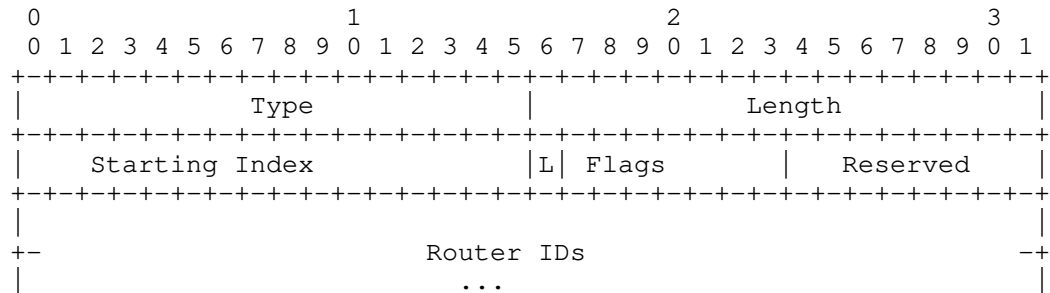
5.2.4. OSPF Area Router IDs TLV

The OSPF Area Router IDs TLV is a top level TLV of the OSPFv2 Dynamic Flooding Opaque LSA and OSPFv3 Dynamic Flooding LSA.

The OSPF Area Router IDs TLV is used by the Area Leader to enumerate the Router IDs that it has used in computing the flooding topology. Conceptually, the Area Leader creates a list of Router IDs for all routers in the area, assigning indices to each router, starting with index 0.

Because the space in a single OSPF Area Router IDs TLV is limited, more than one TLV may be required to encode all of the Router IDs in the area. This TLV may also recur in multiple OSPFv2 Dynamic Flooding Opaque LSAs or OSPFv3 Dynamic Flooding LSA, so that all Router IDs can be advertised.

The format of the Area Router IDs TLV is:



OSPF Area Router IDs TLV

TLV Type: 1

TLV Length: 4 + (Router ID length * (number of Router IDs))

Starting index: The index of the first Router ID that appears in this TLV.

L (Last): This bit is set if the index of the last system ID that appears in this TLV is equal to the last index in the full list of Router IDs for the area.

Router IDs: A concatenated list of Router IDs for the area.

If there are multiple OSPF Area Router IDs TLVs with the L bit set advertised by the same router, the TLV which specifies the smaller maximum index is used and the other TLV(s) with L bit set are ignored. TLVs which specify Router IDs with indices greater than that specified by the TLV with the L bit set are also ignored.

5.2.5. OSPF Flooding Path TLV

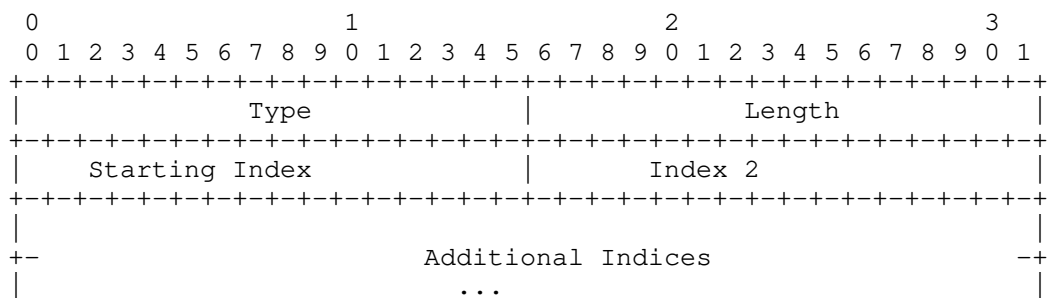
The OSPF Flooding Path TLV is a top level TLV of the OSPFv2 Dynamic Flooding Opaque LSAs and OSPFv3 Dynamic Flooding LSA.

The usage of the OSPF Flooding Path TLV is identical to IS-IS and is described in Section 5.1.3.

The OSPF Flooding Path TLV contains a list of Router ID indices relative to the Router IDs advertised through the OSPF Area Router IDs TLV. At least 2 indices must be included in the TLV.

Multiple OSPF Flooding Path TLVs can be advertised in a single OSPFv2 Dynamic Flooding Opaque LSA or OSPFv3 Dynamic Flooding LSA. OSPF Flooding Path TLVs can also be advertised in multiple OSPFv2 Dynamic Flooding Opaque LSAs or OSPFv3 Dynamic Flooding LSA, if they all can not fit in a single LSA.

The Flooding Path TLV has the format:



OSPF Flooding Path TLV

TLV Type: 2

TLV Length: 2 * (number of indices in the path)

Starting index: The index of the first Router ID in the path.

Index 2: The index of the next Router ID in the path.

Additional indices (optional): A sequence of additional indices to Router IDs along the path.

6. Behavioral Specification

In this section, we specify the detailed behaviors of the nodes participating in the IGP.

6.1. Leader Election

Any node that is capable MAY advertise its eligibility to become Area Leader.

Nodes that are not reachable are not eligible as Area Leader. Nodes that do not advertise their eligibility to become Area Leader are not eligible. Amongst the eligible nodes, the node with the numerically highest priority is the Area Leader. If multiple nodes all have the highest priority, then the node with the numerically highest system identifier in the case of IS-IS, or Router-ID in the case of OSPFv2 and OSPFv3 is the Area Leader.

6.2. Area Leader Responsibilities

If the Area Leader operates in centralized mode, it MUST advertise algorithm 0 in its Area Leader Sub-TLV. It also MUST compute and advertise a flooding topology for the area. The Area Leader MAY update the flooding topology at any time, however, it should not destabilize the network with undue or overly frequent topology changes.

The flooding topology MUST include all reachable nodes in the area. If nodes become unreachable on the flooding topology, the flooding topology MUST be recalculated. In centralized mode, the Area Leader MUST advertise a new flooding topology.

The flooding topology MAY be bi-connected. This is strongly RECOMMENDED but not required.

6.3. Distributed Flooding Topology Calculation

If the Area Leader advertises a non-zero algorithm in its Area Leader Sub-TLV, all routers in the area that support Dynamic Flooding and the value of algorithm advertised by the Area Leader MUST compute the flooding topology based on the Area Leader's advertised algorithm. Routers that do not support the value of algorithm advertised by the Area Leader MUST continue to use legacy flooding mechanism as defined by the protocol.

If the value of the algorithm advertised by the Area Leader is from the range 128-254 (Private distributed algorithms), it is the responsibility of the network operator to guarantee that all nodes in the area have a common understanding of what the given algorithm value represents.

6.4. Flooding Behavior

Nodes that support Dynamic Flooding MUST use the flooding topology for flooding. The flooding topology is calculated locally in the case of distributed mode. In centralized mode the flooding topology is advertised in the area link state database. Link state updates received on one link in the flooding topology MUST be flooded on all other links in the flooding topology other than the link on which the update has been received. Link state updates received on a link not in the flooding topology MUST be flooded on all links in the flooding topology.

In centralized mode, if multiple flooding topologies are present in the area link state database, the node SHOULD flood on the union of the topologies.

When the flooding topology changes on a node, either as a result of the local computation in distributed mode or as a result of the advertisement from the Area Leader in centralized mode, the node MUST continue to flood on the union of the old and new flooding topology for a limited amount of time. This is required to provide all nodes sufficient time to migrate to the new flooding topology.

When failures occur, nodes will learn about them from link state updates and can compare those to the existing flooding topology. If the flooding topology becomes disconnected, then the nodes at the edges of the flooding topology should perform a database synchronization on all links. While the flooding topology is disconnected, if a new link state update is received on a link not in the flooding topology, then the node SHOULD temporarily consider the link as part of the flooding topology. When a new flooding topology is received or locally calculated, this MUST be discontinued.

7. IANA Considerations

7.1. IS-IS

This document requests the following code point from the "sub-TLVs for TLV 242" registry (IS-IS Router CAPABILITY TLV).

Type: TBD1

Description: IS-IS Area Leader Sub-TLV

Reference: This document (Section 5.1.1)

This document requests that IANA allocate and assign two code points from the "IS-IS TLV Codepoints" registry. One for each of the following TLVs:

Type: TBD2

Description: IS-IS Area System IDs TLV

Reference: This document (Section 5.1.2)

Type: TBD3

Description: IS-IS Flooding Path TLV

Reference: This document (Section 5.1.3)

7.2. OSPF

This document requests the following code point from the "OSPF Router Information (RI) TLVs" registry:

Type: TBD4

Description: OSPF Area Leader Sub-TLV

Reference: This document (Section 5.2.1)

This document requests the following code point from the "Opaque Link-State Advertisements (LSA) Option Types" registry:

Type: TBD5

Description: OSPFv2 Dynamic Flooding Opaque LSA

Reference: This document (Section 5.2.2)

This document requests the following code point from the "OSPFv3 LSA Function Codes" registry:

Type: TBD6

Description: OSPFv3 Dynamic Flooding LSA

Reference: This document (Section 5.2.3)

7.2.1. OSPF Dynamic Flooding LSA TLVs Registry

This specification also requests one new registry - "OSPF Dynamic Flooding LSA TLVs". New values can be allocated via IETF Review or IESG Approval

The "OSPF Dynamic Flooding LSA TLVs" registry will define top-level TLVs for the OSPFv2 Dynamic Flooding Opaque LSA and OSPFv3 Dynamic Flooding LSAs. It should be added to the "Open Shortest Path First (OSPF) Parameters" registries group.

The following initial values are allocated:

Type: 0

Description: Reserved

Reference: This document

Type: 1

Description: OSPF Area Router IDs TLV

Reference: This document (Section 5.2.4)

Type: 2

Description: OSPF Flooding Path TLV

Reference: This document (Section 5.2.5)

Types in the range 32768-33023 are for experimental use; these will not be registered with IANA, and MUST NOT be mentioned by RFCs.

Types in the range 33024-65535 are not to be assigned at this time. Before any assignments can be made in the 33024-65535 range, there MUST be an IETF specification that specifies IANA Considerations that covers the range being assigned.

7.3. IGP

IANA is requested to set up a registry called "IGP Algorithm Type For Computing Flooding Topology" under an existing "Interior Gateway Protocol (IGP) Parameters" IANA registries. The registration policy for this registry is "Standards Action" ([RFC8126] and [RFC7120]).

Values in this registry come from the range 0-255.

The initial values in the IGP Algorithm Type For Computing Flooding Topology registry are:

0: Reserved for centralized mode.

1-127: Available for standards action.

128-254: Reserved for private use.

255: Reserved.

8. Security Considerations

This document introduces no new security issues. Security of routing within a domain is already addressed as part of the routing protocols themselves. This document proposes no changes to those security architectures.

It is possible that an attacker could become Area Leader and introduce a flawed flooding algorithm into the network thus compromising the operation of the protocol. Authentication methods as describe in [RFC5304] and [RFC5310] for IS-IS, [RFC2328] and [RFC7474] for OSPFv2 and [RFC5340] and [RFC4552] for OSPFv3 SHOULD be used to prevent such attack.

9. Acknowledgements

The authors would like to thank Les Ginsberg, Zeqing (Fred) Xia, Naiming Shen, Adam Sweeney and Olufemi Komolafe for their helpful comments.

The authors would like to thank Tom Edsall for initially introducing them to the problem.

10. References

10.1. Normative References

- [ISO10589]
International Organization for Standardization,
"Intermediate System to Intermediate System Intra-Domain
Routing Exchange Protocol for use in Conjunction with the
Protocol for Providing the Connectionless-mode Network
Service (ISO 8473)", ISO/IEC 10589:2002, Nov. 2002.

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC4552] Gupta, M. and N. Melam, "Authentication/Confidentiality for OSPFv3", RFC 4552, DOI 10.17487/RFC4552, June 2006, <<https://www.rfc-editor.org/info/rfc4552>>.
- [RFC5250] Berger, L., Bryskin, I., Zinin, A., and R. Coltun, "The OSPF Opaque LSA Option", RFC 5250, DOI 10.17487/RFC5250, July 2008, <<https://www.rfc-editor.org/info/rfc5250>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC7120] Cotton, M., "Early IANA Allocation of Standards Track Code Points", BCP 100, RFC 7120, DOI 10.17487/RFC7120, January 2014, <<https://www.rfc-editor.org/info/rfc7120>>.
- [RFC7474] Bhatia, M., Hartman, S., Zhang, D., and A. Lindem, Ed., "Security Extension for OSPFv2 When Using Manual Key Management", RFC 7474, DOI 10.17487/RFC7474, April 2015, <<https://www.rfc-editor.org/info/rfc7474>>.
- [RFC7770] Lindem, A., Ed., Shen, N., Vasseur, JP., Aggarwal, R., and S. Shaffer, "Extensions to OSPF for Advertising Optional Router Capabilities", RFC 7770, DOI 10.17487/RFC7770, February 2016, <<https://www.rfc-editor.org/info/rfc7770>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.

- [RFC8126] Cotton, M., Leiba, B., and T. Narten, "Guidelines for Writing an IANA Considerations Section in RFCs", BCP 26, RFC 8126, DOI 10.17487/RFC8126, June 2017, <<https://www.rfc-editor.org/info/rfc8126>>.

10.2. Informative References

- [Clos] Clos, C., "A Study of Non-Blocking Switching Networks", The Bell System Technical Journal Vol. 32(2), DOI 10.1002/j.1538-7305.1953.tb01433.x, March 1953, <<http://dx.doi.org/10.1002/j.1538-7305.1953.tb01433.x>>.
- [Leiserson] Leiserson, C., "Fat-Trees: Universal Networks for Hardware-Efficient Supercomputing", IEEE Transactions on Computers 34(10):892-901, 1985.
- [RFC2973] Balay, R., Katz, D., and J. Parker, "IS-IS Mesh Groups", RFC 2973, DOI 10.17487/RFC2973, October 2000, <<https://www.rfc-editor.org/info/rfc2973>>.
- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", RFC 3630, DOI 10.17487/RFC3630, September 2003, <<https://www.rfc-editor.org/info/rfc3630>>.
- [RFC7938] Lapukhov, P., Premji, A., and J. Mitchell, Ed., "Use of BGP for Routing in Large-Scale Data Centers", RFC 7938, DOI 10.17487/RFC7938, August 2016, <<https://www.rfc-editor.org/info/rfc7938>>.

Authors' Addresses

Tony Li
Arista Networks
5453 Great America Parkway
Santa Clara, California 95054
USA

Email: tony.li@tony.li

Peter Psenak
Cisco Systems, Inc.
Eurovea Centre, Central 3
Pribinova Street 10
Bratislava 81109
Slovakia

Email: ppsenak@cisco.com

Networking Working Group
Internet-Draft
Intended status: Standards Track
Expires: April 19, 2019

N. Shen
L. Ginsberg
Cisco Systems
S. Thyamagundalu
October 16, 2018

IS-IS Routing for Spine-Leaf Topology
draft-shen-isis-spine-leaf-ext-07

Abstract

This document describes a mechanism for routers and switches in a Spine-Leaf type topology to have non-reciprocal Intermediate System to Intermediate System (IS-IS) routing relationships between the leafs and spines. The leaf nodes do not need to have the topology information of other nodes and exact prefixes in the network. This extension also has application in the Internet of Things (IoT).

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 April 19, 2019.

Copyright Notice

Copyright (c) 2018 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. Requirements Language	3
2. Motivations	3
3. Spine-Leaf (SL) Extension	4
3.1. Topology Examples	4
3.2. Applicability Statement	5
3.3. Spine-Leaf TLV	6
3.3.1. Spine-Leaf Sub-TLVs	7
3.3.1.1. Leaf-Set Sub-TLV	7
3.3.1.2. Info-Req Sub-TLV	8
3.3.2. Advertising IPv4/IPv6 Reachability	8
3.3.3. Advertising Connection to RF-Leaf Node	8
3.4. Mechanism	8
3.4.1. Pure CLOS Topology	10
3.5. Implementation and Operation	11
3.5.1. CSNP PDU	11
3.5.2. Overload Bit	11
3.5.3. Spine Node Hostname	11
3.5.4. IS-IS Reverse Metric	11
3.5.5. Spine-Leaf Traffic Engineering	12
3.5.6. Other End-to-End Services	12
3.5.7. Address Family and Topology	12
3.5.8. Migration	13
4. IANA Considerations	13
5. Security Considerations	14
6. Acknowledgments	14
7. Document Change Log	14
7.1. Changes to draft-shen-isis-spine-leaf-ext-05.txt	14
7.2. Changes to draft-shen-isis-spine-leaf-ext-04.txt	14
7.3. Changes to draft-shen-isis-spine-leaf-ext-03.txt	14
7.4. Changes to draft-shen-isis-spine-leaf-ext-02.txt	14
7.5. Changes to draft-shen-isis-spine-leaf-ext-01.txt	15
7.6. Changes to draft-shen-isis-spine-leaf-ext-00.txt	15
8. References	15
8.1. Normative References	15
8.2. Informative References	16
Authors' Addresses	17

1. Introduction

The IS-IS routing protocol defined by [ISO10589] has been widely deployed in provider networks, data centers and enterprise campus environments. In the data center and enterprise switching networks, a Spine-Leaf topology is commonly used. This document describes a mechanism where IS-IS routing can be optimized for a Spine-Leaf topology.

In a Spine-Leaf topology, normally a leaf node connects to a number of spine nodes. Data traffic going from one leaf node to another leaf node needs to pass through one of the spine nodes. Also, the decision to choose one of the spine nodes is usually part of equal cost multi-path (ECMP) load sharing. The spine nodes can be considered as gateway devices to reach destinations on other leaf nodes. In this type of topology, the spine nodes have to know the topology and routing information of the entire network, but the leaf nodes only need to know how to reach the gateway devices to which are the spine nodes they are uplinked.

This document describes the IS-IS Spine-Leaf extension that allows the spine nodes to have all the topology and routing information, while keeping the leaf nodes free of topology information other than the default gateway routing information. The leaf nodes do not even need to run a Shortest Path First (SPF) calculation since they have no topology information.

1.1. Requirements Language

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 [RFC2119].

2. Motivations

- o The leaf nodes in a Spine-Leaf topology do not require complete topology and routing information of the entire domain since their forwarding decision is to use ECMP with spine nodes as default gateways
- o The spine nodes in a Spine-Leaf topology are richly connected to leaf nodes, which introduces significant flooding duplication if they flood all Link State PDUs (LSPs) to all the leaf nodes. It saves both spine and leaf nodes' CPU and link bandwidth resources if flooding is blocked to leaf nodes. For small Top of the Rack (ToR) leaf switches in data centers, it is meaningful to prevent full topology routing information and massive database flooding through those devices.

- o When a spine node advertises a topology change, every leaf node connected to it will flood the update to all the other spine nodes, and those spine nodes will further flood them to all the leaf nodes, causing a $O(n^2)$ flooding storm which is largely redundant.
- o Similar to some of the overlay technologies which are popular in data centers, the edge devices (leaf nodes) may not need to contain all the routing and forwarding information on the device's control and forwarding planes. "Conversational Learning" can be utilized to get the specific routing and forwarding information in the case of pure CLOS topology and in the events of link and node down.
- o Small devices and appliances of Internet of Things (IoT) can be considered as leafs in the routing topology sense. They have CPU and memory constrains in design, and those IoT devices do not have to know the exact network topology and prefixes as long as there are ways to reach the cloud servers or other devices.

3. Spine-Leaf (SL) Extension

3.1. Topology Examples

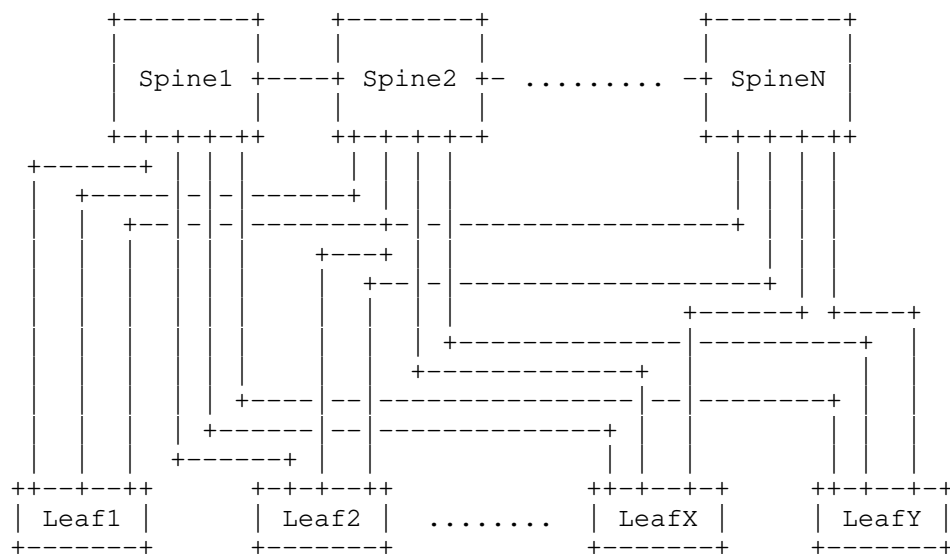


Figure 1: A Spine-Leaf Topology

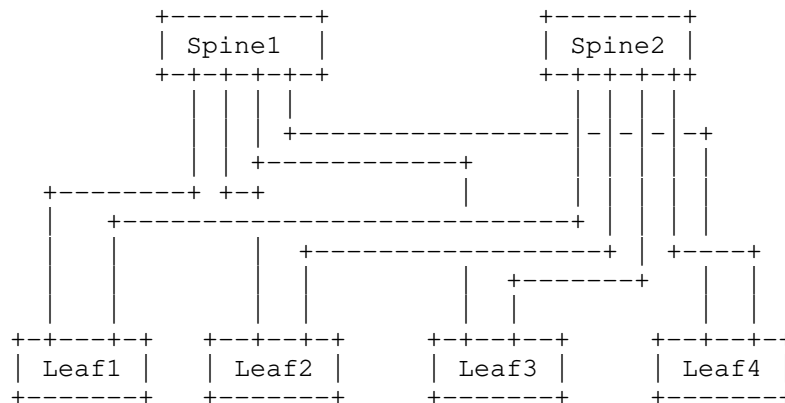


Figure 2: A CLOS Topology

3.2. Applicability Statement

This extension assumes the network is a Spine-Leaf topology, and it should not be applied in an arbitrary network setup. The spine nodes can be viewed as the aggregation layer of the network, and the leaf nodes as the access layer of the network. The leaf nodes use a load sharing algorithm with spine nodes as nexthops in routing and forwarding.

This extension works when the spine nodes are inter-connected, and it works with a pure CLOS or Fat Tree topology based network where the spines are NOT horizontally interconnected.

Although the example diagram in Figure 1 shows a fully meshed Spine-Leaf topology, this extension also works in the case where they are partially meshed. For instance, leaf1 through leaf10 may be fully meshed with spine1 through spine5 while leaf11 through leaf20 is fully meshed with spine4 through spine8, and all the spines are inter-connected in a redundant fashion.

This extension can also work in multi-level spine-leaf topology. The lower level spine node can be a 'leaf' node to the upper level spine node. A spine-leaf 'Tier' can be exchanged with IS-IS hello packets to allow tier X to be connected with tier X+1 using this extension. Normally tier-0 will be the TOR routers and switches if provisioned.

This extension also works with normal IS-IS routing in a topology with more than two layers of spine and leaf. For instance, in example diagrams Figure 1 and Figure 2, there can be another Core layer of routers/switches on top of the aggregation layer. From an IS-IS routing point of view, the Core nodes are not affected by this

extension and will have the complete topology and routing information just like the spine nodes. To make the network even more scalable, the Core layer can operate as a level-2 IS-IS sub-domain while the Spine and Leaf layers operate as stays at the level-1 IS-IS domain.

This extension assumes the link between the spine and leaf nodes are point-to-point, or point-to-point over LAN [RFC5309]. The links connecting among the spine nodes or the links between the leaf nodes can be any type.

3.3. Spine-Leaf TLV

This extension introduces a new TLV, the Spine-Leaf TLV, which may be advertised in IS-IS Hello (IIH) PDUs, LSPs, or in Circuit Scoped Link State PDUs (CS-LSP) [RFC7356]. It is used by both spine and leaf nodes in this Spine-Leaf mechanism.

```

      0               1               2               3
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          |      Length      |          SL Flag          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          .. Optional Sub-TLVs          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

The fields of this TLV are defined as follows:

Type: 1 octet Suggested value 150 (to be assigned by IANA)

Length: 1 octet (2 + length of sub-TLVs).

SL Flags: 16 bits

```

0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Tier |      Reserved      | T | R | L |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Tier: A value from 0 to 15. It represents the spine-leaf tier level. The value 15 is reserved to indicate the tier level is unknown. This value is only valid when the 'T' bit (see below) is set. If the 'T' bit is clear, this value MUST be set to zero on transmission, and it MUST be ignored on receipt.

L bit (0x01): Only leaf node sets this bit. If the L bit is set in the SL flag, the node indicates it is in 'Leaf-Mode'.

R bit (0x02): Only Spine node sets this bit. If the R bit is set, the node indicates to the leaf neighbor that it can be used as the default route gateway.

T bit (0x04): If set, the value in the "Tier" field (see above) is valid.

Optional Sub-TLV: Not defined in this document, for future extension

sub-TLVs MAY be included when the TLV is in a CS-LSP.
sub-TLVs MUST NOT be included when the TLV is in an IIH

3.3.1. Spine-Leaf Sub-TLVs

If the data center topology is a pure CLOS or Fat Tree, there are no link connections among the spine nodes. If we also assume there is not another Core layer on top of the aggregation layer, then the traffic from one leaf node to another may have a problem if there is a link outage between a spine node and a leaf node. For instance, in the diagram of Figure 2, if Leaf1 sends data traffic to Leaf3 through Spine1 node, and the Spine1-Leaf3 link is down, the data traffic will be dropped on the Spine1 node.

To address this issue spine and leaf nodes may send/request specific reachability information via the sub-TLVs defined below.

Two Spine-Leaf sub-TLVs are defined. The Leaf-Set sub-TLV and the Info-Req sub-TLV.

3.3.1.1. Leaf-Set Sub-TLV

This sub-TLV is used by spine nodes to optionally advertise Leaf neighbors to other Leaf nodes. The fields of this sub-TLV are defined as follows:

Type: 1 octet Suggested value 1 (to be assigned by IANA)

Length: 1 octet MUST be a multiple of 6 octets.

Leaf-Set: A list of IS-IS System-ID of the leaf node neighbors of this spine node.

3.3.1.2. Info-Req Sub-TLV

This sub-TLV is used by leaf nodes to request the advertisement of more specific prefix information from a selected spine node. The list of leaf nodes in this sub-TLV reflects the current set of leaf-nodes for which not all spine node neighbors have indicated the presence of connectivity in the Leaf-Set sub-TLV (See Section 3.3.1.1). The fields of this sub-TLV are defined as follows:

Type: 1 octet Suggested value 2 (to be assigned by IANA)

Length: 1 octet. It MUST be a multiple of 6 octets.

Info-Req: List of IS-IS System-IDs of leaf nodes for which connectivity information is being requested.

3.3.2. Advertising IPv4/IPv6 Reachability

In cases where connectivity between a leaf node and a spine node is down, the leaf node MAY request reachability information from a spine node as described in Section 3.3.1.2. The spine node utilizes TLVs 135 [RFC5305] and TLVs 236 [RFC5308] to advertise this information. These TLVs MAY be included either in IIHs or CS-LSPs [RFC7356] sent from the spine to the requesting leaf node. Sending such information in IIHs has limited scale - all reachability information MUST fit within a single IIH. It is therefore recommended that CS-LSPs be used.

3.3.3. Advertising Connection to RF-Leaf Node

For links between Spine and Leaf Nodes on which the Spine Node has set the R-bit and the Leaf node has set the L-bit in their respective Spine-Leaf TLVs, spine nodes may advertise the link with a bit in the "link-attribute" sub-TLV [RFC5029] to express this link is not used for LSP flooding. This information can be used by nodes computing a flooding topology e.g., [DYNAMIC-FLOODING], to exclude the RF-Leaf nodes from the computed flooding topology.

3.4. Mechanism

Leaf nodes in a spine-leaf application using this extension are provisioned with two attributes:

1) Tier level of 0. This indicates the node is a Leaf Node. The value 0 is advertised in the Tier field of Spine-Leaf TLV defined above.

2) Flooding reduction enabled/disabled. If flooding reduction is enabled the L-bit is set to one in the Spine-Leaf TLV defined above

A spine node does not need explicit configuration. Spine nodes can dynamically discover their tier level by computing the number of hops to a leaf node. Until a spine node determines its tier level it MUST advertise level 15 (unknown tier level) in the Spine-Leaf TLV defined above. Each tier level can also be statically provisioned on the node.

When a spine node receives an IIH which includes the Spine-Leaf TLV with Tier level 0 and 'L' bit set, it labels the point-to-point interface and adjacency to be a 'Reduced Flooding Leaf-Peer (RF-Leaf)'. IIHs sent by a spine node on a link to an RF-Leaf include the Spine-Leaf TLV with the 'R' bit set in the flags field. The 'R' bit indicates to the RF-Leaf neighbor that the spine node can be used as a default routing nexthop.

There is no change to the IS-IS adjacency bring-up mechanism for Spine-Leaf peers.

A spine node blocks LSP flooding to RF-Leaf adjacencies, except for the LSP PDUs in which the IS-IS System-ID matches the System-ID of the RF-Leaf neighbor. This exception is needed since when the leaf node reboots, the spine node needs to forward to the leaf node non-purged LSPs from the RF-Leaf's previous incarnation.

Leaf nodes will perform IS-IS LSP flooding as normal over all of its IS-IS adjacencies, but in the case of RF-Leafs only self-originated LSPs will exist in its LSP database.

Spine nodes will receive all the LSP PDUs in the network, including all the spine nodes and leaf nodes. It will perform Shortest Path First (SPF) as a normal IS-IS node does. There is no change to the route calculation and forwarding on the spine nodes.

The LSPs of a node only floods north bound towards the upper layer spine nodes. The default route is generated with loadsharing also towards the upper layer spine nodes.

RF-Leaf nodes do not have any LSP in the network except for its own. Therefore there is no need to perform SPF calculation on the RF-Leaf node. It only needs to download the default route with the nexthops of those Spine Neighbors which have the 'R' bit set in the Spine-Leaf TLV in IIH PDUs. IS-IS can perform equal cost or unequal cost load sharing while using the spine nodes as nexthops. The aggregated metric of the outbound interface and the 'Reverse Metric' [REVERSE-METRIC] can be used for this purpose.

3.4.1. Pure CLOS Topology

In a data center where the topology is pure CLOS or Fat Tree, there is no interconnection among the spine nodes, and there is not another Core layer above the aggregation layer with reachability to the leaf nodes. When flooding reduction to RF-Leafs is in use, if the link between a spine and a leaf goes down, there is then a possibility of black holing the data traffic in the network.

As in the diagram Figure 2, if the link Spine1-Leaf3 goes down, there needs to be a way for Leaf1, Leaf2 and Leaf4 to avoid the Spine1 if the destination of data traffic is to Leaf3 node.

In the above example, the Spine1 and Spine2 are provisioned to advertise the Leaf-Set sub-TLV of the Spine-Leaf TLV. Originally both Spines will advertise Leaf1 through Leaf4 as their Leaf-Set. When the Spine1-Leaf3 link is down, Spine1 will only have Leaf1, Leaf2 and Leaf4 in its Leaf-Set. This allows the other leaf nodes to know that Spine1 has lost connectivity to the leaf node of Leaf3.

Each RF-Leaf node can select another spine node to request for some prefix information associated with the lost leaf node. In this diagram of Figure 2, there are only two spine nodes (Spine-Leaf topology can have more than two spine nodes in general). Each RF-Leaf node can independently select a spine node for the leaf information. The RF-Leaf nodes will include the Info-Req sub-TLV in the Spine-Leaf TLV in hellos sent to the selected spine node, Spine2 in this case.

The spine node, upon receiving the request from one or more leaf nodes, will find the IPv6/IPv4 prefixes advertised by the leaf nodes listed in the Info-Req sub-TLV. The spine node will use the mechanism defined in Section 3.3.2 to advertise these prefixes to the RF-Leaf node. For instance, it will include the IPv4 loopback prefix of leaf3 based on the policy configured or administrative tag attached to the prefixes. When the leaf nodes receive the more specific prefixes, they will install the advertised prefixes towards the other spine nodes (Spine2 in this example).

For instance in the data center overlay scenario, when any IP destination or MAC destination uses the leaf3's loopback as the tunnel nexthop, the overlay tunnel from leaf nodes will only select Spine2 as the gateway to reach leaf3 as long as the Spine1-Leaf3 link is still down.

In cases where multiple links or nodes fail at the same time, the RF-leaf node may need to send the Info-Req to multiple upper layer spine

nodes in order to obtain reachability information for all the partially connected nodes.

This negative routing is more useful between tier 0 and tier 1 spine-leaf levels in a multi-level spine-leaf topology when the reduced flooding extension is in use. Nodes in tiers 1 or greater may have much richer topology information and alternative paths.

3.5. Implementation and Operation

3.5.1. CSNP PDU

In Spine-Leaf extension, Complete Sequence Number PDU (CSNP) does not need to be transmitted over the Spine-Leaf link to an RF-Leaf. Some IS-IS implementations send periodic CSNPs after the initial adjacency bring-up over a point-to-point interface. There is no need for this optimization here since the RF-Leaf does not need to receive any other LSPs from the network, and the only LSPs transmitted across the Spine-Leaf link is the leaf node LSP.

Also in the graceful restart case[RFC5306], for the same reason, there is no need to send the CSNPs over the Spine-Leaf interface to an RF-Leaf. Spine nodes only need to set the SRMflag on the LSPs belonging to the RF-Leaf.

3.5.2. Overload Bit

The leaf node SHOULD set the 'overload' bit on its LSP PDU, since if the spine nodes were to forward traffic not meant for the local node, the leaf node does not have the topology information to prevent a routing/forwarding loop.

3.5.3. Spine Node Hostname

This extension creates a non-reciprocal relationship between the spine node and leaf node. The spine node will receive leaf's LSP and will know the leaf's hostname, but the leaf does not have spine's LSP. This extension allows the Dynamic Hostname TLV [RFC5301] to be optionally included in spine's IIH PDU when sending to a 'Leaf-Peer'. This is useful in troubleshooting cases.

3.5.4. IS-IS Reverse Metric

This metric is part of the aggregated metric for leaf's default route installation with load sharing among the spine nodes. When a spine node is in 'overload' condition, it should use the IS-IS Reverse Metric TLV in IIH [REVERSE-METRIC] to set this metric to maximum to discourage the leaf using it as part of the loadsharing.

In some cases, certain spine nodes may have less bandwidth in link provisioning or in real-time condition, and it can use this metric to signal to the leaf nodes dynamically.

In other cases, such as when the spine node loses a link to a particular leaf node, although it can redirect the traffic to other spine nodes to reach that destination leaf node, but it MAY want to increase this metric value if the inter-spine connection becomes over utilized, or the latency becomes an issue.

In the leaf-leaf link as a backup gateway use case, the 'Reverse Metric' SHOULD always be set to very high value.

3.5.5. Spine-Leaf Traffic Engineering

Besides using the IS-IS Reverse Metric by the spine nodes to affect the traffic pattern for leaf default gateway towards multiple spine nodes, the IPv6/IPv4 Info-Advertise sub-TLVs can be selectively used by traffic engineering controllers to move data traffic around the data center fabric to alleviate congestion and to reduce the latency of a certain class of traffic pairs. By injecting more specific leaf node prefixes, it will allow the spine nodes to attract more traffic on some underutilized links.

3.5.6. Other End-to-End Services

Losing the topology information will have an impact on some of the end-to-end network services, for instance, MPLS TE or end-to-end segment routing. Some other mechanisms such as those described in PCE [RFC4655] based solution may be used. In this Spine-Leaf extension, the role of the leaf node is not too much different from the multi-level IS-IS routing while the level-1 IS-IS nodes only have the default route information towards the node which has the Attach Bit (ATT) set, and the level-2 backbone does not have any topology information of the level-1 areas. The exact mechanism to enable certain end-to-end network services in Spine-Leaf network is outside the scope of this document.

3.5.7. Address Family and Topology

IPv6 Address families[RFC5308], Multi-Topology (MT)[RFC5120] and Multi-Instance (MI)[RFC8202] information is carried over the IIH PDU. Since the goal is to simplify the operation of IS-IS network, for the simplicity of this extension, the Spine-Leaf mechanism is applied the same way to all the address families, MTs and MIs.

3.5.8. Migration

For this extension to be deployed in existing networks, a simple migration scheme is needed. To support any leaf node in the network, all the involved spine nodes have to be upgraded first. So the first step is to migrate all the involved spine nodes to support this extension, then the leaf nodes can be enabled with 'Leaf-Mode' one by one. No flag day is needed for the extension migration.

4. IANA Considerations

A new TLV codepoint is defined in this document and needs to be assigned by IANA from the "IS-IS TLV Codepoints" registry. It is referred to as the Spine-Leaf TLV and the suggested value is 150. This TLV is only to be optionally inserted either in the IIH PDU or in the Circuit Flooding Scoped LSP PDU. IANA is also requested to maintain the SL-flag bit values in this TLV, and 0x01, 0x02 and 0x04 bits are defined in this document.

Value	Name	IIH	LSP	SNP	Purge	CS-LSP
-----	-----	---	---	---	-----	-----
150	Spine-Leaf	y	y	n	n	y

This extension also proposes to have the Dynamic Hostname TLV, already assigned as code 137, to be allowed in IIH PDU.

Value	Name	IIH	LSP	SNP	Purge
-----	-----	---	---	---	-----
137	Dynamic Name	y	y	n	y

Two new sub-TLVs are defined in this document and needs to be added assigned by IANA from the "IS-IS TLV Codepoints". They are referred to in this document as the Leaf-Set sub-TLV and the Info-Req sub-TLV. It is suggested to have the values 1 and 2 respectively.

This document also requests that IANA allocate from the registry of link-attribute bit values for sub-TLV 19 of TLV 22 (Extended IS reachability TLV). This new bit is referred to as the "Connect to RF-Leaf Node" bit.

Value	Name	Reference
-----	-----	-----
0x3	Connect to RF-Leaf Node	This document

5. Security Considerations

Security concerns for IS-IS are addressed in [ISO10589], [RFC5304], [RFC5310], and [RFC7602]. This extension does not raise additional security issues.

6. Acknowledgments

The authors would like to thank Tony Przygienda for his discussion and contributions. The authors also would like to thank Acee Lindem, Russ White and Christian Hopps for their review and comments of this document.

7. Document Change Log

7.1. Changes to draft-shen-isis-spine-leaf-ext-05.txt

- o Submitted January 2018.
- o Just a refresh.

7.2. Changes to draft-shen-isis-spine-leaf-ext-04.txt

- o Submitted June 2017.
- o Added the Tier level information to handle the multi-level spine-leaf topology using this extension.

7.3. Changes to draft-shen-isis-spine-leaf-ext-03.txt

- o Submitted March 2017.
- o Added the Spine-Leaf sub-TLVs to handle the case of data center pure CLOS topology and mechanism.
- o Added the Spine-Leaf TLV and sub-TLVs can be optionally inserted in either IIH PDU or CS-LSP PDU.
- o Allow use of prefix Reachability TLVs 135 and 236 in IIHs/CS-LSPs sent from spine to leaf.

7.4. Changes to draft-shen-isis-spine-leaf-ext-02.txt

- o Submitted October 2016.
- o Removed the 'Default Route Metric' field in the Spine-Leaf TLV and changed to using the IS-IS Reverse Metric in IIH.

7.5. Changes to draft-shen-isis-spine-leaf-ext-01.txt

- o Submitted April 2016.
- o No change. Refresh the draft version.

7.6. Changes to draft-shen-isis-spine-leaf-ext-00.txt

- o Initial version of the draft is published in November 2015.

8. References

8.1. Normative References

[ISO10589]

ISO "International Organization for Standardization",
"Intermediate system to Intermediate system intra-domain
routing information exchange protocol for use in
conjunction with the protocol for providing the
connectionless-mode Network Service (ISO 8473), ISO/IEC
10589:2002, Second Edition.", Nov 2002.

[REVERSE-METRIC]

Shen, N., Amante, S., and M. Abrahamsson, "IS-IS Routing
with Reverse Metric", draft-ietf-isis-reverse-metric-07
(work in progress), 2017.

[RFC2119]

Bradner, S., "Key words for use in RFCs to Indicate
Requirement Levels", BCP 14, RFC 2119,
DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.

[RFC5029]

Vasseur, JP. and S. Previdi, "Definition of an IS-IS Link
Attribute Sub-TLV", RFC 5029, DOI 10.17487/RFC5029,
September 2007, <<https://www.rfc-editor.org/info/rfc5029>>.

[RFC5120]

Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi
Topology (MT) Routing in Intermediate System to
Intermediate Systems (IS-ISs)", RFC 5120,
DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.

[RFC5301]

McPherson, D. and N. Shen, "Dynamic Hostname Exchange
Mechanism for IS-IS", RFC 5301, DOI 10.17487/RFC5301,
October 2008, <<https://www.rfc-editor.org/info/rfc5301>>.

- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5306] Shand, M. and L. Ginsberg, "Restart Signaling for IS-IS", RFC 5306, DOI 10.17487/RFC5306, October 2008, <<https://www.rfc-editor.org/info/rfc5306>>.
- [RFC5308] Hopps, C., "Routing IPv6 with IS-IS", RFC 5308, DOI 10.17487/RFC5308, October 2008, <<https://www.rfc-editor.org/info/rfc5308>>.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<https://www.rfc-editor.org/info/rfc5310>>.
- [RFC7356] Ginsberg, L., Previdi, S., and Y. Yang, "IS-IS Flooding Scope Link State PDUs (LSPs)", RFC 7356, DOI 10.17487/RFC7356, September 2014, <<https://www.rfc-editor.org/info/rfc7356>>.
- [RFC7602] Chunduri, U., Lu, W., Tian, A., and N. Shen, "IS-IS Extended Sequence Number TLV", RFC 7602, DOI 10.17487/RFC7602, July 2015, <<https://www.rfc-editor.org/info/rfc7602>>.
- [RFC8202] Ginsberg, L., Previdi, S., and W. Henderickx, "IS-IS Multi-Instance", RFC 8202, DOI 10.17487/RFC8202, June 2017, <<https://www.rfc-editor.org/info/rfc8202>>.

8.2. Informative References

- [DYNAMIC-FLOODING] Li, T., "Dynamic Flooding on Dense Graphs", draft-li-dynamic-flooding (work in progress), 2018.
- [RFC4655] Farrel, A., Vasseur, J., and J. Ash, "A Path Computation Element (PCE)-Based Architecture", RFC 4655, DOI 10.17487/RFC4655, August 2006, <<https://www.rfc-editor.org/info/rfc4655>>.

[RFC5309] Shen, N., Ed. and A. Zinin, Ed., "Point-to-Point Operation over LAN in Link State Routing Protocols", RFC 5309, DOI 10.17487/RFC5309, October 2008, <<https://www.rfc-editor.org/info/rfc5309>>.

Authors' Addresses

Naiming Shen
Cisco Systems
560 McCarthy Blvd.
Milpitas, CA 95035
US

Email: naiming@cisco.com

Les Ginsberg
Cisco Systems
821 Alder Drive
Milpitas, CA 95035
US

Email: ginsberg@cisco.com

Sanjay Thyamagundalu

Email: tsanjay@gmail.com

LSR Working Group
Internet-Draft
Intended status: Standards Track
Expires: December 30, 2018

A. Wang
China Telecom
June 28, 2018

OSPF Extend for Inter-Area Topology Retrieval
draft-wang-lsr-ospf-inter-area-topology-ext-00

Abstract

This document describes method to transfer the source router id of inter-area prefixes for OSPFv2 [RFC2328] and OSPFv3 [RFC5340], which is needed in topology retrieval processing for inter-area scenario.

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 <https://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 December 30, 2018.

Copyright Notice

Copyright (c) 2018 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 (<https://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	2
2. Conventions used in this document	2
3. Inter-Area Topology Retrieval Scenario	3
3.1. OSPFv2 Extend Solution (IPv4 Source Router ID)	4
3.2. OSPFv3 Extend Solution (IPv6 Source Router ID)	5
3.3. Prefix Source Router ID sub TLV	7
3.4. Extend LSA generate process	8
3.5. Inter-Area Topology Retrieval Process	8
4. Security Considerations	8
5. IANA Considerations	9
6. References	9
6.1. Normative References	9
6.2. Informative References	10
Author's Address	10

1. Introduction

BGP-LS [RFC7752] describes the methodology that using BGP protocol to transfer the Link-State information. Such method can enable SDN controller to collect the underlay network topology automatically.

But if the underlay network is divided into multi area and running OSPF protocol, it is not easy for the SDN controller to rebuild the multi-area topology, because normally the ABR that locates on the boundary of different area will hide the detail topology information in non-backbone area, and the router in backbone area that runs BGP-LS protocol can only get and report the summary network information in non-backbone area.

[RFC7794] introduces "IPv4/IPv6 Source Router IDs" TLV to label the source of the prefixes redistributed from different Level, this TLV can be used to reconstruct the detail overall topology within level 1 and level 2. Such solution can also be applied into network that run OSPF protocol, but the related LSP message must be redefined.

This draft gives such solution for the OSPF v2 and OSPF v3 protocol.

2. Conventions used in this document

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 [RFC2119] .

3. Inter-Area Topology Retrieval Scenario

Fig.1 illustrates the topology retrieval scenario when OSPF is running in multi-area. R0-R4 are routers in backbone area, S1-S4,T1-T4 are internal router in area 1 and area 2 respectively. R1 and R3 are border routers between area 0 and area 1; R2 and R4 are border routers between area 0 and area 2. N1 is the network between router S1 and S2, N2 is the network between router T1 and T2.

Normally, ABR router R1 or R3 will send the summary LSA(for OSPFv2) or Inter-Area-Prefix-LSAs(for OSPFv3) for network N1. When R0 receives such LSA, it can only know network N1 locates behind R1, and does not know where it is originated. When R0 reports the summary LSA information via BGP-LS protocol, the IP SDN controller can't certainly deduce the detail network topology within area 1. The situation is same as that in Area 2.

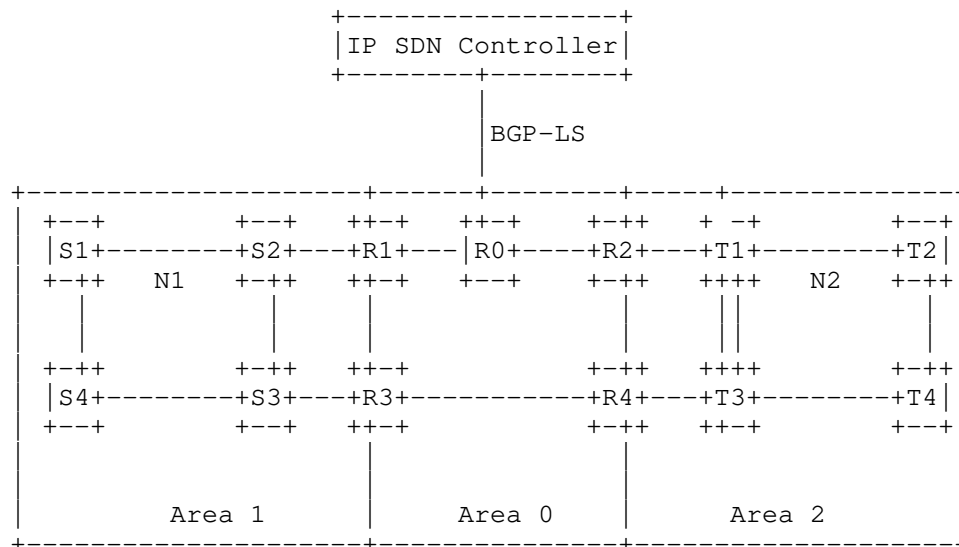


Fig.1 OSPF Inter-Area Topology Retrieval Scenario

If R0 has some methods to know the originator of network N1 and reports such information to IP SDN controller, then it is easy for the controller to retrieval the detail topology in non-backbone area.

Because traditional OSPFv2/v3 packet is not in the TLV format, we need to find some solutions to reuse or redefine the existing fields in summary LSA (OSPFv2) and Inter-Area-Prefix-LSAs(for OSPFv3)to transfer the additional information. The extend methods should not conflict with the usage of existing semantics.

Section 3.1 and section 3.2 give the proposed solutions for OSPFv2 and OSPFv3 respectively.

3.1. OSPFv2 Extend Solution (IPv4 Source Router ID)

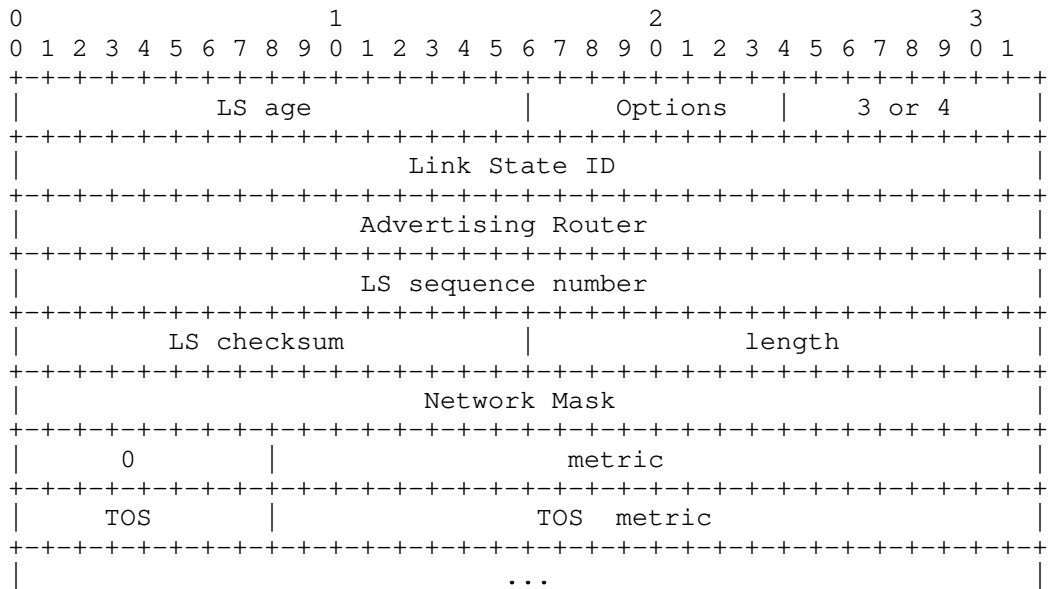


Fig.2 Summary LSA Format

Fig.2 illustrates the format of summary LSA. There is one byte that originally defined for the number of TOS types but in actually this feature does not applied in real network or implemented in the main stream router.

To transfer the additional information, this draft proposes to reuse/ redefine this field. In order to prevent possible conflict, even it is in very rare event, we can start the usage of this field from the upper limit, for example, 0xFE. Then the proposed extend summary LSA format is the followings:

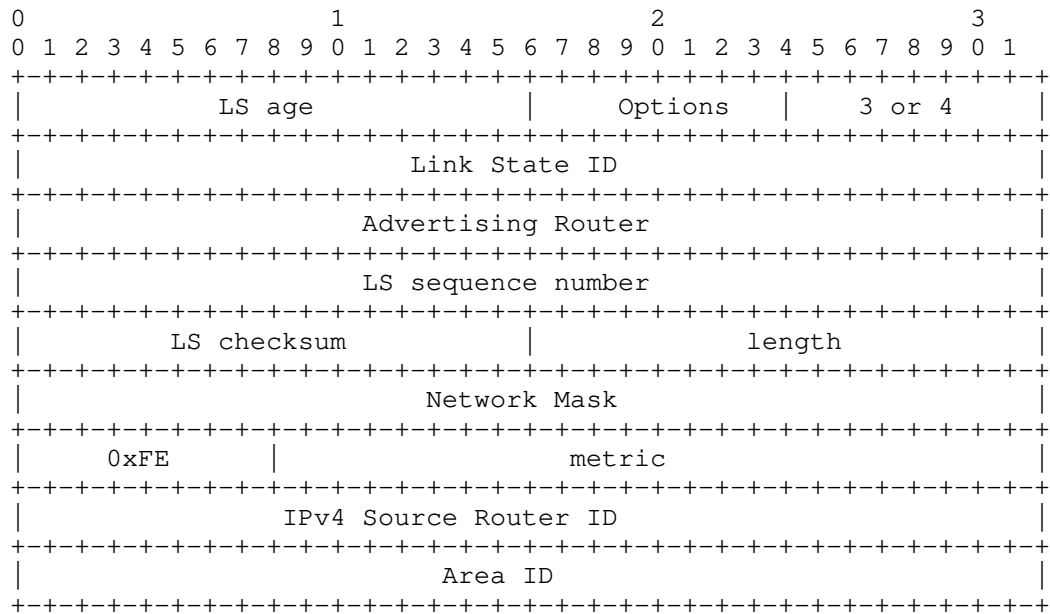


Fig.3 Extended Summary LSA Format

That is to say, if the field of "Numbers of TOS" equal "0xFE", then the "IPv4 Source Router ID" (4 bytes) of the inter-area network reported in summary LSA and its associated area id (4 bytes) are included in the field that follows the "metric" field.

3.2. OSPFv3 Extend Solution (IPv6 Source Router ID)

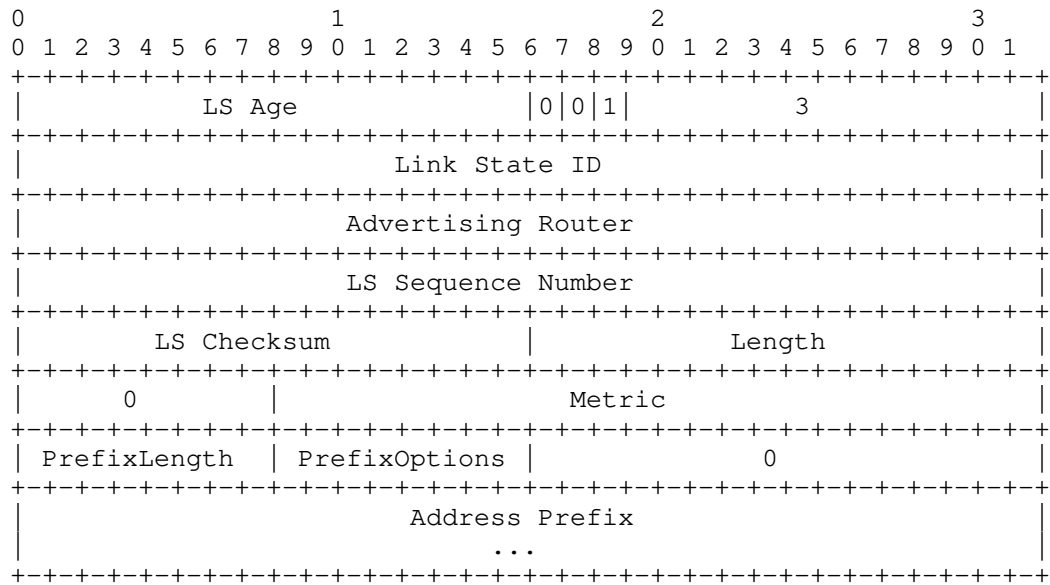


Fig.4 Inter-Area-Prefix-LSA Format

For OSPFv3, this draft proposes the similar method, because the semantic of the Inter-Area-Prefix-LSA format is almost same as the summary LSA format.

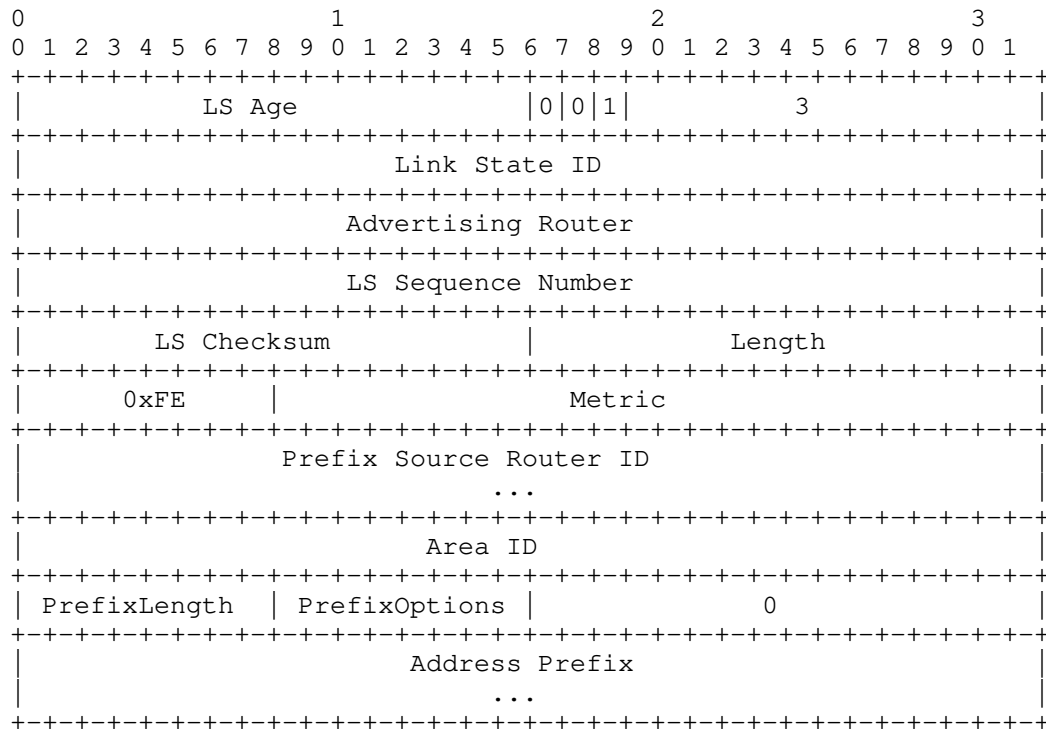


Fig.5 Extended Inter-Area-Prefix-LSA Format

If the value of "Numbers of TOS" equal "0xFE", then the "IPv6 source router ID" (16 bytes) and its corresponding area ID (4 bytes) information are inserted in the "Inter-Area-Prefix-LSA" after the field "Metric". After this, the normal Prefix information is followed as shown in Fig.5

3.3. Prefix Source Router ID sub TLV

[RFC7684] and [RFC8362] define the TLV format extension for OSPFv2 and OSPFv3 respectively. These documents give the flexibility to add new attributes for the prefixes and links. Based on these formats, we can define new sub TLV to transfer the "Prefix Source Router ID", as that defined in [RFC7794].

The proposed "Prefix Source Router ID" format is the following:

For IPv4 network, it is the following:

- o Pv4 Source Router ID Type: TBD
- o Length: 4

- o Value: IPv4 Router ID of the source of the advertisement

This sub TLV should be included in the "OSPFv2 Extended Prefix Opaque LSA" that defined in [RFC7684]

For IPv6 network, it is the following:

- o IPv6 Source Router ID Type: TBD
- o Length: 16
- o Value: IPv6 Router ID of the source of the advertisement

This sub TLV should be included in "E-Inter-Area-Prefix-LSA" that defined in [RFC8362]

3.4. Extend LSA generate process

When ABR(for example R1 in Fig.1)receives the "Router LSA" announcement in area 1, it should generate the corresponding extend "Summary LSA" or "Inter-Area-Prefix-LSA" that includes the "Source Router ID" of the network prefixes, which labels the corresponding link and the "area ID" that the source router belongs to.

When R0 receives such extend LSA, it then strips this additional information, put it into the corresponding part that in BGP-LS protocol as described in[I-D.wang-idr-bgpls-inter-as-topology-ext] and reports them to the IP SDN Controller.

3.5. Inter-Area Topology Retrieval Process

When IP SDN Controller receives this information, it should compare the prefix NLRI that included in the BGP-LS packet. When it encounters the same prefix but with different source router ID, it should extract the corresponding area ID, rebuild the link between these two different source router in non-backbone area.

Iterating the above process continuously, the IP SDN controller can then retrieve the detail topology that span multi-area.

4. Security Considerations

TBD.

5. IANA Considerations

TBD.

6. References

6.1. Normative References

- [I-D.ietf-pce-pcep-extension-native-ip]
Wang, A., Khasanov, B., Cheruathur, S., and C. Zhu, "PCEP Extension for Native IP Network", draft-ietf-pce-pcep-extension-native-ip-00 (work in progress), June 2018.
- [I-D.ietf-teas-native-ip-scenarios]
Wang, A., Huang, X., Qou, C., Huang, L., and K. Mi, "CCDR Scenario, Simulation and Suggestion", draft-ietf-teas-native-ip-scenarios-00 (work in progress), February 2018.
- [I-D.ietf-teas-pcecc-use-cases]
Zhao, Q., Li, Z., Khasanov, B., Ke, Z., Fang, L., Zhou, C., Communications, T., and A. Rachitskiy, "The Use Cases for Using PCE as the Central Controller(PCECC) of LSPs", draft-ietf-teas-pcecc-use-cases-01 (work in progress), May 2017.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2328] Moy, J., "OSPF Version 2", STD 54, RFC 2328, DOI 10.17487/RFC2328, April 1998, <<https://www.rfc-editor.org/info/rfc2328>>.
- [RFC5340] Coltun, R., Ferguson, D., Moy, J., and A. Lindem, "OSPF for IPv6", RFC 5340, DOI 10.17487/RFC5340, July 2008, <<https://www.rfc-editor.org/info/rfc5340>>.
- [RFC7684] Psenak, P., Gredler, H., Shakir, R., Henderickx, W., Tantsura, J., and A. Lindem, "OSPFv2 Prefix/Link Attribute Advertisement", RFC 7684, DOI 10.17487/RFC7684, November 2015, <<https://www.rfc-editor.org/info/rfc7684>>.
- [RFC7752] Gredler, H., Ed., Medved, J., Previdi, S., Farrel, A., and S. Ray, "North-Bound Distribution of Link-State and Traffic Engineering (TE) Information Using BGP", RFC 7752, DOI 10.17487/RFC7752, March 2016, <<https://www.rfc-editor.org/info/rfc7752>>.

- [RFC7794] Ginsberg, L., Ed., Decraene, B., Previdi, S., Xu, X., and U. Chunduri, "IS-IS Prefix Attributes for Extended IPv4 and IPv6 Reachability", RFC 7794, DOI 10.17487/RFC7794, March 2016, <<https://www.rfc-editor.org/info/rfc7794>>.
- [RFC8362] Lindem, A., Roy, A., Goethals, D., Reddy Vallem, V., and F. Baker, "OSPFv3 Link State Advertisement (LSA) Extensibility", RFC 8362, DOI 10.17487/RFC8362, April 2018, <<https://www.rfc-editor.org/info/rfc8362>>.

6.2. Informative References

- [I-D.wang-idr-bgpls-inter-as-topology-ext]
Wang, A., "BGP-LS extend for inter-AS topology retrieval", draft-wang-idr-bgpls-inter-as-topology-ext-00 (work in progress), March 2018.

Author's Address

Aijun Wang
China Telecom
Beiqijia Town, Changping District
Beijing, Beijing 102209
China

Email: wangaj.bri@chinatelecom.cn

Network Working Group
Internet-Draft
Intended status: Informational
Expires: May 9, 2019

R. White, Ed.
S. Zandi, Ed.
LinkedIn
November 5, 2018

IS-IS Support for Openfabric
draft-white-openfabric-07

Abstract

Spine and leaf topologies are widely used in hyperscale and cloud scale networks. In most of these networks, configuration is automated, but difficult, and topology information is extracted through broad based connections. Policy is often integrated into the control plane, as well, making configuration, management, and troubleshooting difficult. Openfabric is an adaptation of an existing, widely deployed link state protocol, Intermediate System to Intermediate System (IS-IS) that is designed to:

- o Provide a full view of the topology from a single point in the network to simplify operations
- o Minimize configuration of each Intermediate System (IS) (also called a router or switch) in the network
- o Optimize the operation of IS-IS within a spine and leaf fabric to enable scaling

This document begins with an overview of openfabric, including a description of what may be removed from IS-IS to enable scaling. The document then describes an optimized adjacency formation process; an optimized flooding scheme; some thoughts on the operation of openfabric, metrics, and aggregation; and finally a description of the changes to the IS-IS protocol required for openfabric.

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 <https://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 May 9, 2019.

Copyright Notice

Copyright (c) 2018 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 (<https://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. Goals	3
1.2. Contributors	3
1.3. Simplification	3
1.4. Additions and Requirements	4
1.5. Sample Network	4
2. Modified Adjacency Formation	6
2.1. Level 2 Adjacencies Only	6
2.2. Point-to-point Adjacencies	6
2.3. Three Way Handshake Support	7
2.4. Adjacency Formation Optimization	7
3. Advertisement of Reachability Information	8
4. Determining and Advertising Location on the Fabric	9
5. Flooding Optimization	10
5.1. Flooding Failures	11
6. Other Optimizations	12
6.1. Transit Link Reachability	12
6.2. Transiting T0 Intermediate Systems	12
7. Openfabric and Route Aggregation	13
8. Security Considerations	13
9. References	13
9.1. Normative References	13
9.2. Informative References	15
Appendix A. Flooding Optimization Operation	17
Appendix B. Fabric Location Calculation	19
Authors' Addresses	20

1. Introduction

1.1. Goals

Spine and leaf fabrics are often used in large scale data centers; in this application, they are commonly called a fabric because of their regular structure and predictable forwarding and convergence properties. This document describes modifications to the IS-IS protocol to enable it to run efficiently on a large scale spine and leaf fabric, openfabric. The goals of this control plane are:

- o Provide a full view of the topology from a single point in the network to simplify operations
- o Minimize configuration of each IS in the network
- o Optimize the operation of IS-IS within a spine and leaf fabric to enable scaling

1.2. Contributors

The following people have contributed to this draft: Nikos Triantafyllis (reflected flooding optimization), Ivan Pepelnjak (fabric locality calculation modifications), Christian Franke (fabric locality calculation modification), Hannes Gredler (do not reflood optimizations), Les Ginsberg (capabilities encoding, circuit local reflooding), Naiming Shen (capabilities encoding, circuit local reflooding), Uma Chunduri (failure mode suggestions, flooding), Nick Russo, and Rodny Molina.

See [RFC5449], [RFC5614], and [RFC7182] for similar solutions in the Mobile Ad Hoc Networking (MANET) solution space.

1.3. Simplification

In building any scalable system, it is often best to begin by removing what is not needed. In this spirit, openfabric implementations MAY remove the following from IS-IS:

- o External metrics. There is no need for external metrics in large scale spine and leaf fabrics; it is assumed that metrics will be properly configured by the operator to account for the correct order of route preference at any route redistribution point.
- o Tags and traffic engineering processing. Openfabric is only designed to provide topology and reachability information. It is not designed to provide for traffic engineering, route preference through tags, or other policy mechanisms. It is assumed that all

routing policy will be provided through an overlay system which communicates directly with each IS in the fabric, such as PCEP [RFC5440] or I2RS [RFC7921]. Traffic engineering is assumed to be provided through Segment Routing (SR) [I-D.ietf-spring-segment-routing].

1.4. Additions and Requirements

To create a scalable link state fabric, openfabric includes the following:

- o A slightly modified adjacency formation process.
- o Mechanisms for determining which tier within a spine and leaf fabric in which the IS is located.
- o A mechanism that reduces flooding to the minimum possible, while still ensuring complete database synchronization among the intermediate systems within the fabric.

Three general requirements are placed here; more specific requirements are considered in the following sections. Openfabric implementations:

- o MUST support [RFC5301] and enable hostname advertisement by default if a hostname is configured on the intermediate system.
- o SHOULD support [RFC6232], purge originator identification for IS-IS.
- o MUST NOT be mixed with standard IS-IS implementations in operational deployments. Openfabric and standard IS-IS implementations SHOULD be treated as two separate protocols.

1.5. Sample Network

The following spine and leaf fabric will be used to describe these modifications.

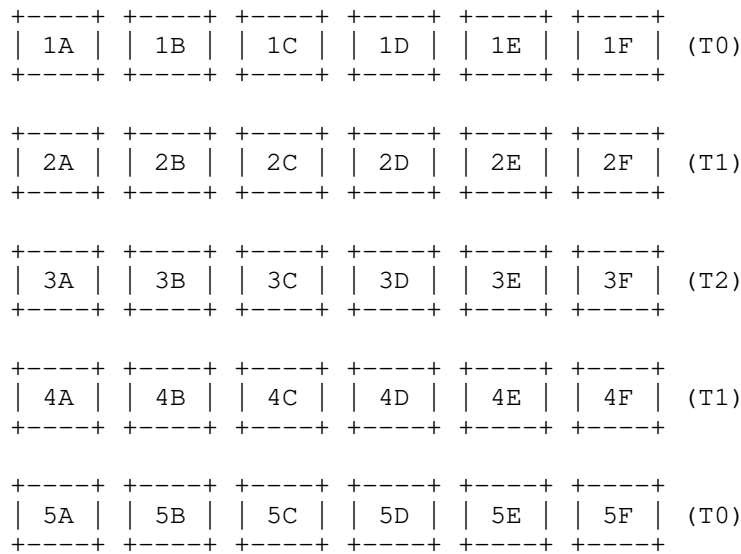


Figure 1

To reduce confusion (spine and leaf fabrics are difficult to draw in plain text art), this diagram does not contain the connections between devices. The reader should assume that each device in a given layer is connected to every device in the layer above it. For instance:

- o 5A is connected to 4A, 4B, 4C, 4D, 4E, and 4F
- o 5B is connected to 4A, 4B, 4C, 4D, 4E, and 4F
- o 4A is connected to 3A, 3B, 3C, 3D, 3E, 3F, 5A, 5B, 5C, 5D, 5E, and 5F
- o 4B is connected to 3A, 3B, 3C, 3D, 3E, 3F, 5A, 5B, 5C, 5D, 5E, and 5F
- o etc.

The tiers or stages of the fabric are also marked for easier reference. T0 is assumed to be connected to application servers, or rather they are Top of Rack (ToR) intermediate systems. The remaining tiers, T1 and T2, are connected only to the fabric itself. Note there are no "cross links," or "east west" links in the illustrated fabric. The fabric locality detection mechanism described here will not work if there are cross links running east/

west through the fabric. Locality detection may be possible in such a fabric; this is an area for further study.

2. Modified Adjacency Formation

Because Openfabric operates in a tightly controlled data center environment, various modifications can be made to the IS-IS neighbor formation process to increase efficiency and simplify the protocol. Specifically, Openfabric implementations SHOULD support [RFC3719], section 4, hello padding for IS-IS. Variable hello padding SHOULD NOT be used, as data center fabrics are built using high speed links on which padded hellos will have little performance impact. Further modifications to the neighbor formation process are considered in the following sections.

2.1. Level 2 Adjacencies Only

Openfabric is designed to work in a single flooding domain over a single data center fabric at the scale of thousands of routers with hundreds of thousands of routes (so a moderate scale in router and route count terms). Because of the way Openfabric optimizes operation in this environment, it is not necessary nor desirable to build multiple flooding domains. For instance, the flooding optimizations described later in this document require a full view of the topology, as does any proposed overlay to inject policy into the forwarding plane. In light of this, the following changes SHOULD BE to IS-IS implementations to support Openfabric:

- o IIH PDU 17 (level 2 point-to-point circuit hello) should be the only IIH PDU type transmitted (see section 9.7 of ISO 10589)
- o In IIH PDU 17 (level 2 point-to-point circuit hello), the Circuit Type field should be set to 2 (see section 9.7 of ISO 10589)
- o Support for IIH PDU 15 (level 1 broadcast hello) should be removed (see section 9.5 of ISO 10589)
- o Support for IIH PDU 16 (level 2 broadcast hello) should be removed (see section 9.6 of ISO 10589)

2.2. Point-to-point Adjacencies

Data center network fabrics only contain point-to-point links; because of this, there is no reason to support any broadcast link types, nor to support the Designated Intermediate System processing, including pseudonode creation. In light of this, processing related to sections 7.2.3 (broadcast networks), 7.3.8 (generation of level 1 pseudonode LSPs), 7.3.10 (generation of level 2 pseudonode LSPs), and

section 8.4.5 (LAN designated intermediate systems) in [ISO10589] SHOULD BE removed.

2.3. Three Way Handshake Support

It is important that two way connectivity be established before synchronizing the link state database, or routing through a link in a data center fabric. To reject optical failures that cause a one way connection between two routers, fabricDC must support the three way handshake mechanism described in [RFC5303].

2.4. Adjacency Formation Optimization

While adjacency formation is not considered particularly burdensome in IS-IS, it may still be useful to reduce the amount of state transferred across the network when connecting a new IS to the fabric. In its simplest form, the process is:

- o An IS connected to the fabric will send hellos on all links.
- o The IS will only complete the three-way handshake with one newly discovered neighbor; this would normally be the first neighbor which sends the newly connected intermediate system's ID back in the three-way handshake process.
- o The IS will complete its database exchange with this one newly adjacent neighbor.
- o Once this process is completed, the IS will continue processing the remaining neighbors as normal.
- o If synchronization is not achieved within twice the dead timer on the local interface, the newly connected IS will repeat this process with the second neighbor with which it forms a three-way adjacency.

This process allows each IS newly added to the fabric to exchange a full table once; a very minimal amount of information will be transferred with the remaining neighbors to reach full synchronization.

Any such optimization is bound to present a tradeoff between several factors; the mechanism described here increases the amount of time required to form adjacencies slightly in order to reduce the total state carried across the network. An alternative mechanism could provide a better balance of the amount of information carried across the network for initial synchronization and the time required to synchronize a new IS. For instance, an IS could choose to

synchronize its database with two or three adjacent intermediate systems, which could speed the synchronization process up at the cost of carrying additional data on the network. A locally determined balance between the speed of synchronization and the amount of data carried on the network can be achieved by adjusting the number of adjacent intermediate systems the newly attached IS synchronizes with.

3. Advertisement of Reachability Information

IS-IS describes the topology in two different sets of TLVs; the first describes the set of neighbors connected to an IS, the second describes the set of reachable destination connected to an IS. There are two different forms of both of these descriptions, one of which carries what are widely called narrow metrics, the other of which carries what are widely called wide metrics. In a tightly controlled data center fabric implementation, such as the ones Openfabric is designed to support, no IS that supports narrow metrics will ever be deployed or supported; hence there is no reason to support any metric type other than wide metrics.

- o The Level 2 Link State PDU (type 20 in section 9.9 of [ISO10589]) and the scoped flooding PDU (type 10 in section 3.1 of [RFC7356]) SHOULD BE the only PDU types used to carry link state information in a Openfabric implementation
- o Processing related to the Level 1 Link State PDU (type 18) MAY BE removed from Openfabric implementations (see section 9.8 of [ISO10589])
- o Neighbor reachability MUST BE carried in TLV type 22 (see section 3 of [RFC5305])
- o IPv4 reachability SHOULD BE carried in TLV type 135 (see section 4 of [RFC5305]), or TLV type 235 for multitopology implementations (see [RFC5120])
- o IPv6 reachability SHOULD BE carried in TLV type 236 (see [RFC5308]), or TLV type 237 for multitopology implemenations (see [RFC5120])
- o Processing related to the neighbor reachability TLV (type 2, see sections 9.8 and 9.9 of [ISO10589]) SHOULD BE removed
- o Processing related to the narrow metric IP reachability TLV (types 128 and 130) SHOULD BE removed

Further, if segment routing support is desired, Openfabric MAY support the Prefix Segment Identifier sub-TLV and other TLVs as required in [I-D.ietf-isis-segment-routing-extensions].

4. Determining and Advertising Location on the Fabric

The tier to which a IS is connected is useful to enable autoconfiguration of intermediate systems connected to the fabric and to reduce flooding. Once the tier of an intermediate system within the fabric has been determined, it MUST be advertised using the 4 bit Tier field described in section 3.3 of [I-D.shen-isis-spine-leaf-ext]. This section describes a method of calculating the tier number, assuming the tier numbers rise in value from the edge of the fabric.

This method begins with two of the T0 intermediate systems advertising their location in the fabric. This information can either be obtained through:

- o Two T0 intermediate systems are manually configured to advertise 0x00 in their IS reachability tier sub-TLV, indicating they are at the edge of the fabric (a ToR IS).
- o The T0 intermediate systems detect they are T0 through the presence connected hosts (i.e. through a request for address assignment or some other means). If such detection is used, and the IS determines it is located at T0, it should advertise 0x00 in its IS reachability tier sub-TLV.

If the first method is used, the two T0 routers MUST be "maximally separated" on the fabric. They must be a maximal number of hops apart, or rather they MUST NOT be connected to the same T1 device as their "upstream" towards the superspines in a 5 ary fabric.

The second method above SHOULD be used with care, as it may not be secure, and it may not work in all data center environments. For instance, if a host is mistakenly (or intentionally, as a form of attack) attached to a spine IS, or a request for address assignment is transmitted to a spine IS during the bootup phase of the device or fabric, it is possible to cause a spine IS to advertise itself as a T0. Unless the autodetection of the T0 devices is secured, the manual mechanism SHOULD BE used (configuring at least one T0 device manually).

Given the correct configuration of two T0 devices, maximally spaced on the fabric, the remaining intermediate systems calculate their tier number as follows:

- o The local IS calculates an SPT (using SPF) setting the cost of every link to 1; this effectively calculates a topology only view of the network, without considering any configured link costs
- o Ensure that at least two T0 are in the calculated SPT; otherwise abort
- o Find the furthest T0; call this node A and set LD to the cost; the "furthest T0" is the T0 with the largest metric, or the furthest distance from the local calculating node
- o Calculate an SPT (using SPF) from the perspective of A (above) setting the cost of every link to 1
- o Find the furthest IS in A's SPT; call this node B and set RD to the cost from A to B
- o Calculate the tier number of the local IS by subtracting LD from RD

In the example network, assume 5A and 1C are manually configured as a T0, and are advertising their tier numbers. From here:

- o From 1A the path to 5A is 4 hops; this is LD
- o Run SPF from the perspective of 5A with all link metrics set to 1
- o From 5A the path length to 1C is 4; this is RD
- o $RD - LD$ is 0 at 1A, so 1A is T0, or a ToR

This process will work for any spine and leaf fabric without "cross links."

5. Flooding Optimization

Flooding is perhaps the most challenging scaling issue for a link state protocol running on a dense, large scale fabric. To reduce the flooding of link state information in the form of Link State Protocol Data Units (LSPs), Openfabric takes advantage of information already available in the link state protocol, the list of the local intermediate system's neighbor's neighbors, and the fabric locality computed above. The following tables are required to compute a set of reflooders:

- o Neighbor List (NL) list: The set of neighbors

- o Neighbor's Neighbors (NN) list: The set of neighbor's neighbors; this can be calculated by running SPF truncated to two hops
- o Do Not Reflood (DNR) list: The set of neighbors who should have LSPs (or fragments) who should not reflood LSPs
- o Reflood (RF) list: The set of neighbors who should flood LSPs (or fragments) to their adjacent neighbors to ensure synchronization

NL is set to contain all neighbors, and sorted deterministically (for instance, from the highest IS identifier to the lowest). All intermediate systems within a single fabric SHOULD use the same mechanism for sorting the NL list. NN is set to contain all neighbor's neighbors, or all intermediate systems that are two hops away, as determined by performing a truncated SPF. The DNR and RF tables are initially empty. To begin, the following steps are taken to reduce the size of NN and NL:

- o Move any IS in NL with its tier (or fabric location) set to T0 to DNR
- o Remove all intermediate systems from NL and NN that in the shortest path to the IS that originated the LSP

Then, for every IS in NL:

- o If the current entry in NL is connected to any entries in NN:
 - * Move the IS to RF
 - * Remove the intermediate systems connected to the IS from NN
- o Else move the IS to DNR

The calculation terminates when the NL is empty.

When flooding, LSPs transmitted to adjacent neighbors on the RF list will be transmitted normally. Adjacent intermediate systems on this list will reflood received LSPs into the next stage of the topology, ensuring database synchronization. LSPs transmitted to adjacent neighbors on the DNR list, however, MUST be transmitted using a circuit scope PDU as described in [RFC7356].

5.1. Flooding Failures

It is possible in some failure modes for flooding to be incomplete because of the flooding optimizations outlined. Specifically, if a reflooder fails, or is somehow disconnected from all the links across

which it should be reflooding, it is possible an LSP is only partially flooded through the fabric. To prevent such situations, any IS receiving an LSP transmitted using DNR SHOULD:

- o Set a short timer; the default should be less than one second
- o When the timer expires, send a Complete Sequence Number Packet (CSNP) to all neighbors
- o Process any Partial Sequence Number Packets (PSNPs) as required to resynchronize
- o If a resynchronization is required, notify the network operator through a network management system

6. Other Optimizations

6.1. Transit Link Reachability

In order to reduce the amount of control plane state carried on large scale spine and leaf fabrics, openfabric implementations SHOULD NOT advertise reachability for transit links. These links MAY remain unnumbered, as IS-IS does not require layer 3 IP addresses to operate. Each IS SHOULD be configured with a single loopback address, which is assigned an IPv6 address, to provide reachability to intermediate systems which make up the fabric.

[RFC3277] SHOULD be supported on devices supporting openfabric with unnumbered interface in order to support traceability and network management.

6.2. Transiting T0 Intermediate Systems

In data center fabrics, ToR intermediate systems SHOULD NOT be used to transit between two T1 (or above) spine intermediate systems. The simplest way to prevent this is to set the overload bit [RFC3277] for all the LSPs originated from T0 intermediate systems. However, this solution would have the unfortunate side effect of causing all reachability beyond any T0 IS to have the same metric, and many implementations treat a set overload bit as a metric of 0xFFFF in calculating the Shortest Path Tree (SPT). This document proposes an alternate solution which preserves the leaf node metric, while still avoiding transiting T0 intermediate systems.

Specifically, all T0 intermediate systems SHOULD advertise their metric to reach any T1 adjacent neighbor with a cost of 0XFFE. T1 intermediate systems, on the other hand, will advertise T0 intermediate systems with the actual interface cost used to reach the

T0 IS. Hence, links connecting T0 and T1 intermediate systems will be advertised with an asymmetric cost that discourages transiting T0 intermediate systems, while leaving reachability to the destinations attached to T0 devices the same.

7. Openfabric and Route Aggregation

While schemes may be designed so reachability information can be aggregated in Openfabric deployments, this is not a recommended configuration.

8. Security Considerations

This document outlines modifications to the IS-IS protocol for operation on large scale data center fabrics. While it does add new TLVs, and some local processing changes, it does not add any new security vulnerabilities to the operation of IS-IS. However, openfabric implementations SHOULD implement IS-IS cryptographic authentication, as described in [RFC5304], and should enable other security measures in accordance with best common practices for the IS-IS protocol.

If T0 intermediate systems are auto-detected using information outside Openfabric, it is possible to attack the calculations used for flooding reduction and auto-configuration of intermediate systems. For instance, if a request for an address pool is used as an indicator of an attached host, and hence receiving such a request causes an intermediate system to advertise itself as T0, it is possible for an attacker (or a simple mistake) to cause auto-configuration to fail. Any such auto-detection mechanisms SHOULD BE secured using appropriate techniques, as described by any protocols or mechanisms used.

9. References

9.1. Normative References

[I-D.shen-isis-spine-leaf-ext]

Shen, N., Ginsberg, L., and S. Thyamagundalu, "IS-IS Routing for Spine-Leaf Topology", draft-shen-isis-spine-leaf-ext-07 (work in progress), October 2018.

- [ISO10589] International Organization for Standardization, "Intermediate system to Intermediate system intra-domain routeing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)", ISO/IEC 10589:2002, Second Edition, Nov 2002.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2629] Rose, M., "Writing I-Ds and RFCs using XML", RFC 2629, DOI 10.17487/RFC2629, June 1999, <<https://www.rfc-editor.org/info/rfc2629>>.
- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<https://www.rfc-editor.org/info/rfc5120>>.
- [RFC5301] McPherson, D. and N. Shen, "Dynamic Hostname Exchange Mechanism for IS-IS", RFC 5301, DOI 10.17487/RFC5301, October 2008, <<https://www.rfc-editor.org/info/rfc5301>>.
- [RFC5303] Katz, D., Saluja, R., and D. Eastlake 3rd, "Three-Way Handshake for IS-IS Point-to-Point Adjacencies", RFC 5303, DOI 10.17487/RFC5303, October 2008, <<https://www.rfc-editor.org/info/rfc5303>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [RFC5308] Hopps, C., "Routing IPv6 with IS-IS", RFC 5308, DOI 10.17487/RFC5308, October 2008, <<https://www.rfc-editor.org/info/rfc5308>>.
- [RFC5309] Shen, N., Ed. and A. Zinin, Ed., "Point-to-Point Operation over LAN in Link State Routing Protocols", RFC 5309, DOI 10.17487/RFC5309, October 2008, <<https://www.rfc-editor.org/info/rfc5309>>.

- [RFC5311] McPherson, D., Ed., Ginsberg, L., Previdi, S., and M. Shand, "Simplified Extension of Link State PDU (LSP) Space for IS-IS", RFC 5311, DOI 10.17487/RFC5311, February 2009, <<https://www.rfc-editor.org/info/rfc5311>>.
- [RFC5316] Chen, M., Zhang, R., and X. Duan, "ISIS Extensions in Support of Inter-Autonomous System (AS) MPLS and GMPLS Traffic Engineering", RFC 5316, DOI 10.17487/RFC5316, December 2008, <<https://www.rfc-editor.org/info/rfc5316>>.
- [RFC7356] Ginsberg, L., Previdi, S., and Y. Yang, "IS-IS Flooding Scope Link State PDUs (LSPs)", RFC 7356, DOI 10.17487/RFC7356, September 2014, <<https://www.rfc-editor.org/info/rfc7356>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

9.2. Informative References

- [I-D.ietf-isis-segment-routing-extensions]
Previdi, S., Ginsberg, L., Filsfils, C., Bashandy, A., Gredler, H., Litkowski, S., Decraene, B., and J. Tantsura, "IS-IS Extensions for Segment Routing", draft-ietf-isis-segment-routing-extensions-19 (work in progress), July 2018.
- [I-D.ietf-spring-segment-routing]
Filsfils, C., Previdi, S., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", draft-ietf-spring-segment-routing-15 (work in progress), January 2018.
- [RFC3277] McPherson, D., "Intermediate System to Intermediate System (IS-IS) Transient Blackhole Avoidance", RFC 3277, DOI 10.17487/RFC3277, April 2002, <<https://www.rfc-editor.org/info/rfc3277>>.
- [RFC3719] Parker, J., Ed., "Recommendations for Interoperable Networks using Intermediate System to Intermediate System (IS-IS)", RFC 3719, DOI 10.17487/RFC3719, February 2004, <<https://www.rfc-editor.org/info/rfc3719>>.

- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC5440] Vasseur, JP., Ed. and JL. Le Roux, Ed., "Path Computation Element (PCE) Communication Protocol (PCEP)", RFC 5440, DOI 10.17487/RFC5440, March 2009, <<https://www.rfc-editor.org/info/rfc5440>>.
- [RFC5449] Baccelli, E., Jacquet, P., Nguyen, D., and T. Clausen, "OSPF Multipoint Relay (MPR) Extension for Ad Hoc Networks", RFC 5449, DOI 10.17487/RFC5449, February 2009, <<https://www.rfc-editor.org/info/rfc5449>>.
- [RFC5614] Ogier, R. and P. Spagnolo, "Mobile Ad Hoc Network (MANET) Extension of OSPF Using Connected Dominating Set (CDS) Flooding", RFC 5614, DOI 10.17487/RFC5614, August 2009, <<https://www.rfc-editor.org/info/rfc5614>>.
- [RFC5820] Roy, A., Ed. and M. Chandra, Ed., "Extensions to OSPF to Support Mobile Ad Hoc Networking", RFC 5820, DOI 10.17487/RFC5820, March 2010, <<https://www.rfc-editor.org/info/rfc5820>>.
- [RFC5837] Atlas, A., Ed., Bonica, R., Ed., Pignataro, C., Ed., Shen, N., and JR. Rivers, "Extending ICMP for Interface and Next-Hop Identification", RFC 5837, DOI 10.17487/RFC5837, April 2010, <<https://www.rfc-editor.org/info/rfc5837>>.
- [RFC6232] Wei, F., Qin, Y., Li, Z., Li, T., and J. Dong, "Purge Originator Identification TLV for IS-IS", RFC 6232, DOI 10.17487/RFC6232, May 2011, <<https://www.rfc-editor.org/info/rfc6232>>.
- [RFC7182] Herberg, U., Clausen, T., and C. Dearlove, "Integrity Check Value and Timestamp TLV Definitions for Mobile Ad Hoc Networks (MANETs)", RFC 7182, DOI 10.17487/RFC7182, April 2014, <<https://www.rfc-editor.org/info/rfc7182>>.
- [RFC7921] Atlas, A., Halpern, J., Hares, S., Ward, D., and T. Nadeau, "An Architecture for the Interface to the Routing System", RFC 7921, DOI 10.17487/RFC7921, June 2016, <<https://www.rfc-editor.org/info/rfc7921>>.

Appendix A. Flooding Optimization Operation

Recent testing has shown that flooding is largely a "non-issue" in terms of scaling when using high speed links connecting intermediate systems with reasonable processing power and memory. However, testing has also shown that flooding will impact convergence speed even in such environments, and flooding optimization has a major impact on the performance of a link state protocol in resource constrained environments. Some thoughts on flooding optimization in general, and the flooding optimization contained in this document, follow.

There are two general classes of flooding optimization available for link state protocols. The first class of optimization relies on a centralized service or server to gather the link state information and redistribute it back into the intermediate systems making up the fabric. Such solutions are attractive in many, but not all, environments; hence these systems compliment, rather than compete with, the system described here. Systems relying on a service or server necessarily also rely on connectivity to that service or server, either through an out-of-band network or connectivity through the fabric itself. Because of this, these mechanisms do not apply to all deployments; some deployments require underlying reachability regardless of connectivity to an outside service or server.

The second possibility is to create a fully distributed system that floods the minimal amount of information possible to every intermediate system. The system described in this draft is an example of such a system. Again, there are many ways to accomplish this goal, but simplicity is a primary goal of the system described in this draft.

The system described here divides the work into two different parts; forward and reverse optimization. The forward optimization begins by finding the set of intermediate systems two hops away from the flooding device, and choosing a subset of connected neighbors that will successfully reach this entire set of intermediate systems, as shown in the diagram below.

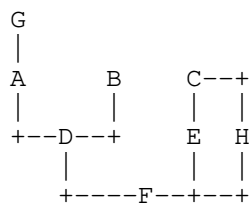


Figure 2

If F is flooding some piece of information, then it will find the entire set of intermediate systems within two hops by discovering its neighbors and their neighbors from the local LSDB. This will include A, B, C, D, and E--but not G. From this set, F can determine that D can reach A and B, while a single flood to either E or H will reach C. Hence F can flood to D and either E or H to reach C. F can choose to flood to D and E normally. Because H still needs to receive this new LSP (or fragment!), but does not need to reflood to C, F can send the LSP using link local signaling. In this case, H will receive and process the new LSP, but not reflood it.

Rather than carrying the information necessary through hello extensions, as is done in [RFC5820], the neighbors are allowed to complete initial synchronization, and then a truncated shortest path tree is built to determine the "two hop neighborhood." This has the advantage of using mechanisms already used in IS-IS, rather than adding new processes. The risk with this process is any LSPs flooded through the network before this initial calculation takes place will be suboptimal. This "two hop neighborhood" process has been used in OSPF deployments for a number of years, and has proven stable in practice.

Rather than setting a timer for reflooding, the implementation described here uses IS-IS' ability to describe the entire database using a CSNP to ensure flooding is successful. This adds some small amount of overhead, so there is some balance between optimal flooding and ensuring flooding is complete.

The reverse optimization is simpler. It relies on the observation that any intermediate system between the local IS and the origin of the LSP, other than in the case of floods removing an LSP from the shared LSDB, should have already received a copy of the LSP. For instance, if F originates an LSP in the figure above, and E refloods the LSP to C, C does not need to reflood back to F if F is on its shortest path tree towards F. It is obvious this is not a "perfect" optimization. A perfect optimization would block flooding back along a directed acyclic graph towards the originator. Using the SPT, however, is a quick way to reduce flooding without performing more calculations.

The combination of these two optimizations have been seen, in testing, to reduce the number of copies any IS receives from the tens to precisely one.

Appendix B. Fabric Location Calculation

Determining the location of a device in a symmetric topology is quite challenging. The authors of this draft worked through a number of possible solutions to this problem, each of which was found to either not work in some topology, or was found to be liable to unacceptable errors. For instance:

- o Method 1:

- * Caculate the maximum distance through the fabric, and the distance from one of those points to the local intermediate system
- * This works in a five stage Clos spine and leaf, but not in a three stage, nor in some other five stage spine and leaf fabrics, such as the common butterfly or Benes fabric

- o Method 2:

- * Manually mark one edge leaf node in the fabric as T0
- * Calculate maximum distance through the fabric from this point
- * Calculate local position based on this maximum distance the distance to the single marked device
- * This works in three and five stage Clod fabrics, but does not work from every location in other spine and leaf fabrics, such as the common butterfly or Benes fabric

In the end, marking two devices located as far from one another topologically as possible provides the anchor points necessary to calculate the total distance through the fabric, and then from those points to the location of the calculating device.

The information obtained in this way can also be combined with other forms of location calculation, such as whether a device requesting an address through some mechanism is attached to the local device, or other indications of fabric locality. It generally true that having more than one method to determine fabric location will be better than any single method to account for errors, failures, and other problems that can arise with any mechanism.

Authors' Addresses

Russ White (editor)
LinkedIn

Email: russ@riw.us

Shawn Zandi (editor)
LinkedIn

Email: szandi@linkedin.com