```
Workgroup: Network Working Group
Internet-Draft:
draft-chen-lsvr-flood-reduction-00
Published: 22 October 2021
Intended Status: Standards Track
Expires: 25 April 2022
Authors: H. Chen
                  G. Mishra
                                 A. Wang
        Futurewei Verizon Inc.
                                  China Telecom
        Y. Liu
                                Y. Fan
                       H. Wang
        China Mobile Huawei
                                Casa Systems
                      BGP-SPF Flooding Reduction
```

### Abstract

This document describes extensions to Border Gateway Protocol (BGP) for flooding the link states on a topology that is a subgraph of the complete topology of a BGP-SPF domain, so that the amount of flooding traffic in the domain is greatly reduced. This would reduce convergence time with a more stable and optimized routing environment.

## **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 <u>RFC 2119</u> [<u>RFC2119</u>].

## 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 <u>https://datatracker.ietf.org/drafts/current/</u>.

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 25 April 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

(<u>https://trustee.ietf.org/license-info</u>) 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

- <u>1</u>. <u>Introduction</u>
- <u>2</u>. <u>Terminologies</u>
- 3. Overview of BGP-SPF Link State Flooding
  - 3.1. Flooding in RR Model
  - 3.2. Flooding in Node Connections Model
  - 3.3. Flooding in Directly-Connected Nodes Model
- <u>4. Revised Flooding Procedures</u>
  - 4.1. Revised Flooding Procedure for RR Model
  - 4.2. <u>Revised Flooding Procedure for Node Connections Model</u>
- 5. BGP Extensions for Flooding Reduction
  - 5.1. Extensions for RR Model
  - 5.2. Extensions for Node Connections Model
    - <u>5.2.1</u>. <u>New TLVs</u>
    - 5.2.2. Flooding Topology Distribution in Centralized Mode
    - 5.2.3. An Algorithm for Distributed Mode
- <u>6</u>. <u>Security Considerations</u>
- <u>7</u>. <u>Acknowledgements</u>
- <u>8</u>. <u>IANA Considerations</u>
- <u>9</u>. <u>References</u>
  - <u>9.1</u>. <u>Normative References</u>
  - 9.2. Informative References

<u>Authors' Addresses</u>

## 1. Introduction

For some networks such as dense Data Center (DC) networks with BGP-SPF, the existing Link State (LS) flooding mechanism defined in [<u>I-D.ietf-lsvr-bgp-spf</u>] for a BGP-SPF domain may not be 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 describes extensions to Border Gateway Protocol (BGP) for flooding the link states on a topology that is a subgraph of the

complete topology of a BGP-SPF domain, so that the amount of flooding traffic in the domain is greatly reduced.

## 2. Terminologies

The following terminologies are used in this document.

**BGP:** Border Gateway Protocol

LS: Link State

**SPF:** Shortest Path First

RR: Route Reflector

## 3. Overview of BGP-SPF Link State Flooding

[<u>I-D.ietf-lsvr-bgp-spf</u>] defines three BGP peering models:

\*BGP Peering in Route-Reflector or Controller Topology (RR or Sparse model for short).

\*BGP Single-Hop Peering on Network Node Connections (Node Connections model for short), and

\*BGP Peering Between Directly-Connected Nodes (Directly-Connected Nodes model for short).

This section briefs the BGP-SPF Link State Flooding in each of these models.

#### 3.1. Flooding in RR Model

In RR model, BGP-SPF speakers/nodes peer solely with one or more Route Reflectors (RRs) or controllers. A BGP-SPF speaker sends/ advertises its BGP-LS-SPF Link NLRI in a BGP update message to the RRs or controllers that the speaker peers with when it discovers that its corresponding link is up. After receiving the Link NLRI, each of the RRs or controllers sends the NLRI in a BGP update message to the other BGP-SPF speakers that peer with the RRs or controllers.

For example, <u>Figure 1</u> shows a BGP-SPF domain, which contains two RRs RR1 and RR2, and three network nodes A, B and C. RR1 peers with all three nodes A, B and C in the network. RR2 also peers with all three nodes A, B and C in the network. There is a link between A and B, a link between A and C, and a link between B and C.



Figure 1: BGP-SPF Domain with two RRs

Each of the nodes A, B and C in the network sends/advertises its link NLRIs in BGP update messages to both RR1 and RR2. After receiving a link NLRI in a BGP update message from a node (e.g., node A), each of RR1 and RR2 sends the NLRI in a BGP update message to the other nodes (e.g., nodes B and C). Each of the other nodes receives two copies of the same NLRI, one from RR1 and the other from RR2. One copy is enough, the other redundant copy should be reduced.

## 3.2. Flooding in Node Connections Model

In Node Connections model, EBGP single-hop sessions are established over direct point-to-point links interconnecting the nodes in the BGP-SPF routing domain. Once the session has been established and the BGP-LS-SPF AFI/SAFI capability has been exchanged for the corresponding session, then the link is considered up from a BGP-SPF perspective and the corresponding BGP-LS-SPF Link NLRI is advertised to all the nodes in the domain through all the BGP sessions over the links. If the session goes down, the corresponding Link NLRI will be withdrawn. The withdrawal is done through advertising a BGP update containing the NLRI in MP\_UNREACH\_NLRI to all the nodes in the domain using all BGP sessions over the links.

For example, <u>Figure 2</u> shows a BGP-SPF domain, which contains four nodes A, B, C and D. These four nodes are connected by six links. There are two parallel links between A and B, a link between A and C, a link between A and D, a link between B and C and a link between C and D.



Figure 2: BGP-SPF Domain with parallel links

Suppose that the BGP sessions over all the links except for the session over the link between A and D have been established and the BGP-LS-SPF AFI/SAFI capability has been exchanged for the corresponding sessions. When the BGP session over the link between A and D is established and the BGP-LS-SPF AFI/SAFI capability is exchanged for the corresponding session, node A considers that the link from A to D is up and sends the BGP-LS-SPF Link NLRI for the link through its four BGP sessions (i.e., the session between A and B over the first parallel link between A and B, the session between A and B over the second parallel link between A and B, the session between A and C over the link between A and C, and the session between A and D over the link between A and D) to nodes B, C and D. After receiving the NLRI from node A, each of the nodes B, C and D sends the NLRI to the other nodes that have BGP sessions with the node. Node B sends the NLRI to node C. Node C sends the NLRI to nodes B and D. Node D sends the NLRI to node C.

Similarly, when the BGP session over the link between A and D is established and the BGP-LS-SPF AFI/SAFI capability is exchanged for the corresponding session, node D considers that the link from D to A is up and sends the BGP-LS-SPF Link NLRI for the link through its two BGP sessions (i.e., the session between D and C over the link between D and C, and the session between D and A over the link between D and A) to nodes C and A. After receiving the NLRI from node D, each of the nodes A and C sends the NLRI to the other nodes that have BGP sessions with the node. Node C sends the NLRI to nodes A and B. Node A sends the NLRI to nodes B and C through two parallel BGP sessions to B and the BGP session to C.

### 3.3. Flooding in Directly-Connected Nodes Model

In Directly-Connected Nodes model, BGP-SPF speakers peer with all directly-connected nodes but the sessions may be between loopback addresses. Consequently, there will be a single BGP session even if there are multiple direct connections between BGP-SPF speakers. BGP- LS-SPF Link NLRI is advertised as long as a BGP session has been established, the BGP-LS-SPF AFI/SAFI capability has been exchanged. Since there are BGP sessions between every directly-connected nodes in the BGP-SPF routing domain, there is only a reduction in BGP sessions when there are parallel links between nodes comparing to node connections model.

### 4. Revised Flooding Procedures

#### 4.1. Revised Flooding Procedure for RR Model

- In RR model, the revised flooding procedure is as follows:
  - \*A BGP-SPF speaker/node sends its BGP-LS-SPF Link NLRI to some such as one of the RRs or controllers that the speaker peers with when it discovers that its corresponding link is up.
  - \*After receiving the Link NLRI, the RR or controller sends the NLRI to the other BGP-SPF speakers that peer with the RR or controller.

For example, for the BGP-SPF domain in Figure 1, using the revised flooding procedure, speaker/Node A sends its Link NLRI for link A to B to one RR RR1 when A discovers that link A to B is up. Node A does not send the NLRI to RR2. After receiving the Link NLRI for link A to B from speaker/node A, RR1 sends the NLRI to the other nodes B and C. Each of the other nodes receives only one copy of the same NLRI, which is from RR1. There is no redundant copy of the same NLRI. Comparing to the normal flooding in RR model as illustrated in Figure 1, the revised flooding procedure reduces the amount of link states flooding by half.

In an option, for a number of RRs or controllers that peer with all the nodes/speakers in a network, the nodes are evenly divided into the number of groups. A first group of nodes send their link NLRIs to a first RR; a second group of nodes send their link NLRIs to a second RR; and so on. After receiving a NLRI from a node, a RR sends the NLRI to the other nodes in the network. This option may be used if each node peers with every RR or controller; otherwise, it should not be used.

In one implementation, the nodes (supposing there are m nodes in total) are divided into N groups through ordering the nodes by their IDs in ascending order and grouping the nodes. Each of the N groups has m/N nodes. The first m/N nodes in the ordered nodes are in the first group; the m/N nodes following the first group are in the second group; the m/N nodes following the second group are in the third group; and so on. The nodes following the second last group are in the N-th group (i.e., the last group).

For example, for the BGP-SPF domain in <u>Figure 1</u>, there are two RRs and three nodes, the nodes in the network are evenly divided into two groups. The first group contains one (3/2 = 1) node: node A. The second group contains the rest nodes: nodes B and C.

Node A in the first group sends its link NLRIs to RR1. After receiving a Link NLRI from node A, RR1 sends the NLRI to the other nodes B and C in the network. Nodes B and C in the second group send their link NLRIs to RR2. After receiving a Link NLRI from node B, RR2 sends the NLRI to the other nodes A and C in the network. After receiving a Link NLRI from node C, RR2 sends the NLRI to the other nodes A and B in the network.

Each of the other nodes receives only one copy of the same NLRI, which is from RR1 or RR2. There is no redundant copy of the same NLRI.

In this option, every group of nodes has about the same number of nodes as each of the other groups, the workload is balanced among the RRs (i.e., each of RRs has almost the same workload as any other RR).

In another option, for a number of RRs or controllers that peer with all the nodes/speakers in a network, the nodes in the network sends their link NLRIs to the same one or more of the RRs.

For example, for the BGP-SPF domain in Figure 1, nodes A, B and C in the network send their link NLRIS to the same RR1. After receiving the Link NLRI from a node, RR1 sends the NLRI to the other nodes in the network. For example, after receiving the Link NLRI from node A, RR1 sends the NLRI to the other nodes B and C in the network. After receiving the Link NLRI from node B, RR1 sends the NLRI to the other nodes A and C in the network. After receiving the Link NLRI from node C, RR1 sends the NLRI to the other nodes A and B in the network.

### 4.2. Revised Flooding Procedure for Node Connections Model

In Node Connections model, the revised flooding procedure is as follows:

\*A BGP-SPF speaker/node has a flooding topology of the BGP-SPF domain. In an option, the flooding topology is computed in a distributed mode, where every BGP-SPF speaker computes a flooding topology for the domain using a same algorithm. In another option, the flooding topology is computed in a centralized mode, where one BGP-SPF speaker elected as a leader computes a flooding topology for the domain and advertises the flooding topology to every BGP-SPF speaker in the domain. \*A BGP-SPF speaker/node sends its link NLRI in a BGP update message for its link up or down to its peers that are directly connected on the flooding topology, and sends its link NLRI in a BGP update message for its link down to all its peers. When receiving the NLRI in a new BGP update message for a link up or down from a peer, the speaker sends the NLRI in a BGP update message to its other peers that are directly connected on the flooding topology.

\*When a BGP-SPF session is down, the BGP-SPF speaker/node that was connected to the session will not withdraw the link NLRIs received from the session right away. It keeps the NLRIs for some time.

Given a real network topology (RT), a flooding topology (FT) of the RT is a sub network topology of the RT and connects all the nodes in the RT.

For example, <u>Figure 3</u> shows a flooding topology of the real topology in <u>Figure 2</u>.



Figure 3: A Flooding Topology

The flooding topology in <u>Figure 3</u> is a sub network topology of the RT in <u>Figure 2</u> and connects all the nodes (i.e., nodes A, B, C and D) in the RT in <u>Figure 2</u>.

Figure 4 shows a reduced flooding flow of a link NLRI in a BGP update message for a link up or down in the BGP-SPF domain, which is the same as the one in Figure 2.



Figure 4: A Reduced Link State Flooding Flow

Speaker/Node A sends the NLRI in a BGP update message for its link to its peers B and D. Nodes B and D are peers of node A and are directly connected to A on the flooding topology (FT). Node A does not send the NLRI to its peer C since C is not directly connected to A on the FT.

After receiving the NLRI in the message from A, node B sends the NLRI in a BGP update message to B's other peer C (which is directly connected to B on the FT). After receiving the NLRI in a BGP update message from A, node D sends the NLRI in a BGP update message to D's other peer C (which is directly connected to D on the FT).

The number of NLRIs in messages flooded in <u>Figure 4</u> is much less than that in <u>Figure 2</u>. The performance of network is improved using the revised flooding procedure.

## 5. BGP Extensions for Flooding Reduction

This section specifies BGP extensions for flooding reduction in two models: RR model and Node Connections model. The extensions for Directly-Connected Node model are included in the extensions for Node Connections model.

## 5.1. Extensions for RR Model

A single RR for a BGP-SPF domain is elected as a leader RR of the domain. The leader RR is the RR with the highest priority to become a leader in the domain. If there are more than one RRs having the same highest priority, the RR with the highest Node ID and the highest priority is the leader RR in the domain. In a deployment, only every RR advertises its priority for becoming a leader using a Leader Priority TLV defined below.

Two new TLVs are defined for flooding reduction in RR model.

\*Leader Priority TLV: A node uses it to advertise its priority for becoming a leader.

\*Node Flood TLV: A RR or controller uses it to tell every node the flooding behavior the node needs to follow.

The format of Leader Priority TLV is illustrated in Figure 5.

Figure 5: Leader Priority TLV

Type: It is to be assigned by IANA.

Length: 4.

**Reserved:** MUST be set to zero in transmission and should be ignored on reception.

**Priority:** A unsigned integer from 0 to 255 in one octet indicating priority to become a leader.

The format of Node Flood TLV is illustrated in Figure 6.

Figure 6: Node Flood TLV

Type: It is to be assigned by IANA.

Length: 4.

**Reserved:** MUST be set to zero in transmission and should be ignored on reception.

#### Flood-behavior:

The following flooding behavior are defined.

0 - Reserved.

- 1 send link states to the RR with the minimum ID
- 2 send link states to the RR with the maximum ID
- 3 balanced groups
- 4 send link states to 2 RRs with smaller IDs
- 5 send link states to 2 RRs with larger IDs
- 6 balanced groups with redundancy of 2
- 7-127 Standardized flooding behaviors for RR Model

128-254 - Private flooding behaviors for RR Model.

In a deployment, the flooding behavior for every node is configured on a RR or controller such as the leader RR and the RR advertises the behavior to the other RRs and every node in the network though using a Node Flood TLV.

For example, if we want every node in the network to send its link states to only one RR, we configure this behavior on a RR and the RR advertises the behavior to every node using a Node Flood TLV with Flood-behavior set to one, which tells every node to send its link states to the RR with the minimum ID. If we want every node in the network to send its link states to two RRs for redundancy, we configure this behavior on a RR and the RR advertises the behavior to every node using a Node Flood TLV with Flood-behavior set to 4, which tells every node to send its link states to the two RRs with smaller IDs (i.e., the RR with the minimum ID and the RR with the second minimum ID).

If we want to balance the traffic among RRs or controllers through dividing the nodes into groups and letting each group send their link states to a RR, we configure this behavior on a RR and the RR advertises the behavior to every node using a Node Flood TLV with Flood-behavior set to 3, which tells every node to divide the nodes in the network into a number of groups. A node in a group sends its link states to the RR corresponding to the group.

### 5.2. Extensions for Node Connections Model

There are two modes for the flooding topology computation: centralized mode and distributed mode. In a centralized mode, one BGP-SPF node is elected as a leader. The leader computes a flooding topology for the BGP-SPF domain and advertises the flooding topology to every BGP-SPF node in the domain. In a distributed mode, every BGP-SPF node computes a flooding topology for the BGP-SPF domain using a same algorithm. There is not any flooding topology distribution. This section defines the new TLVs for the two modes, describes the flooding topology distribution in centralized mode and an algorithm that can be used by every node to compute its flooding topology in distributed mode.

## 5.2.1. New TLVs

Five new TLVs are defined for flooding reduction in Node Connections model.

- \*Node Algorithm TLV: A leader uses this TLV to tell every node the algorithm to be used to compute a flooding topology.
- \*Algorithms Support TLV: A node uses this TLV to indicate the algorithms that it supports for distributed mode.
- \*Node IDs TLV: A leader uses this TLV to indicate the mapping from nodes to their indices for centralized mode.
- \*Paths TLV: A leader uses this TLV to advertise a part of flooding topology for centralized mode.

\*Connection Used for Flooding TLV: A node uses this TLV to indicate that a connection/link is a part of the flooding topology and used for flooding.

# 5.2.1.1. Node Algorithm TLV

The format of Node Algorithm TLV is illustrated in Figure 7.

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
+	+ - +	+	+	+ - +	+ - +	+	+ - +	+ - +	+ - +	+ - +	+ - +	+	+	+ - +	+ - +	+ - +	+ - +	+ - +		+ - +	+ - +	+	+	+ - +	+ - +	+ - +	+	+ - +			+-+
L	Type = TBD3								Lei						eng	ngth = 4					I										
+-																															
I											Re	ese	er۱	/ec	k										ŀ	410	goi	rit	hn	n	
+-																															

Figure 7: Node Algorithm TLV

Type: It is to be assigned by IANA.

Length: 4.

**Reserved:** MUST be set to zero in transmission and should be ignored on reception.

Algorithm:

- O The leader computes a flooding topology using its own algorithm and advertises the flooding topology to every node.
- 1-127 Every node computes its flooding topology using this standardized distributed algorithm.
- 128-254 Private distributed algorithms.

A node such as the leader node can use this TLV to tell every node in the domain to use the flooding topology from the leader for flooding the link states through advertising the TLV with the Algorithm field set to zero, or to tell every node to compute its own flooding topology using the algorithm given by the Algorithm field in the TLV containing an identifier of an algorithm when the Algorithm field is not zero.

## 5.2.1.2. Algorithms Support TLV

The format of Algorithms Support TLV is illustrated in Figure 8.

Figure 8: Algorithms Support TLV

Type: It is to be assigned by IANA.

Length: The number of Algorithms in the TLV.

**Algorithm:** A numeric identifier in the range 0-255 indicating the algorithm that can be used to compute the flooding topology.

# 5.2.1.3. Node IDs TLV

The format of Node IDs TLV is illustrated in Figure 9.

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 = TBD5Length (variable) 1 Reserved |L| Starting Index Node ID . . . . . . Node ID 

Figure 9: Node IDs TLV

Type: It is to be assigned by IANA.

**Length:** 4 \* (number of Node IDs + 1).

- **Reserved:** MUST be set to zero in transmission and should be ignored on reception.
- L: This bit is set to one if the index of the last node ID in this TLV is equal to the last index in the full list of node IDs for the BFP-SPF domain.
- Starting Index: The index of the first node ID in this TLV is
   Starting Index; the index of the second node ID in this TLV is
   Starting Index + 1; the index of the third node ID in this TLV is
   Starting Index + 2; and so on.

Node ID: The BGP identifier of a node in the BGP-SPF domain.

## 5.2.1.4. Paths TLV

The format of Paths TLV is illustrated in Figure 10. A leader uses this TLV to advertise a part of flooding topology for centralized mode. A path may be described as a sequence of indices: (Index 1, Index 2, Index 3, ...), denoting a connection between the node with index 1 and the node with index 2, a connection between the node with index 2 and the node with index 3, and so on. A single link/ connection is a simple case of a path that only connects two nodes. A single link path may be encoded in a paths TLV of 8 bytes with two indices.

Θ	1	2	3						
012	3 4 5 6 7 8 9 0 1 2 3 4	5 6 7 8 9 0 1 2 3 4	5678901						
+ - + - + - +	-+	-+	+ - + - + - + - + - + - + - +						
	Type = TBD6	Length	(variable)						
+-									
1	Index 1	Index	< 2						
+-									
~			~						
+-+-+	-+	-+	+-+-+-+-+-+-+						

Figure 10: Paths TLV

Type: It is to be assigned by IANA.

**Length:** 2 \* (number of indices in the path) when the TLV contains the indices for one path.

**Index 1:** The index of the first node in the path.

Index 2: The index of the second (next) node in the path.

Multiple such as N paths may be encoded in one paths TLV. Each of the multiple paths is represented as a sequence of indices of the nodes on the path, and two paths (i.e., two sequences of indices for the two paths) are separated by a special index value such as 0xFFFF. In this case, there are (N - 1) special indices as separators to separate N paths, and the Length field has a value of 2 \* (number of indices in N paths + N - 1).

When there are a number such as N of single link paths, using one paths TLV to represent them is more efficient than using N paths TLVs to represent them (i.e., each paths TLV represents a single link path). Using one TLV consumes 4 + 2 \* (2\*N + N - 1) = 6\*N + 2bytes. Using N TLVs occupies N \* (4 + 4) = 8\*N bytes. The space used by the former is about three quarters of the space used by the latter for a big N such as 30.

### 5.2.1.5. Connection Used for Flooding TLV

The format of Connection Used for Flooding TLV is illustrated in Figure 11.

Θ	1	2	3						
0123	4 5 6 7 8 9 0 1 2 3 4	4 5 6 7 8 9 0 1 2	3 4 5 6 7 8 9 0 1						
+ - + - + - + -	+ - + - + - + - + - + - + - + - + - + -	-+-+-+-+-+-+-+-+	·-+-+-+-+-+-+-+-+-+						
	Type = TBD7	Leng	jth = 8						
+-									
Local Node ID									
+-									
Remote Node ID									
+-									

Figure 11: Connection Used for Flooding TLV

Type: It is to be assigned by IANA.

Length: 8.

- **Local Node ID:** The BGP ID of the local node of the session over the connection on the flooding topology which is used for flooding link states.
- **Remote Node ID:** The BGP ID of the remote node of the session over the connection on the flooding topology which is used for flooding link states.

## 5.2.2. Flooding Topology Distribution in Centralized Mode

In centralized mode, the leader computes a flooding topology for the domain whenever there is a change in the real network topology of the domain and advertises the flooding topology to every node in the domain.

After the current leader has failed, a new leader is elected. The new leader computes a flooding topology for the domain and advertises the flooding topology to every node in the domain.

For a brand new flooding topology of the domain computed, the leader advertises the whole flooding topology to every node in the domain. The leader advertises the mappings between all the node IDs and their indices to every node in the domain using a number of node IDs TLVs first. These node IDs TLVs contain the IDs of all the nodes in the domain and indicates the index corresponding to each of the node IDs and are advertised under MP\_REACH\_NLRI in BGP update messages. And then the leader advertises the connections/links on the flooding topology to every node in the domain using a number of paths TLVs. These paths TLVs contain all the connections/links on the flooding topology and are advertised under MP\_REACH\_NLRI in BGP update messages.

After advertising a flooding topology to every node in the domain, which is called the current flooding topology, for a new flooding

topology computed for the updated real network topology of the domain, the leader advertises only the changes in the new flooding topology comparing to the current flooding topology to every node in the domain. The leader advertises the changes in the mappings between all the node IDs and their indices to every node in the domain using node IDs TLVs first, and then advertises the changes in the flooding topology to every node in the domain using paths TLVs.

For the new nodes added into the domain, the leader advertises the mappings between the IDs of the new nodes and their indices using a node IDs TLV under MP\_REACH\_NLRI in a BGP update message to add the mappings. For the dead nodes removed from the domain, the leader advertises the mappings between the IDs of the dead nodes and their indices using a node IDs TLV under MP\_UNREACH\_NLRI in a BGP update message to withdraw the mappings.

For the new connections/links added into the current flooding topology, the leader advertises the new connections/links using a paths TLV under MP\_REACH\_NLRI in a BGP update message to add the new connections/inks to the current flooding topology. For the old connections/links removed from the current flooding topology, the leader advertises the old connections/links using a paths TLV under MP\_UNREACH\_NLRI in a BGP update message to withdraw the old connections/links from the current flooding topology.

### 5.2.3. An Algorithm for Distributed Mode

This section specifies an algorithm that can be used by every node to compute its flooding topology.

The algorithm for computing a flooding topology of a BGP-SPF domain (real topology) is described as follows.

\*Select a node R0 with the smallest node ID and without the status indicating that the node does not support transit;

\*Build a tree using R0 as root of the tree (details below);

\*And then connect a leaf to the tree to have a flooding topology (details follow).

The algorithm starts from

\*a variable MaxD with an initial value 3,

\*an initial flooding topology FT = {(R0, D=0, PHs={})} with node R0 as root, where R0's Degree D = 0, Previous Hops PHs = { };

\*an initial candidate queue Cq =  $\{(R1, D=0, PHs=\{R0\}), (R2, D=0, PHs=\{R0\}), \ldots, (Rm, D=0, PHs=\{R0\})\}$ , where each of nodes R1 to Rm

is connected to R0, its Degree D = 0 and Previous Hops PHs ={R0}, R1 to Rm are in increasing order by their IDs.

- Find and remove the first element with node A from Cq that is not on FT and one PH's D in PHs < MaxD, and add the element with A into FT; Set A's D to one, increase A's PH's D by one. If no element in Cq satisfies the conditions, algorithm is restarted with ++MaxD, the initial FT and Cq.
- 2. If all the nodes are on the FT, then goto step 4;
- 3. Suppose that node Xi (i = 1, 2, ..., n) is connected to node A and not on FT, and X1, X2, ..., Xn are in increasing order by their IDs (i.e., X1's ID < X2's ID < ... < Xn's ID). If they are not ordered, then make them in the order. If Xi is not in Cq, then add it into the end of Cq with D = 0 and PHs = {A}; otherwise (i.e., Xi is in Cq), add A into the end of Xi's PHs; Goto step 1.</p>
- 4. For each node B on FT whose D is one (from minimum to maximum node ID), find a link L attached to B such that L's remote node R can transit traffic and has minimum D and ID (if there is no node R which can transit traffic, then find a link L to node R whose D and ID are minimum), add link L between B and R into FT and increase B's D and R's D by one. Return FT.

### 6. Security Considerations

TBD

## 7. Acknowledgements

The authors of this document would like to thank Donald E. Eastlake for the comments.

#### 8. IANA Considerations

TBD

## 9. References

### 9.1. Normative References

[I-D.ietf-lsvr-bgp-spf] Patel, K., Lindem, A., Zandi, S., and W. Henderickx, "BGP Link-State Shortest Path First (SPF) Routing", Work in Progress, Internet-Draft, draft-ietflsvr-bgp-spf-15, 1 July 2021, <<u>https://www.ietf.org/</u> archive/id/draft-ietf-lsvr-bgp-spf-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, <<u>https://www.rfc-editor.org/info/</u> rfc2119>.
- [RFC4721] Perkins, C., Calhoun, P., and J. Bharatia, "Mobile IPv4 Challenge/Response Extensions (Revised)", RFC 4721, DOI 10.17487/RFC4721, January 2007, <<u>https://www.rfc-</u> editor.org/info/rfc4721>.
- [RFC4760] Bates, T., Chandra, R., Katz, D., and Y. Rekhter, "Multiprotocol Extensions for BGP-4", RFC 4760, DOI 10.17487/RFC4760, January 2007, <<u>https://www.rfc-</u> editor.org/info/rfc4760>.

# 9.2. Informative References

#### [I-D.ietf-lsr-dynamic-flooding]

Li, T., Psenak, P., Ginsberg, L., Chen, H., Przygienda, T., Cooper, D., Jalil, L., Dontula, S., and G. S. Mishra, "Dynamic Flooding on Dense Graphs", Work in Progress, Internet-Draft, draft-ietf-lsr-dynamic-flooding-09, 9 June 2021, <<u>https://www.ietf.org/archive/id/draft-ietf-lsr-dynamic-flooding-09.txt</u>>.

### [I-D.ietf-lsr-flooding-topo-min-degree]

Chen, H., Toy, M., Yang, Y., Wang, A., Liu, X., Fan, Y., and L. Liu, "Flooding Topology Minimum Degree Algorithm", Work in Progress, Internet-Draft, draft-ietf-lsrflooding-topo-min-degree-02, 1 June 2021, <<u>https://</u> www.ietf.org/archive/id/draft-ietf-lsr-flooding-topo-mindegree-02.txt>.

## Authors' Addresses

Huaimo Chen Futurewei Boston, MA, United States of America

Email: <u>huaimo.chen@futurewei.com</u>

Gyan S. Mishra Verizon Inc. 13101 Columbia Pike Silver Spring, MD 20904 United States of America Phone: <u>301 502-1347</u> Email: gyan.s.mishra@verizon.com Aijun Wang China Telecom Beiqijia Town, Changping District Beijing 102209 China Email: wangaj3@chinatelecom.cn Yisong Liu China Mobile Email: liuyisong@chinamobile.com Haibo Wang Huawei Huawei Bld., No.156 Beiqing Rd. Beijing 100095 China Email: <a href="mailto:rainsword.wang@huawei.com">rainsword.wang@huawei.com</a> Yanhe Fan Casa Systems

Email: yfan@casa-systems.com

United States of America