
Stream:	Internet Engineering Task Force (IETF)			
RFC:	9815			
Category:	Standards Track			
Published:	July 2025			
ISSN:	2070-1721			
Authors:	K. Patel	A. Lindem	S. Zandi	W. Henderickx
	<i>Arrcus, Inc.</i>	<i>LabN Consulting, LLC</i>	<i>LinkedIn</i>	<i>Nokia</i>

RFC 9815

BGP Link-State Shortest Path First (SPF) Routing

Abstract

Many Massively Scaled Data Centers (MSDCs) have converged on simplified Layer 3 (L3) routing. Furthermore, requirements for operational simplicity have led many of these MSDCs to converge on BGP as their single routing protocol for both fabric routing and Data Center Interconnect (DCI) routing. This document describes extensions to BGP for use with BGP - Link State (BGP-LS) distribution and the Shortest Path First (SPF) algorithm. In doing this, it allows BGP to be efficiently used as both the underlay protocol and the overlay protocol in MSDCs.

Status of This Memo

This is an Internet Standards Track document.

This document is a product of the Internet Engineering Task Force (IETF). It represents the consensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Further information on Internet Standards is available in Section 2 of RFC 7841.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <https://www.rfc-editor.org/info/rfc9815>.

Copyright Notice

Copyright (c) 2025 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	4
1.1. Terminology	4
1.2. BGP Shortest Path First (SPF) Motivation	5
1.3. Document Overview	6
1.4. Requirements Language	6
2. Base BGP Protocol Relationship	6
3. BGP - Link State (BGP-LS) Relationship	7
4. BGP SPF Peering Models	7
4.1. BGP Single-Hop Peering on Network Node Connections	7
4.2. BGP Peering Between Directly Connected Nodes	8
4.3. BGP Peering in Route-Reflector or Controller Topology	8
5. BGP Shortest Path Routing (SPF) Protocol Extensions	9
5.1. BGP-LS Shortest Path Routing (SPF) SAFI	9
5.1.1. BGP-LS-SPF NLRI TLVs	9
5.1.2. BGP-LS Attribute	9
5.2. Extensions to BGP-LS	10
5.2.1. Node NLRI Usage	10
5.2.1.1. BGP-LS-SPF Node NLRI Attribute SPF Status TLV	10
5.2.2. Link NLRI Usage	11
5.2.2.1. BGP-LS Link NLRI Address Family Link Descriptor TLV	12
5.2.2.2. BGP-LS-SPF Link NLRI Attribute SPF Status TLV	13
5.2.3. IPv4/IPv6 Prefix NLRI Usage	14
5.2.3.1. BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV	14
5.2.4. BGP-LS Attribute Sequence Number TLV	15
5.3. BGP-LS-SPF End of RIB (EoR) Marker	15
5.4. BGP Next-Hop Information	16

6. Decision Process with the SPF Algorithm	16
6.1. BGP SPF NLRI Selection	17
6.1.1. BGP Self-Originated NLRI	17
6.2. Dual-Stack Support	18
6.3. SPF Calculation Based on BGP-LS-SPF NLRI	18
6.4. IPv4/IPv6 Unicast Address Family Interaction	21
6.5. NLRI Advertisement	22
6.5.1. Link/Prefix Failure Convergence	22
6.5.2. Node Failure Convergence	22
7. Error Handling	23
7.1. Processing of BGP-LS-SPF TLVs	23
7.2. Processing of BGP-LS-SPF NLRIs	24
7.3. Processing of BGP-LS Attributes	24
7.4. BGP-LS-SPF Link State NLRI Database Synchronization	24
8. IANA Considerations	24
8.1. BGP-LS-SPF Allocation in the SAFI Values Registry	24
8.2. BGP-LS-SPF Assignments in the BGP-LS NLRI and Attribute TLVs Registry	25
8.3. BGP-LS-SPF Node NLRI Attribute SPF Status TLV Status Registry	25
8.4. BGP-LS-SPF Link NLRI Attribute SPF Status TLV Status Registry	26
8.5. BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV Status Registry	26
8.6. Assignment in the BGP Error (Notification) Codes Registry	26
9. Security Considerations	27
10. Management Considerations	27
10.1. Configuration	27
10.2. Link Metric Configuration	28
10.3. Unnumbered Link Configuration	28
10.4. Adjacency End-of-RIB (EOR) Marker Requirement	28
10.5. backoff-config	28
10.6. BGP-LS-SPF NLRI Readvertisement Delay	28
10.7. Operational Data	29

10.8. BGP-LS-SPF Address Family Session Isolation	29
11. References	29
11.1. Normative References	29
11.2. Informative References	31
Acknowledgements	31
Contributors	32
Authors' Addresses	33

1. Introduction

Many Massively Scaled Data Centers (MSDCs) have converged on simplified Layer 3 (L3) routing. Furthermore, requirements for operational simplicity has led many of these MSDCs to converge on BGP [[RFC4271](#)] as their single routing protocol for both fabric routing and Data Center Interconnect (DCI) routing [[RFC7938](#)]. This document describes an alternative solution that leverages BGP-LS [[RFC9552](#)] and the Shortest Path First (SPF) algorithm used by Interior Gateway Protocols (IGPs).

This document leverages both the BGP protocol [[RFC4271](#)] and BGP-LS extensions [[RFC9552](#)]. The relationship as well as the scope of changes are described in Sections 2 and 3, respectively. The modifications to [[RFC4271](#)] for BGP SPF described herein only apply to IPv4 and IPv6 as underlay unicast Subsequent Address Family Identifiers (SAFIs). Operations for any other BGP SAFIs are outside the scope of this document.

This solution avails the benefits of both BGP and SPF-based IGPs. These include TCP-based flow-control, no periodic link-state refresh, and completely incremental Network Layer Reachability Information (NLRI) advertisement. These advantages can reduce the overhead in MSDCs where there is a high degree of Equal-Cost Multipath (ECMP) load balancing. Additionally, using an SPF-based computation can support fast convergence and the computation of Loop-Free Alternatives (LFAs). The SPF LFA extensions defined in [[RFC5286](#)] can be similarly applied to BGP SPF calculations. However, the details are a matter of implementation detail and out of scope for this document. Furthermore, a BGP-based solution lends itself to multiple peering models including those incorporating route reflectors [[RFC4456](#)] or controllers.

1.1. Terminology

This specification reuses terms defined in Section 1.1 of [[RFC4271](#)].

Additionally, this document introduces the following terms:

BGP SPF Routing Domain:

A set of BGP routers under a single administrative domain that exchange link-state information using the BGP-LS-SPF SAFI and compute routes using BGP SPF, as described herein.

BGP-LS-SPF NLRI: The BGP-LS Network Layer Reachability Information (NLRI) that is being advertised in the BGP-LS-SPF SAFI ([Section 5.1](#)) and is being used for BGP SPF route computation.

Dijkstra Algorithm: An algorithm for computing the shortest path from a given node in a graph to every other node in the graph.

Prefix NLRI: In the context of BGP SPF, this term refers to the IPv4 Topology Prefix NLRI and/or the IPv6 Topology Prefix NLRI.

1.2. BGP Shortest Path First (SPF) Motivation

Given that [\[RFC7938\]](#) already describes how BGP could be used as the sole routing protocol in an MSDC, one might question the motivation for defining an alternative BGP deployment model when a mature solution exists. For both alternatives, BGP offers the operational benefits of a single routing protocol as opposed to the combination of IGP for the underlay and BGP as the overlay. However, BGP SPF offers some unique advantages above and beyond standard BGP path-vector routing. With BGP SPF, the simple single-hop peering model recommended in [Section 5.2.1](#) of [\[RFC7938\]](#) is augmented with peering models requiring fewer BGP sessions.

A primary advantage is that all BGP speakers in the BGP SPF routing domain have a complete view of the topology. This allows support for ECMP, IP fast-reroute (e.g., Loop-Free Alternatives (LFAs) [\[RFC5286\]](#), Shared Risk Link Groups (SRLGs) [\[RFC4202\]](#), and other routing enhancements without advertisement of additional BGP paths [\[RFC7911\]](#) or other extensions.

With the BGP SPF decision process as defined in [Section 6](#), NLRI changes can be disseminated throughout the BGP routing domain much more rapidly. The added advantage of BGP using TCP for reliable transport leverages TCP's inherent flow-control and guaranteed in-order delivery.

Another primary advantage is a potential reduction in NLRI advertisement. With standard BGP path-vector routing, a single link failure may impact 100s or 1000s of prefixes and result in the withdrawal or readvertisement of the attendant NLRI. With BGP SPF, only the BGP speakers originating the Link NLRI need to withdraw the corresponding BGP-LS-SPF Link NLRI. Additionally, the changed NLRI is advertised immediately as opposed to normal BGP where it is only advertised after the best route selection. These advantages provide NLRI dissemination throughout the BGP SPF routing domain with efficiencies similar to link-state protocols.

With controller and route-reflector peering models, BGP SPF advertisement and distributed computation require a minimal number of sessions and copies of the NLRI as only the latest version of the NLRI from the originator is required (see [Section 4](#)). Given that verification of whether or not to advertise a link (with a BGP-LS-SPF Link NLRI) is done outside of BGP, each BGP speaker only needs as many sessions and copies of the NLRI as required for redundancy. Additionally, a controller could inject topology (i.e., BGP-LS-SPF NLRI) that is learned outside the BGP SPF routing domain.

Given that BGP-LS NLRI is already defined [RFC9552], this functionality can be reused for BGP-LS-SPF NLRI.

Another advantage of BGP SPF is that both IPv6 and IPv4 can be supported using the BGP-LS-SPF SAFI with the same BGP-LS-SPF Link NRIs. In many MSDC fabrics, the IPv4 and IPv6 topologies are congruent (refer to [Section 5.2.2](#)). However, beyond the scope of this document, BGP-LS-SPF NLRI multi-topology extensions could be defined to support separate IPv4, IPv6, unicast, and multicast topologies while sharing the same NLRI.

Finally, the BGP SPF topology can be used as an underlay for other BGP SAFIs (using the existing model) and realize all the above advantages.

1.3. Document Overview

This document begins with [Section 2](#) defining the precise relationship that BGP SPF has with the base BGP protocol [RFC4271] and [Section 3](#) defining the BGP - Link State (BGP-LS) extensions [RFC9552]. The BGP peering models as well as their respective trade-offs are then discussed in [Section 4](#). The remaining sections, which make up the bulk of the document, define the protocol enhancements necessary to support BGP SPF including BGP-LS extensions ([Section 5](#)), replacement of the base BGP decision process with the SPF computation ([Section 6](#)), and BGP SPF error handling ([Section 7](#)).

1.4. 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.

2. Base BGP Protocol Relationship

With the exception of the decision process, BGP SPF extensions leverage the BGP protocol [RFC4271] without change. This includes the BGP protocol Finite State Machine, BGP messages and their encodings, the processing of BGP messages, BGP attributes and path attributes, BGP NLRI encodings, and any error handling defined in [RFC4271], [RFC4760], and [RFC7606].

Due to changes in the decision process, there are mechanisms and encodings that are no longer applicable. Unless explicitly specified in the context of BGP SPF, all optional path attributes **SHOULD NOT** be advertised. If received, all path attributes **MUST** be accepted, validated, and propagated consistently with the BGP protocol [RFC4271], even if not needed by BGP SPF.

[Section 9.1](#) of [RFC4271] defines the decision process that is used to select routes for subsequent advertisement by applying the policies in the local Policy Information Base (PIB) to the routes stored in its Adj-RIBs-In. The output of the Decision Process is the set of routes that are announced by a BGP speaker to its peers. These selected routes are stored by a BGP speaker in the speaker's Adj-RIBs-Out, according to policy.

The BGP SPF extension fundamentally changes the decision process, as described herein. Specifically:

1. BGP advertisements are readvertised to neighbors immediately without waiting or dependence on the route computation, as specified in phase 3 of the base BGP decision process. Multiple peering models are supported, as specified in [Section 4](#).
2. Determining the degree of preference for BGP routes for the SPF calculation as described in phase 1 of the base BGP decision process is replaced with the mechanisms in [Section 6.1](#).
3. Phase 2 of the base BGP protocol decision process is replaced with the SPF algorithm, also known as the Dijkstra algorithm.

3. BGP - Link State (BGP-LS) Relationship

[\[RFC9552\]](#) describes a mechanism by which link-state and Traffic Engineering (TE) information can be collected from networks and shared with external entities using BGP. This is achieved by defining NLRI that are advertised using the BGP-LS AFI. The BGP-LS extensions defined in [\[RFC9552\]](#) make use of the decision process defined in [\[RFC4271\]](#). Rather than reusing the BGP-LS SAFI, the BGP-LS-SPF SAFI ([Section 5.1](#)) is introduced to ensure backward compatibility for BGP-LS SAFI usage.

The "BGP-LS NLRI and Attribute TLVs" registry [\[RFC9552\]](#) is shared between the BGP-LS SAFI and the BGP-LS-SPF SAFI. However, the TLVs defined in this document may not be applicable to the BGP-LS SAFI. As specified in [Section 5.1](#) of [\[RFC9552\]](#), the presence of unknown or unexpected TLVs is required so that the NLRI or BGP-LS Attribute will not be considered malformed ([Section 5.2](#) of [\[RFC9552\]](#)). The list of BGP-LS TLVs applicable to the BGP-LS-SPF SAFI are described in [Section 5.2](#). By default, the usage of other BGP-LS TLVs or extensions are ignored for the BGP-LS-SPF SAFI. However, this doesn't preclude the usage specification of these TLVs for the BGP-LS-SPF SAFI in future documents.

4. BGP SPF Peering Models

Depending on the topology, scaling, capabilities of the BGP speakers, and redundancy requirements, various peering models are supported. The only requirement is that all BGP speakers in the BGP SPF routing domain adhere to this specification.

The choice of the deployment model is up to the operator and their requirements and policies. Deployment model choice is out of scope for this document and is discussed in [\[RFC9816\]](#). The subsections below describe several BGP SPF deployment models. However, this doesn't preclude other deployment models.

4.1. BGP Single-Hop Peering on Network Node Connections

The simplest peering model is the one where External BGP (EBGP) single-hop sessions are established over direct point-to-point links interconnecting the nodes in the BGP SPF routing domain. Once the single-hop BGP session has been established and the Multiprotocol Extensions

capabilities have been exchanged with the BGP-LS-SPF AFI/SAFI [RFC4760] for the corresponding session, then the link is considered up and available from a BGP SPF perspective, and the corresponding BGP-LS-SPF Link NLRI is advertised.

An End-of-RIB (EoR) marker (Section 5.3) for the BGP-LS-SPF SAFI **MAY** be required from a peer prior to advertising the BGP-LS-SPF Link NLRI for the corresponding link to that peer. When required, the default wait indefinitely for the EoR marker prior to advertising the BGP-LS-SPF Link NLRI. Refer to Section 10.4.

A failure to consistently configure the use of the EoR marker can result in transient micro-loops and dropped traffic due to incomplete forwarding state.

If the session goes down, the corresponding Link NLRIs are withdrawn. Topologically, this would be equivalent to the peering model in [RFC7938] where there is a BGP session on every link in the data center switch fabric. The content of the Link NLRI is described in Section 5.2.2.

4.2. BGP Peering Between Directly Connected Nodes

In this model, BGP speakers peer with all directly connected nodes but the sessions may be between loopback addresses (i.e., two-hop sessions), and the direct connection discovery and liveness detection for the interconnecting links are independent of the BGP protocol. The Bidirectional Forwarding Detection (BFD) protocol [RFC5880] is **RECOMMENDED** for liveness detection. Usage of other liveness connection mechanisms is outside the scope of this document. Consequently, there is a single BGP session even if there are multiple direct connections between BGP speakers. The 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 [RFC4760], the link is operational as determined using liveness detection mechanisms, and, optionally, the EoR marker has been received as described in Section 5.3. This is much like the previous peering model, except peering is between loopback addresses and the interconnecting links can be unnumbered. However, since there are BGP sessions between every directly connected node in the BGP SPF routing domain, there is a reduction in BGP sessions when there are parallel links between nodes. Hence, this peering model is **RECOMMENDED** over the single-hop peering model Section 4.1.

4.3. BGP Peering in Route-Reflector or Controller Topology

In this model, BGP speakers peer solely with one or more route reflectors [RFC4456] or controllers. As in the previous model, direct connection discovery and liveness detection for those links in the BGP SPF routing domain are done outside of the BGP protocol. BGP-LS-SPF Link NLRI is advertised as long as the corresponding link is considered up and available as per the chosen liveness detection mechanism (thus, the BFD protocol [RFC5880] is **RECOMMENDED**).

This peering model, known as "sparse peering", allows for fewer BGP sessions and, consequently, fewer instances of the same NLRI received from multiple peers. Ideally, the route reflectors or controller BGP sessions would be on directly connected links to avoid dependence

on another routing protocol for session connectivity. However, multi-hop peering is not precluded. The number of BGP sessions is dependent on the redundancy requirements and the stability of the BGP sessions.

The controller may use constraints to determine when to advertise BGP-LS-SPF NLRI for BGP-LS peers. For example, a controller may delay advertisement of a link between two peers until the EoR marker [Section 5.3](#) has been received from both BGP peers and the BGP-LS Link NLRI for the link(s) between the two nodes has been received from both BGP peers.

5. BGP Shortest Path Routing (SPF) Protocol Extensions

5.1. BGP-LS Shortest Path Routing (SPF) SAFI

This document introduces the BGP-LS-SPF SAFI with a value of 80. The SPF-based decision process ([Section 6](#)) applies only to the BGP-LS-SPF SAFI and **MUST NOT** be used with other combinations of the BGP-LS AFI (16388). In order for two BGP speakers to exchange BGP-LS-SPF NLRI, they **MUST** exchange Multiprotocol Extensions capabilities [[RFC4760](#)] to ensure that they are both capable of properly processing such an NLRI. This is done with AFI 16388 / SAFI 80. The BGP-LS-SPF SAFI is used to advertise IPv4 and IPv6 prefix information in a format facilitating an SPF-based decision process.

5.1.1. BGP-LS-SPF NLRI TLVs

All the TLVs defined for BGP-LS [[RFC9552](#)] are applicable and can be used with the BGP-LS-SPF SAFI to describe links, nodes, and prefixes comprising BGP SPF Link State Database (LSDB) information.

The NLRI and comprising TLVs **MUST** be encoded as specified in [Section 5.1](#) of [[RFC9552](#)]. TLVs specified as mandatory in [[RFC9552](#)] are considered mandatory for the BGP-LS-SPF SAFI as well. If a mandatory TLV is not present, the NLRI **MUST NOT** be used in the BGP SPF route calculation. All the other TLVs are considered as optional TLVs. Documents specifying usage of optional TLVs for BGP SPF **MUST** address backward compatibility.

5.1.2. BGP-LS Attribute

The BGP-LS attribute of the BGP-LS-SPF SAFI uses the exact same format as the BGP-LS AFI [[RFC9552](#)]. In other words, all the TLVs used in the BGP-LS attribute of the BGP-LS AFI are applicable and are used for the BGP-LS attribute of the BGP-LS-SPF SAFI. This attribute is an optional, non-transitive BGP attribute that is used to carry link, node, and prefix properties and attributes. The BGP-LS attribute is a set of TLVs.

All the TLVs defined for the BGP-LS Attribute [[RFC9552](#)] are applicable and can be used with the BGP-LS-SPF SAFI to carry link, node, and prefix properties and attributes.

The BGP-LS attribute may potentially be quite large depending on the amount of link-state information associated with a single BGP-LS-SPF NLRI. The BGP specification [[RFC4271](#)] mandates a maximum BGP message size of 4096 octets. It is **RECOMMENDED** that an implementation support [[RFC8654](#)] in order to accommodate a greater amount of information

within the BGP-LS Attribute. BGP speakers **MUST** ensure that they limit the TLVs included in the BGP-LS Attribute to ensure that a BGP update message for a single BGP-LS-SPF NLRI does not cross the maximum limit for a BGP message. The determination of the types of TLVs to be included by the BGP speaker originating the attribute is outside the scope of this document. If, due to the limits on the maximum size of an UPDATE message, a single route doesn't fit into the message, the BGP speaker **MUST NOT** advertise the route to its peer and **MAY** choose to log an error locally [RFC4271].

5.2. Extensions to BGP-LS

[RFC9552] describes a mechanism by which link-state and TE information can be collected from IGP and shared with external components using the BGP protocol. It describes both the definition of the BGP-LS NLRI that advertise links, nodes, and prefixes comprising IGP link-state information and the definition of a BGP path attribute (BGP-LS attribute) that carries link, node, and prefix properties and attributes, such as the link and prefix metric or auxiliary Router-IDs of nodes, etc. This document extends the usage of BGP-LS NLRI for the purpose of BGP SPF calculation via advertisement in the BGP-LS-SPF SAFI.

The protocol identifier specified in the Protocol-ID field [RFC9552] represents the origin of the advertised NLRI. For Node NLRI and Link NLRI, the specified Protocol-ID **MUST** be the direct protocol (4). Node or Link NLRI with a Protocol-ID other than the direct protocol is considered malformed. For Prefix NLRI, the specified Protocol-ID **MUST** be the origin of the prefix. The Local and Remote Node Descriptors for all NLRI **MUST** include the BGP Router-ID (TLV 516) [RFC9086] and the Autonomous System (TLV 512) number [RFC9552]. The BGP Confederation Member (TLV 517) [RFC9086] is not applicable.

5.2.1. Node NLRI Usage

The Node NLRI **MUST** be advertised unconditionally by all routers in the BGP SPF routing domain.

5.2.1.1. BGP-LS-SPF Node NLRI Attribute SPF Status TLV

A BGP-LS Attribute SPF Status TLV of the BGP-LS-SPF Node NLRI is defined to indicate the status of the node with respect to the BGP SPF calculation. This is used to rapidly take a node out of service (refer to Section 6.5.2) or to indicate that the node is not to be used for transit (i.e., non-local) traffic (refer to Section 6.3). If the SPF Status TLV is not included with the Node NLRI, the node is considered to be up and is available for transit traffic. A single TLV type is shared by the Node, Link, and Prefix NLRI. The TLV type is 1184.

```

      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 (1184)   |             Length (1 Octet)             |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| SPF Status      |
+---+---+---+---+---+

```

Value	Description
0	Reserved
1	Node unreachable with respect to BGP SPF
2	Node does not support transit with respect to BGP SPF
3-254	Unassigned
255	Reserved

Table 1: SPF Status Values

If a BGP speaker received the Node NLRI but the SPF Status TLV is not received, then any previously received SPF status information is considered as implicitly withdrawn, and the NLRI is propagated to other BGP speakers. A BGP speaker receiving a BGP Update containing an SPF Status TLV in the BGP-LS attribute [RFC9552] with an unknown value **SHOULD** be advertised to other BGP speakers and **MUST** ignore the Status TLV with an unknown value in the SPF computation. An implementation **MAY** log this condition for further analysis. If the SPF Status TLV contains a reserved value (0 or 255), the TLV is considered malformed and is handled as described in Section 7.1.

5.2.2. Link NLRI Usage

The criteria for advertisement of Link NLRIs are discussed in Section 4.

Link NLRI is advertised with unique Local and Remote Node Descriptors dependent on the IP addressing. For IPv4 links, the link's local IPv4 interface address (TLV 259) and remote IPv4 neighbor address (TLV 260) are used. For IPv6 links, the local IPv6 interface address (TLV 261) and remote IPv6 neighbor address (TLV 262) are used (see Section 5.2.2 of [RFC9552]). IPv6 links without global IPv6 addresses are considered unnumbered links and are handled as described below. For links supporting both IPv4 and IPv6 addresses, both sets of descriptors **MAY** be included in the same Link NLRI.

For unnumbered links, the Link Local/Remote Identifiers (TLV 258) are used. The Link Remote Identifier isn't normally exchanged in BGP, and discovering the Link Remote Identifier is beyond the scope of this document. If the Link Remote Identifier is unknown, a Link Remote Identifier of 0 **MUST** be advertised. When 0 is advertised and there are parallel unnumbered links between a pair of BGP speakers, there may be transient intervals where the BGP speakers don't agree on which of the parallel unnumbered links are operational. For this reason, it is **RECOMMENDED** that the Link Remote Identifiers be known (e.g., discovered using alternate mechanisms or configured) in the presence of parallel unnumbered links.

The link descriptors are described in Table 4 of [RFC9552]. Additionally, the Address Family (AF) Link Descriptor TLV is defined to determine whether an unnumbered link can be used in the IPv4 SPF, the IPv6, or both (refer to Section 5.2.2.1).

For a link to be used in SPF computation for a given address family, i.e., IPv4 or IPv6, both routers connecting the link **MUST** have matching addresses (i.e., router interface addresses must be on the same subnet for numbered interfaces, and the local/remote link identifiers ([Section 6.3](#)) must match for unnumbered interfaces).

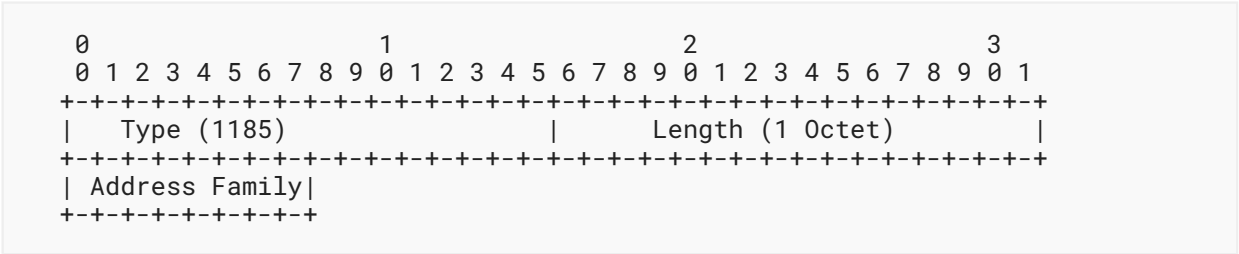
The IGP Metric (TLV 1095) **MUST** be advertised. If a BGP speaker receives a Link NLRI without an IGP Metric attribute TLV, then it **MUST** consider the received NLRI as malformed (refer to [Section 7](#)). The BGP SPF metric length is 4 octets. A metric is associated with the output side of each router interface. This metric is configurable by the system administrator. The lower the metric, the more likely the interface is to be used to forward data traffic. One possible default for the metric would be to give each interface a metric of 1 making it effectively a hop count.

The usage of other link attribute TLVs is beyond the scope of this document.

5.2.2.1. BGP-LS Link NLRI Address Family Link Descriptor TLV

For unnumbered links, the address family cannot be ascertained from the endpoint link descriptors. Hence, the Address Family Link Descriptor **SHOULD** be included with the Link Local/Remote Identifiers TLV for unnumbered links, so that the link can be used in the respective address family SPF. If the Address Family Link Descriptor is not present for an unnumbered link, the link will not be used in the SPF computation for either address family. If the Address Family Link Descriptor is present for a numbered link, the link descriptor will be ignored. If the Address Family Link Descriptor TLV contains an undefined value (3-254), the link descriptor will be ignored. If the Address Family Link Descriptor TLV contains a reserved value (0 or 255), the TLV is considered malformed and is handled as described in [Section 7.1](#).

Note that an unnumbered link can be used for both the IPv4 and IPv6 SPF computation by advertising separate Address Family Link Descriptor TLVs for IPv4 and IPv6.



Value	Description
0	Reserved
1	IPv4 Address Family
2	IPv6 Address Family
3-254	Undefined

Value	Description
255	Reserved

Table 2: Address Family Values

5.2.2.2. BGP-LS-SPF Link NLRI Attribute SPF Status TLV

The BGP-LS-SPF Attribute TLV of the BGP-LS-SPF Link NLRI is defined to indicate the status of the link with respect to the BGP SPF calculation. This is used to expedite convergence for link failures as discussed in [Section 6.5.1](#). If the SPF Status TLV is not included with the Link NLRI, the link is considered up and available. The SPF status is acted upon with the execution of the next SPF calculation ([Section 6.3](#)). A single TLV type is shared by the Node, Link, and Prefix NLRI. The TLV type is 1184.

```

      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 (1184)   |           Length (1 Octet)           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| SPF Status      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Value	Description
0	Reserved
1	Link unreachable with respect to BGP SPF
2-254	Unassigned
255	Reserved

Table 3: BGP Status Values

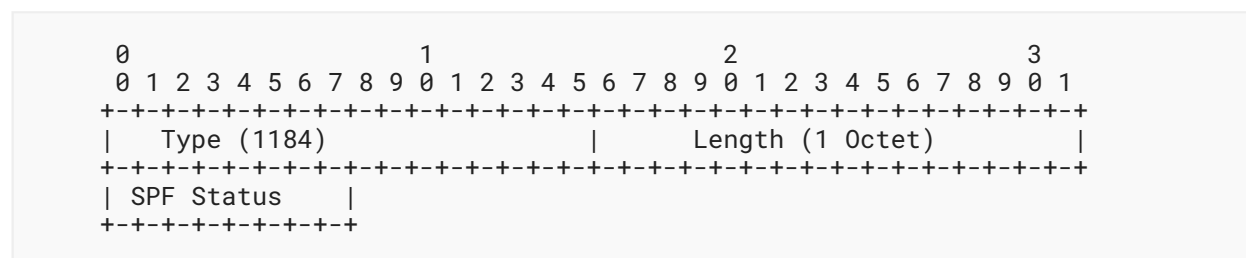
If a BGP speaker received the Link NLRI but the SPF Status TLV is not received, then any previously received SPF status information is considered as implicitly withdrawn, and the NLRI is propagated to other BGP speakers. A BGP speaker receiving a BGP Update containing an SPF Status TLV in the BGP-LS attribute [[RFC9552](#)] with an unknown value **SHOULD** be advertised to other BGP speakers and **MUST** ignore the SPF Status TLV with an unknown value in the SPF computation. An implementation **MAY** log this information for further analysis. If the SPF Status TLV contains a reserved value (0 or 255), the TLV is considered malformed and is handled as described in [Section 7.1](#).

5.2.3. IPv4/IPv6 Prefix NLRI Usage

An IPv4/IPv6 Prefix NLRI is advertised with a Local Node Descriptor and the prefix and length. The Prefix Descriptor field includes IP Reachability Information (TLV 265) as described in [RFC9552]. The Prefix Metric (TLV 1155) **MUST** be advertised to be considered for route calculation. The IGP Route Tag (TLV 1153) **MAY** be advertised. The usage of other BGP-LS attribute TLVs is beyond the scope of this document.

5.2.3.1. BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV

A BGP-LS Attribute SPF Status TLV of the BGP-LS-SPF Prefix NLRI is defined to indicate the status of the prefix with respect to the BGP SPF calculation. This is used to expedite convergence for prefix unreachability, as discussed in Section 6.5.1. If the SPF Status TLV is not included with the Prefix NLRI, the prefix is considered reachable. A single TLV type is shared by the Node, Link, and Prefix NLRI. The TLV type is 1184.



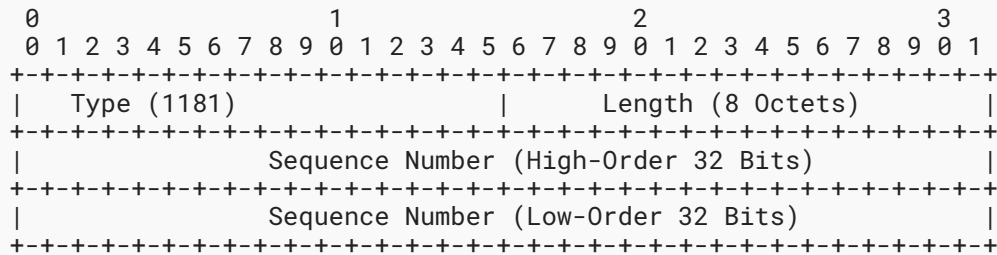
Value	Description
0	Reserved
1	Prefix unreachable with respect to BGP SPF
2-254	Unassigned
255	Reserved

Table 4: BGP Status Values

If a BGP speaker received the Prefix NLRI but the SPF Status TLV is not received, then any previously received SPF status information is considered as implicitly withdrawn, and the NLRI is propagated to other BGP speakers. A BGP speaker receiving a BGP Update containing an SPF Status TLV in the BGP-LS attribute [RFC9552] with an unknown value **SHOULD** be advertised to other BGP speakers and **MUST** ignore the Status TLV with an unknown value in the SPF computation. An implementation **MAY** log this information for further analysis. If the SPF Status TLV contains a reserved value (0 or 255), the TLV is considered malformed and is handled as described in Section 7.1.

5.2.4. BGP-LS Attribute Sequence Number TLV

A BGP-LS Attribute Sequence Number TLV of the BGP-LS-SPF NLRI types is defined to assure the most recent version of a given NLRI is used in the SPF computation. The Sequence Number TLV is mandatory for BGP-LS-SPF NLRI. The TLV type 1181 has been assigned by IANA. The BGP-LS Attribute Sequence Number TLV contains an 8-octet sequence number. The usage of the Sequence Number TLV is described in [Section 6.1](#).



Sequence Number: The 64-bit strictly increasing sequence number **MUST** be incremented for every self-originated version of a BGP-LS-SPF NLRI. BGP speakers implementing this specification **MUST** use available mechanisms to preserve the sequence number's strictly increasing property for the deployed life of the BGP speaker (including cold restarts). One mechanism for accomplishing this would be to use the high-order 32 bits of the sequence number as a wrap/boot count that is incremented any time the BGP router loses its sequence number state or the low-order 32 bits wrap.

When incrementing the sequence number for each self-originated NLRI, the sequence number should be treated as an unsigned 64-bit value. If the lower-order 32-bit value wraps, the higher-order 32-bit value should be incremented and saved in non-volatile storage. If a BGP speaker completely loses its sequence number state (e.g., the BGP speaker hardware is replaced or experiences a cold start), the BGP NLRI selection rules (see [Section 6.1](#)) ensure convergence, albeit not immediately.

If the Sequence Number TLV is not received, then the corresponding NLRI is considered as malformed and **MUST** be handled as 'treat-as-withdraw'. An implementation **SHOULD** log an error for further analysis.

5.3. BGP-LS-SPF End of RIB (EoR) Marker

The usage of the EoR marker [[RFC4724](#)] with the BGP-LS-SPF SAFI is somewhat different than the other BGP SAFIs. Reception of the EoR marker **MAY** optionally be expected prior to advertising a Link NLRI for a given peer.

5.4. BGP Next-Hop Information

The rules for setting the BGP Next-Hop in the MP_REACH_NLRI attribute [RFC4760] for the BGP-LS-SPF SAFI follow the rules in Section 5.5 of [RFC9552]. All BGP peers that support SPF extensions will locally compute the Local-RIB Next-Hop as a result of the SPF process. Hence, the use of the MP_REACH_NLRI Next-Hop as a tiebreaker in the standard BGP path decision processing is not applicable.

6. Decision Process with the SPF Algorithm

The Decision Process described in [RFC4271] takes place in three distinct phases. The Phase 1 decision function of the Decision Process is responsible for calculating the degree of preference for each route received from a BGP speaker's peer. The Phase 2 decision function is invoked on completion of the Phase 1 decision function and is responsible for choosing the best route out of all those available for each distinct destination and for installing each chosen route into the Local-RIB. The combination of the Phase 1 and 2 decision functions is characterized as a Path Vector algorithm.

The SPF-based Decision Process replaces the BGP Decision Process described in [RFC4271]. Since BGP-LS-SPF NLRI always contains the Local Node Descriptor as described in Section 5.2, each NLRI is uniquely originated by a single BGP speaker in the BGP SPF routing domain (the BGP node matching the NLRI's Node Descriptors). Instances of the same NLRI originated by multiple BGP speakers would be indicative of a configuration error or a masquerading attack (refer to Section 9). These selected Node NLRIs and their Link/Prefix NLRIs are used to build a directed graph during the SPF computation as described below. The best routes for BGP prefixes are installed in the RIB as a result of the SPF process.

When BGP-LS-SPF NLRI is received, all that is required is to determine whether it is the most recent by examining the Node-ID and sequence number as described in Section 6.1. If the received NLRI has changed, it is advertised to other BGP-LS-SPF peers. If the attributes have changed (other than the sequence number), a BGP SPF calculation is triggered. However, a changed NLRI **MAY** be advertised immediately to other peers and prior to any SPF calculation. Note that the BGP MinASOriginationIntervalTimer [RFC4271] timer is not applicable to the BGP-LS-SPF SAFI. The MinRouteAdvertisementIntervalTimer is applicable with a suggested default of 5 seconds consistent with Internal BGP (IBGP) (refer to Section 10 of [RFC4271]). The scheduling of the SPF calculation, as described in Section 6.3, is an implementation and/or configuration matter. Scheduling **MAY** be dampened consistent with the SPF Back-Off Delay algorithm specified in [RFC8405].

The Phase 3 decision function of the Decision Process [RFC4271] is also simplified because under normal SPF operation, a BGP speaker **MUST** advertise the changed NLRIs to all BGP peers with the BGP-LS-SPF AFI/SAFI and install the changed routes in the GLOBAL-RIB. The only exceptions are unchanged NLRIs or stale NLRIs, i.e., an NLRI received with a less recent (numerically smaller) sequence number.

6.1. BGP SPF NLRI Selection

For all BGP-LS-SPF NLRI, the selection rules for Phase 1 of the BGP decision process (see [Section 9.1.1](#) of [\[RFC4271\]](#)) no longer apply.

1. NLRI self-originated from directly connected BGP SPF peers are preferred. This condition can be determined by comparing the BGP Identifiers in the received Local Node Descriptor and the BGP OPEN message for an active BGP session. This rule assures that a stale NLRI is updated even if a BGP SPF router loses its sequence number state due to a cold start. Note that once the BGP session goes down, the NLRI received is no longer considered as being from a directly connected BGP SPF peer.
2. Consistent with base BGP [\[RFC4271\]](#), an NLRI received from a peer will always replace the same NLRI received from that peer. Coupled with rule #1, this will ensure that any stale NLRI in the BGP SPF routing domain will be updated.
3. The NLRI with the most recent Sequence Number TLV, i.e., the highest sequence number is selected.
4. The NLRI received from the BGP speaker with the numerically larger BGP Identifier is preferred.

When a BGP speaker completely loses its sequence number state, e.g., due to a cold start, or in the unlikely possibility that a 64-bit sequence number wraps, the BGP routing domain will still converge. This is due to the fact that BGP speakers adjacent to the router always accept self-originated NLRI from the associated speaker as more recent (rule #1). When a BGP speaker reestablishes a connection with its peers, any existing sessions are taken down and stale NLRI are replaced. The adjacent BGP speakers update their NLRI advertisements and advertise to their neighbors until the BGP routing domain has converged.

The modified SPF Decision Process performs an SPF calculation rooted at the local BGP speaker using the metrics from the Link Attribute IGP Metric (TLV 1095) and the Prefix Attribute Prefix Metric (TLV 1155) [\[RFC9552\]](#). These metrics are considered consistently across the BGP SPF domain. As a result, any other BGP attributes that would influence the BGP decision process defined in [\[RFC4271\]](#) including ORIGIN, MULTI_EXIT_DISC, and LOCAL_PREF attributes are ignored by the SPF algorithm. The Next Hop in the MP_REACH_NLRI attribute [\[RFC4760\]](#) is discussed in [Section 5.4](#). The AS_PATH and AS4_PATH attributes [\[RFC6793\]](#) are preserved and used for loop detection [\[RFC4271\]](#). They are ignored during the SPF computation for BGP-LS-SPF NLRI.

6.1.1. BGP Self-Originated NLRI

Nodes, Links, or Prefix NLRI with Node Descriptors matching the local BGP speaker are considered self-originated. When a self-originated NLRI is received and it doesn't match the local node's NLRI content (including the sequence number), special processing is required.

- If a self-originated NLRI is received and the sequence number is more recent (i.e., greater than the local node's sequence number for the NLRI), the NLRI sequence number is

advanced to one greater than the received sequence number, and the NLRI is readvertised to all peers.

- If a self-originated NLRI is received and the sequence number is the same as the local node's sequence number but the attributes differ, the NLRI sequence number is advanced to one greater than the received sequence number, and the NLRI is readvertised to all peers.

The above actions are performed immediately when the first instance of a newer self-originated NLRI is received. In this case, the newer instance is considered to be a stale instance that was advertised by the local node prior to a restart where the NLRI state was lost. However, if subsequent newer self-originated NLRI is received for the same Node, Link, or Prefix NLRI, the readvertisement or withdrawal is delayed by `BGP_LS_SPF_SELF_READVERTISEMENT_DELAY` (default 5) seconds since it is likely being advertised by a misconfigured or rogue BGP speaker (refer to [Section 9](#)).

6.2. Dual-Stack Support

The SPF-based decision process operates on Node, Link, and Prefix NLRIs that support both IPv4 and IPv6 addresses. Whether to run a single SPF computation or multiple SPF computations for separate AFs is an implementation and/or policy matter. Normally, IPv4 next-hops are calculated for IPv4 prefixes, and IPv6 next-hops are calculated for IPv6 prefixes.

6.3. SPF Calculation Based on BGP-LS-SPF NLRI

This section details the BGP-LS-SPF local Routing Information Base (RIB) calculation. The router uses BGP-LS-SPF Node, Link, and Prefix NLRIs to compute routes using the following algorithm. This calculation yields the set of routes associated with the BGP SPF routing domain. A router calculates the shortest-path tree using itself as the root. Optimizations to the BGP-LS-SPF algorithm are possible but **MUST** yield the same set of routes. The algorithm below supports ECMP routes. Weighted Unequal-Cost Multipath (UCMP) routes are out of scope.

The following abstract data structures are defined in order to specify the algorithm.

Local Route Information Base (Local-RIB): A routing table that contains reachability information (i.e., next hops) for all prefixes (both IPv4 and IPv6) as well as BGP-LS-SPF node reachability. Implementations may choose to implement this with separate RIBs for each address family and/or Prefix versus Node reachability.

Global Routing Information Base (GLOBAL-RIB): The RIB containing the current routes that are installed in the router's forwarding plane. This is commonly referred to in networking parlance as "the RIB".

Link State NLRI Database (LSNDB): A database of BGP-LS-SPF NLRIs that facilitate access to all Node, Link, and Prefix NLRIs.

Candidate List (CAN-LIST): A list of candidate Node NLRIs used during the BGP SPF calculation. The list is sorted by the cost to reach the Node NLRI, with the Node NLRI that has the lowest reachability cost at the head of the list. This facilitates the execution of the Dijkstra algorithm, where the shortest paths between the local node and other nodes in the graph are computed. The CAN-LIST is typically implemented as a heap but other data structures have been used.

The Dijkstra algorithm consists of the steps below:

1. The current Local-RIB is invalidated, and the CAN-LIST is initialized to be empty. The Local-RIB is rebuilt during the course of the SPF computation. The existing routing entries are preserved for comparison to determine changes that need to be made to the GLOBAL-RIB in Step 6. These routes are referred to as "stale routes".
2. The cost of the Local-RIB Node route entry for the computing router is set to 0. The computing router's Node NLRI is added to the CAN-LIST (which was previously initialized to be empty in Step 1). The next-hop list is set to the internal loopback next-hop.
3. The Node NLRI with the lowest cost is removed from the CAN-LIST for processing. If the BGP-LS Node attribute includes an SPF Status TLV (refer to [Section 5.2.1.1](#)) indicating the node is unreachable, the Node NLRI is ignored and the next lowest-cost Node NLRI is selected from the CAN-LIST. The Node corresponding to this NLRI is referred to as the "Current-Node". If the CAN-LIST list is empty, the SPF calculation has completed and the algorithm proceeds to Step 6.
4. All the Prefix NLRIs with the same Local Node Descriptors as the Current-Node are considered for installation. The next-hop(s) for these Prefix NLRIs are inherited from the Current-Node. If the Current-Node is for the local BGP Router, the next-hop for the prefix is a direct next-hop. The cost for each prefix is the metric advertised in the Prefix Attribute Prefix Metric (TLV 1155) added to the cost to reach the Current-Node. The following is done for each Prefix NLRI (referred to as the "Current-Prefix"):
 - If the BGP-LS Prefix attribute includes an SPF Status TLV indicating the prefix is unreachable, the Current-Prefix is considered unreachable, and the next Prefix NLRI is examined in Step 4.
 - If the Current-Prefix's corresponding prefix is in the Local-RIB and the Local-RIB metric is less than the Current-Prefix's metric, the Current-Prefix does not contribute to the route, and the next Prefix NLRI is examined in Step 4.
 - If the Current-Prefix's corresponding prefix is not in the Local-RIB, the prefix is installed with the Current-Node's next-hops installed as the Local-RIB route's next-hops and the metric being updated. If the IGP Route Tag (TLV 1153) is included in the Current-Prefix's NLRI Attribute, the tag(s) is installed in the current Local-RIB route's tag(s).
 - If the Current-Prefix's corresponding prefix is in the Local-RIB and the cost is less than the Local-RIB route's metric, the prefix is installed with the Current-Node's next-hops, which replace the Local-RIB route's next-hops and the metric being updated, and any route tags are removed. If the IGP Route Tag (TLV 1153) is included in the Current-Prefix's NLRI Attribute, the tag(s) is installed in the current Local-RIB route's tag(s).
 - If the Current-Prefix's corresponding prefix is in the Local-RIB and the cost is the same as the Local-RIB route's metric, the Current-Node's next-hops are merged with the Local-RIB

route's next-hops. The algorithm below supports ECMP routes. Some platforms or implementations may have limits on the number of ECMP routes that can be supported. The setting or identification of any limitations is outside the scope of this document. Weighted UCMP routes are out of scope as well.

5. All the Link NLRIs with the same Node Identifiers as the Current-Node are considered for installation. Each link is examined and referred to as the "Current-Link" in the following text. The cost of the Current-Link is the advertised IGP Metric (TLV 1095) from the Link NLRI BGP-LS attribute added to the cost to reach the Current-Node. If the Current-Node is for the local BGP Router, the next-hop for the link is a direct next-hop pointing to the corresponding local interface. For any other Current-Node, the next-hop(s) for the Current-Link is inherited from the Current-Node. The following is done for each link:
 - a. If the Current-Link's NLRI attribute includes an SPF Status TLV indicating the link is down, the BGP-LS-SPF Link NLRI is considered down, and the next link for the Current-Node is examined in Step 5.
 - b. If the Current-Node NLRI attributes include the SPF Status TLV (refer to [Section 5.2.1.1](#)) and the status indicates that the Node doesn't support transit, the next link for the Current-Node is processed in Step 5.
 - c. The Current-Link's Remote Node NLRI is accessed (i.e., the Node NLRI with the same Node Identifiers as the Current-Link's Remote Node Descriptors). If it exists, it is referred to as the "Remote-Node" and the algorithm proceeds as follows:
 - If the Remote-Node's NLRI attribute includes an SPF Status TLV indicating the node is unreachable, the next link for the Current-Node is examined in Step 5.
 - All the Link NLRIs corresponding to the Remote-Node are searched for a Link NLRI pointing to the Current-Node. Each Remote-Node's Link NLRI (referred to as the Remote-Link) is examined for Remote Node Descriptors matching the Current-Node and Link Descriptors matching the Current-Link.
 - For IPv4/IPv6 numbered Link Descriptors to match during the IPv4 SPF computation, the Current-Link's IPv4/IPv6 interface address link descriptor **MUST** match the Remote-Link IPv4/IPv6 neighbor address link descriptor, and the Current-Link's IPv4/IPv6 neighbor address **MUST** match the Remote-Link's IPv4/IPv6 interface address.
 - For unnumbered links to match during the IPv4 or IPv6 SPF computation, the Current-Link and Remote-Link's Address Family Link Descriptor TLV must match the address family of the IPv4 or IPv6 SPF computation, the Current-Link's Remote Identifier **MUST** match the Remote-Link's Local Identifier, and the Current-Link's Remote Identifier **MUST** match the Remote-Link's Local Identifier. Since the Link's Remote Identifier may not be known, a value of 0 is considered a wildcard and will match any Current or Remote Link's Local Identifier (see TLV 258 [[RFC9552](#)]). Address Family Link Descriptor TLVs for multiple address families may be advertised so that an unnumbered link can be used in the SPF computation for multiple address families.

If these conditions are satisfied for one of the Remote-Node's links, the bidirectional connectivity check succeeds and the Remote-Node may be processed further. The Remote-Node's Link NLRI providing bidirectional connectivity is referred to as the Remote-Link. If no Remote-Link is found, the next link for the Current-Node is examined in Step 5.

- If the Remote-Link NLRI attribute includes an SPF Status TLV indicating the link is down, the Remote-Link NLRI is considered down, and the next link for the Current-Node is examined in Step 5.
- If the Remote-Node is not on the CAN-LIST, it is inserted based on the cost. The Remote Node's cost is the cost of the Current-Node added to the Current-Link's IGP Metric (TLV 1095). The next-hop(s) for the Remote-Node is inherited from the Current-Link.
- If the Remote-Node NLRI is already on the CAN-LIST with a higher cost, it must be removed and reinserted with the Remote-Node cost based on the Current-Link (as calculated in the previous step). The next-hop(s) for the Remote-Node is inherited from the Current-Link.
- If the Remote-Node NLRI is already on the CAN-LIST with the same cost, it need not be reinserted on the CAN-LIST. However, the Current-Link's next-hop(s) must be merged into the current set of next-hops for the Remote-Node.
- If the Remote-Node NLRI is already on the CAN-LIST with a lower cost, it need not be reinserted on the CAN-LIST.

d. Return to Step 3 to process the next lowest-cost Node NLRI on the CAN-LIST.

6. The Local-RIB is examined and changes (adds, deletes, and modifications) are installed into the GLOBAL-RIB. For each route in the Local-RIB:

- If the route was added during the current BGP SPF computation, install the route into the GLOBAL-RIB.
- If the route was modified during the current BGP SPF computation (e.g., metric, tags, or next-hops), update the route in the GLOBAL-RIB.
- If the route was not installed during the current BGP SPF computation, remove the route from the GLOBAL-RIB.

6.4. IPv4/IPv6 Unicast Address Family Interaction

While the BGP-LS-SPF address family and the BGP unicast address families may install routes into the routing tables of the same device, they operate independently (i.e., "ships-in-the-night" mode). There is no implicit route redistribution between the BGP-LS-SPF address family and the BGP unicast address families.

It is **RECOMMENDED** that BGP-LS-SPF IPv4/IPv6 route computation and installation be given scheduling priority by default over other BGP address families as these address families are considered as underlay SAFIs.

6.5. NLRI Advertisement

6.5.1. Link/Prefix Failure Convergence

A local failure prevents a link from being used in the SPF calculation due to the IGP bidirectional connectivity requirement. Consequently, local link failures **SHOULD** always be communicated as quickly as possible and given priority over other categories of changes to ensure expeditious propagation and optimal convergence.

According to standard BGP procedures, the link would continue to be used until the last copy of the BGP-LS-SPF Link NLRI is withdrawn. In order to avoid this delay, the originator of the Link NLRI **SHOULD** advertise a more recent version with an increased Sequence Number TLV for the BGP-LS-SPF Link NLRI including the SPF Status TLV (refer to [Section 5.2.2.2](#)) indicating the link is down with respect to BGP SPF. The configurable LinkStatusDownAdvertise timer controls the interval that the BGP-LS-LINK NLRI is advertised with SPF Status indicating the link is down prior to withdrawal. If a BGP-LS-SPF Link NLRI has been advertised with the SPF Status TLV and the link becomes available in that period, the originator of the BGP-LS-SPF Link NLRI **MUST** advertise a more recent version of the BGP-LS-SPF Link NLRI without the SPF Status TLV in the BGP-LS Link Attributes. The suggested default value for the LinkStatusDownAdvertise timer is 2 seconds.

Similarly, when a prefix becomes unreachable, a more recent version of the BGP-LS-SPF Prefix NLRI **SHOULD** be advertised with the SPF Status TLV (refer to [Section 5.2.3.1](#)) to indicate that the prefix is unreachable in the BGP-LS Prefix Attributes, and the prefix will be considered unreachable with respect to BGP SPF. The configurable PrefixStatusDownAdvertise timer controls the interval that the BGP-LS-Prefix NLRI is advertised with SPF Status indicating the prefix is unreachable prior to withdrawal. If the BGP-LS-SPF Prefix has been advertised with the SPF Status TLV and the prefix becomes reachable in that period, the originator of the BGP-LS-SPF Prefix NLRI **MUST** advertise a more recent version of the BGP-LS-SPF Prefix NLRI without the SPF Status TLV in the BGP-LS Prefix Attributes. The suggested default value for the PrefixStatusDownAdvertise timer is 2 seconds.

6.5.2. Node Failure Convergence

By default, all the NLRIs advertised by a node are withdrawn when a session failure is detected [[RFC4271](#)]. If fast failure detection such as BFD [[RFC5880](#)] is utilized, and the node is on the fastest converging path, the most recent versions of BGP-LS-SPF NLRI will be withdrawn. This may result in older versions of NLRIs received from one or more peers on a different path(s) in the LSNDDB until they are withdrawn. These stale NLRIs will not delay convergence since the adjacent nodes detect the link failure and advertise a more recent NLRI indicating the link is down with respect to BGP SPF (refer to [Section 6.5.1](#)) and the bidirectional connectivity check fails during the BGP SPF calculation (refer to [Section 6.3](#)).

7. Error Handling

This section describes error-handling actions, as described in [\[RFC7606\]](#), that are specific to BGP-LS-SPF SAFI BGP Update message processing.

7.1. Processing of BGP-LS-SPF TLVs

When a BGP speaker receives a BGP Update containing a malformed Node NLRI SPF Status TLV in the BGP-LS Attribute [\[RFC9552\]](#), the corresponding Node NLRI is considered malformed and **MUST** be handled as 'treat-as-withdraw'. An implementation **SHOULD** log an error (subject to rate limiting) for further analysis.

When a BGP speaker receives a BGP Update containing a malformed Link NLRI SPF Status TLV in the BGP-LS Attribute [\[RFC9552\]](#), the corresponding Link NLRI is considered malformed and **MUST** be handled as 'treat-as-withdraw'. An implementation **SHOULD** log an error (subject to rate limiting) for further analysis.

When a BGP speaker receives a BGP Update containing a malformed Address Family Link Descriptor TLV in the BGP-LS Attribute [\[RFC9552\]](#), the corresponding Link NLRI is considered malformed and **MUST** be handled as 'treat-as-withdraw'. An implementation **SHOULD** log an error (subject to rate limiting) for further analysis.

When a BGP speaker receives a BGP Update containing a malformed Prefix NLRI SPF Status TLV in the BGP-LS Attribute [\[RFC9552\]](#), the corresponding Prefix NLRI is considered malformed and **MUST** be handled as 'treat-as-withdraw'. An implementation **SHOULD** log an error (subject to rate limiting) for further analysis.

When a BGP speaker receives a BGP Update containing a malformed BGP-LS Attribute TE and IGP Metric TLV, the corresponding NLRI is considered malformed and **MUST** be handled as 'treat-as-withdraw' [\[RFC7606\]](#). An implementation **SHOULD** log an error (subject to rate limiting) for further analysis.

The BGP-LS Attribute consists of Node attribute TLVs, Link attribute TLVs, and Prefix attribute TLVs. Node attribute TLVs and their error-handling rules are either defined in [\[RFC9552\]](#) or derived from [\[RFC5305\]](#) and [\[RFC6119\]](#). If a BGP speaker receives a BGP-LS Attribute that is considered malformed based on these error-handling rules, then it **MUST** consider the received NLRI as malformed, and the receiving BGP speaker **MUST** handle such a malformed NLRI as 'treat-as-withdraw' [\[RFC7606\]](#).

Node Descriptor TLVs and their error-handling rules are defined in [Section 5.2.1](#) of [\[RFC9552\]](#). Node Attribute TLVs and their error-handling rules are either defined in [\[RFC9552\]](#) or derived from [\[RFC5305\]](#) and [\[RFC6119\]](#).

Link Descriptor TLVs and their error-handling rules are defined in [Section 5.2.2](#) of [\[RFC9552\]](#). Link Attribute TLVs and their error-handling rules are either defined in [\[RFC9552\]](#) or derived from [\[RFC5305\]](#) and [\[RFC6119\]](#).

Prefix Descriptor TLVs and their error-handling rules are defined in [Section 5.2.3](#) of [\[RFC9552\]](#). Prefix Attribute TLVs and their error-handling rules are either defined in [\[RFC9552\]](#) or derived from [\[RFC5130\]](#) and [\[RFC2328\]](#).

If a BGP speaker receives NLRI with a Node Descriptor TLV, Link Descriptor TLV, or Prefix Descriptor TLV that is considered malformed based on error handling rules defined in the above references, then it **MUST** consider the received NLRI as malformed, and the receiving BGP speaker **MUST** handle such a malformed NLRI as 'treat-as-withdraw' [\[RFC7606\]](#).

When a BGP speaker receives a BGP Update that does not contain any BGP-LS Attributes, it is most likely an indication of 'Attribute Discard' fault handling, and the BGP speaker **SHOULD** preserve and propagate the BGP-LS-SPF NLRI as described in [Section 8.2.2](#) of [\[RFC9552\]](#). However, NLRIs without the BGP-LS attribute **MUST NOT** be used in the SPF calculation ([Section 6.3](#)). How this is accomplished is an implementation matter, but one way would be for these NLRIs not to be returned in LSNDDB lookups.

7.2. Processing of BGP-LS-SPF NLRIs

A BGP speaker supporting the BGP-LS-SPF SAFI **MUST** perform the syntactic validation checks of the BGP-LS-SPF NLRI listed in [Section 8.2.2](#) of [\[RFC9552\]](#) to determine if it is malformed.

7.3. Processing of BGP-LS Attributes

A BGP speaker supporting the BGP-LS-SPF SAFI **MUST** perform the syntactic validation checks of the BGP-LS Attribute listed in [Section 8.2.2](#) of [\[RFC9552\]](#) to determine if it is malformed.

An implementation **SHOULD** log an error for further analysis for problems detected during syntax validation.

7.4. BGP-LS-SPF Link State NLRI Database Synchronization

While uncommon, there may be situations where the LSNDDBs of two BGP speakers support the BGP-LS-SPF SAFI lose synchronization. In these situations, the BGP session **MUST** be reset unless other means of resynchronization are used (beyond the scope of this document). When the session is reset, the BGP speaker **MUST** send a NOTIFICATION message with the BGP error code "Loss of LSDB Synchronization" as described in [Section 3](#) of [\[RFC4271\]](#). The mechanisms to detect loss of synchronization are beyond the scope of this document.

8. IANA Considerations

8.1. BGP-LS-SPF Allocation in the SAFI Values Registry

IANA has assigned value 80 for BGP-LS-SPF from the First Come First Served range [\[RFC8126\]](#) and listed this document as a reference in the "SAFI Values" registry within the "Subsequent Address Family Identifiers (SAFI) Parameters" registry group.

8.2. BGP-LS-SPF Assignments in the BGP-LS NLRI and Attribute TLVs Registry

IANA has assigned six TLVs for BGP-LS-SPF NLRI in the "BGP-LS NLRI and Attribute TLVs" registry. Supported TLV types include Sequence Number, SPF Status, and Address Family Link Descriptor. Deprecated TLV types include SPF Capability, IPv4 Link Prefix Length, and IPv6 Link Prefix Length.

TLV Code Point	Description	Reference
1181	Sequence Number	Section 5.2.4 of RFC 9815
1184	SPF Status	Sections 5.2.1.1 , 5.2.2.2 , and 5.2.3.1 of RFC 9815
1185	Address Family Link Descriptor	Section 5.2.2.1 of RFC 9815

Table 5: NLRI Attribute TLVs

The early allocation assignments for the TLV types SPF Capability (1180), IPv4 Link Prefix Length (1182), and IPv6 Link Prefix Length (1183) are no longer required and have been deprecated.

8.3. BGP-LS-SPF Node NLRI Attribute SPF Status TLV Status Registry

IANA has created the "BGP-LS-SPF Node NLRI Attribute SPF Status TLV Status" registry for status values within the "BGP Shortest Path First (BGP SPF)" registry group. Initial values for this registry are provided below. Future assignments are to be made using the Expert Review registration policy [[RFC8126](#)] with guidance for designated experts as per [Section 7.2](#) of [[RFC9552](#)].

Values	Description
0	Reserved
1	Node unreachable with respect to BGP SPF
2	Node does not support transit traffic with respect to BGP SPF
3-254	Unassigned
255	Reserved

Table 6: BGP-LS-SPF Node NLRI Attribute SPF Status TLV Status Registry Assignments

8.4. BGP-LS-SPF Link NLRI Attribute SPF Status TLV Status Registry

IANA has created the "BGP-LS-SPF Link NLRI Attribute SPF Status TLV Status" registry for status values within the BGP Shortest Path First (BGP SPF) registry group. Initial values for this registry are provided below. Future assignments are to be made using the IETF Review registration policy [RFC8126].

Value	Description
0	Reserved
1	Link unreachable with respect to BGP SPF
2-254	Unassigned
255	Reserved

Table 7: BGP-LS-SPF Link NLRI Attribute SPF Status TLV Status Registry Assignments

8.5. BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV Status Registry

IANA has created the "BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV Status" registry for status values within the "BGP Shortest Path First (BGP SPF)" registry group. Initial values for this registry are provided below. Future assignments are to be made using the IETF Review registration policy [RFC8126].

Value	Description
0	Reserved
1	Prefix unreachable with respect to BGP SPF
2-254	Unassigned
255	Reserved

Table 8: BGP-LS-SPF Prefix NLRI Attribute SPF Status TLV Status Registry Assignments

8.6. Assignment in the BGP Error (Notification) Codes Registry

IANA has assigned value 9 for Loss of LSDB Synchronization in the "BGP Error (Notification) Codes" registry within the "Border Gateway Protocol (BGP) Parameters" registry group.

9. Security Considerations

This document defines a BGP SAFI, i.e., the BGP-LS-SPF SAFI. This document does not change the underlying security issues inherent in the BGP protocol [RFC4271]. The security considerations discussed in [RFC4271] apply to the BGP SPF functionality as well. The analysis of the security issues for BGP mentioned in [RFC4272] and [RFC6952] also applies to this document. The threats and security considerations are similar to the BGP IPv4 unicast SAFI and IPv6 unicast SAFI when utilized in similar deployments, e.g., [RFC7938]. The analysis of generic threats to routing protocols in [RFC4593] is also worth noting.

As the modifications for BGP SPF described in this document apply to IPv4 unicast and IPv6 unicast as underlay SAFIs in a single BGP SPF routing domain, the BGP security solutions described in [RFC6811] and [RFC8205] are out of scope as they are meant to apply for inter-domain BGP, where multiple BGP routing domains are typically involved. The BGP-LS-SPF SAFI NLRI described in this document are typically advertised between EBGP or IBGP speakers under a single administrative domain.

The BGP SPF processing and the BGP-LS-SPF SAFI inherit the encoding from BGP-LS [RFC9552], and consequently, inherit the security considerations for BGP-LS associated with encoding. Additionally, given that BGP SPF processing is used to install IPv4 and IPv6 unicast routes, the BGP SPF processing is vulnerable to attacks to the routing control plane that aren't applicable to BGP-LS. One notable Denial-of-Service attack would be to include malformed BGP attributes in a replicated BGP Update, causing the receiving peer to treat the advertised BGP-LS-SPF to a withdrawal [RFC7606].

In order to mitigate the risk of peering with BGP speakers masquerading as legitimate authorized BGP speakers, it is **RECOMMENDED** that the TCP Authentication Option (TCP-AO) [RFC5925] be used to authenticate BGP sessions. If an authorized BGP peer is compromised, that BGP peer could advertise a modified Node, Link, or Prefix NLRI that results in misrouting, repeating origination of NLRI, and/or excessive SPF calculations. When a BGP speaker detects that its self-originated NLRI is being originated by another BGP speaker, an appropriate error **SHOULD** be logged so that the operator can take corrective action. This exposure is similar to other BGP AFI/SAFIs.

10. Management Considerations

This section includes unique management considerations for the BGP-LS-SPF address family.

10.1. Configuration

All routers in the BGP SPF routing domain are under a single administrative domain allowing for consistent configuration.

10.2. Link Metric Configuration

For loopback prefixes, it is **RECOMMENDED** that the metric be 0. For non-loopback prefixes, the setting of the metric is a local matter and beyond the scope of this document.

Algorithms such as setting the metric inversely to the link speed as supported in some IGP implementations **MAY** be supported. However, the details of how the metric is computed are beyond the scope of this document.

Within a BGP SPF routing domain, the IGP metrics for all advertised links **SHOULD** be configured or defaulted consistently. For example, if a default metric is used for one router's links, then a similar metric should be used for all router's links. Similarly, if the link metric is derived from using the inverse of the link bandwidth on one router, then this **SHOULD** be done for all routers, and the same reference bandwidth **SHOULD** be used to derive the inversely proportional metric. Failure to do so will result in incorrect routing based on the link metric.

10.3. Unnumbered Link Configuration

When parallel unnumbered links between BGP SPF routers are included in the BGP SPF routing domain and the Remote Link Identifiers aren't readily discovered, it is **RECOMMENDED** that the Remote Link Identifiers be configured so that precise Link NLRI matching can be done.

10.4. Adjacency End-of-RIB (EOR) Marker Requirement

Depending on the peering model, topology, and convergence requirements, an EoR marker ([Section 5.3](#)) for the BGP-LS-SPF SAFI **MAY** be required from the peer prior to advertising a BGP-LS Link NLRI for the peer. If configuration is supported, this **MUST** be configurable at the BGP SPF instance level and **MUST** be configured consistently throughout the BGP SPF routing domain.

When this configuration is provided, the default **MUST** be to wait indefinitely prior to advertising a BGP-LS Link NLRI. Configuration of a timer specifying the maximum time to wait prior to advertisement **MAY** be provided.

10.5. backoff-config

In addition to the configuration of the BGP-LS-SPF address family, implementations **SHOULD** support "Shortest Path First (SPF) Back-Off Delay Algorithm for Link-State IGPs" [[RFC8405](#)]. If supported, configuration of the INITIAL_SPF_DELAY, SHORT_SPF_DELAY, LONG_SPF_DELAY, TIME_TO_LEARN, and HOLDDOWN_INTERVAL **MUST** be supported [[RFC8405](#)]. [Section 6](#) of [[RFC8405](#)] recommends consistent configuration of these values throughout the IGP routing domain, and this also applies to the BGP SPF routing domain.

10.6. BGP-LS-SPF NLRI Readvertisement Delay

The configuration parameter that specifies the delay for readvertising a more recent instance of a self-originated NLRI when received more than once in succession is BGP_LS_SPF_SELF_READVERTISEMENT_DELAY. The default is 5 seconds.

10.7. Operational Data

In order to troubleshoot SPF issues, implementations **SHOULD** support an SPF log including entries for previous SPF computations. Each SPF log entry would include the BGP-LS-SPF NLRI SPF triggering the SPF, SPF scheduled time, SPF start time, and SPF end time. Since the size of the log is finite, implementations **SHOULD** also maintain counters for the total number of SPF computations and the total number of SPF triggering events. Additionally, troubleshooting should be available for SPF scheduling and back-off [RFC8405], the current SPF back-off state, the remaining time-to-learn, the remaining hold-down interval, the last trigger event time, the last SPF time, and the next SPF time.

10.8. BGP-LS-SPF Address Family Session Isolation

In common deployment scenarios, the unicast routes installed during BGP-LS-SPF AFI/SAFI SPF computation serve as the underlay for other BGP AFI/SAFIs. To avoid errors encountered in other AFI/SAFIs from impacting the BGP-LS-SPF AFI/SAFI or vice versa, isolation mechanisms such as separate BGP instances or separate BGP sessions (e.g., using different addresses for peering) for BGP SPF Link-State information distribution **SHOULD** be used.

11. References

11.1. 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>>.
- [RFC4202] Kompella, K., Ed. and Y. Rekhter, Ed., "Routing Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 4202, DOI 10.17487/RFC4202, October 2005, <<https://www.rfc-editor.org/info/rfc4202>>.
- [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>>.
- [RFC4760] Bates, T., Chandra, R., Katz, D., and Y. Rekhter, "Multiprotocol Extensions for BGP-4", RFC 4760, DOI 10.17487/RFC4760, January 2007, <<https://www.rfc-editor.org/info/rfc4760>>.
- [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>>.

-
- [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>>.
- [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>>.
- [RFC5925] Touch, J., Mankin, A., and R. Bonica, "The TCP Authentication Option", RFC 5925, DOI 10.17487/RFC5925, June 2010, <<https://www.rfc-editor.org/info/rfc5925>>.
- [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>>.
- [RFC6793] Vohra, Q. and E. Chen, "BGP Support for Four-Octet Autonomous System (AS) Number Space", RFC 6793, DOI 10.17487/RFC6793, December 2012, <<https://www.rfc-editor.org/info/rfc6793>>.
- [RFC6811] Mohapatra, P., Scudder, J., Ward, D., Bush, R., and R. Austein, "BGP Prefix Origin Validation", RFC 6811, DOI 10.17487/RFC6811, January 2013, <<https://www.rfc-editor.org/info/rfc6811>>.
- [RFC7606] Chen, E., Ed., Scudder, J., Ed., Mohapatra, P., and K. Patel, "Revised Error Handling for BGP UPDATE Messages", RFC 7606, DOI 10.17487/RFC7606, August 2015, <<https://www.rfc-editor.org/info/rfc7606>>.
- [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>>.
- [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>>.
- [RFC8205] Lepinski, M., Ed. and K. Sriram, Ed., "BGPsec Protocol Specification", RFC 8205, DOI 10.17487/RFC8205, September 2017, <<https://www.rfc-editor.org/info/rfc8205>>.
- [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 IGPs", RFC 8405, DOI 10.17487/RFC8405, June 2018, <<https://www.rfc-editor.org/info/rfc8405>>.
- [RFC8654] Bush, R., Patel, K., and D. Ward, "Extended Message Support for BGP", RFC 8654, DOI 10.17487/RFC8654, October 2019, <<https://www.rfc-editor.org/info/rfc8654>>.
- [RFC9086] Previdi, S., Talaulikar, K., Ed., Filsfils, C., Patel, K., Ray, S., and J. Dong, "Border Gateway Protocol - Link State (BGP-LS) Extensions for Segment Routing BGP Egress Peer Engineering", RFC 9086, DOI 10.17487/RFC9086, August 2021, <<https://www.rfc-editor.org/info/rfc9086>>.
-

- [RFC9552] Talaulikar, K., Ed., "Distribution of Link-State and Traffic Engineering Information Using BGP", RFC 9552, DOI 10.17487/RFC9552, December 2023, <<https://www.rfc-editor.org/info/rfc9552>>.

11.2. Informative References

- [RFC4272] Murphy, S., "BGP Security Vulnerabilities Analysis", RFC 4272, DOI 10.17487/RFC4272, January 2006, <<https://www.rfc-editor.org/info/rfc4272>>.
- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", RFC 4456, DOI 10.17487/RFC4456, April 2006, <<https://www.rfc-editor.org/info/rfc4456>>.
- [RFC4593] Barbir, A., Murphy, S., and Y. Yang, "Generic Threats to Routing Protocols", RFC 4593, DOI 10.17487/RFC4593, October 2006, <<https://www.rfc-editor.org/info/rfc4593>>.
- [RFC4724] Sangli, S., Chen, E., Fernando, R., Scudder, J., and Y. Rekhter, "Graceful Restart Mechanism for BGP", RFC 4724, DOI 10.17487/RFC4724, January 2007, <<https://www.rfc-editor.org/info/rfc4724>>.
- [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>>.
- [RFC6952] Jethanandani, M., Patel, K., and L. Zheng, "Analysis of BGP, LDP, PCEP, and MSDP Issues According to the Keying and Authentication for Routing Protocols (KARP) Design Guide", RFC 6952, DOI 10.17487/RFC6952, May 2013, <<https://www.rfc-editor.org/info/rfc6952>>.
- [RFC7911] Walton, D., Retana, A., Chen, E., and J. Scudder, "Advertisement of Multiple Paths in BGP", RFC 7911, DOI 10.17487/RFC7911, July 2016, <<https://www.rfc-editor.org/info/rfc7911>>.
- [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>>.
- [RFC9816] Patel, K., Lindem, A., Zandi, S., Dawra, G., and J. Dong, "Usage and Applicability of BGP Link-State Shortest Path Routing (BGP-SPF) in Data Centers", RFC 9816, DOI 10.17487/RFC9816, July 2025, <<https://www.rfc-editor.org/info/rfc9816>>.

Acknowledgements

The authors would like to thank Sue Hares, Jorge Rabadan, Boris Hassanov, Dan Frost, Matt Anderson, Fred Baker, Lukas Krattiger, Yingzhen Qu, and Haibo Wang for their reviews and comments. Thanks to Pushpasis Sarkar for discussions on preventing a BGP SPF Router from being used for non-local traffic (i.e., transit traffic).

The authors extend a special thanks to Eric Rosen for fruitful discussions on BGP-LS-SPF convergence as compared to IGPs.

The authors would also like to thank the following people:

- Alvaro Retana for multiple AD reviews and discussions.
- Ketan Talaulikar for an extensive shepherd review.
- Adrian Farrel, Li Zhang, and Jie Dong for WG Last Call review comments.
- Jim Guichard for his AD review and discussion.
- David Dong for his IANA review.
- Joel Halpern for his GENART review.
- Erik Kline, Eric Vyncke, Mahesh Jethanandani, and Roman Danyliw for IESG review comments.
- John Scudder for his detailed IESG review and specifically for helping align the document with BGP documents.

Contributors

The following people contributed substantially to the content of this document and should be considered coauthors:

Derek Yeung

Arrcus, Inc.

Email: derek@arrcus.com

Gunter Van de Velde

Nokia

Email: gunter.van_de_velde@nokia.com

Abhay Roy

Arrcus, Inc.

Email: abhay@arrcus.com

Venu Venugopal

Cisco Systems

Email: venuv@cisco.com

Chaitanya Yadlapalli

AT&T

Email: cy098d@att.com

Authors' Addresses

Keyur Patel

Arrcus, Inc.

Email: keyur@arrcus.com

Acee Lindem

LabN Consulting, LLC

301 Midenhall Way

Cary, NC 27513

United States of America

Email: acee.ietf@gmail.com

Shawn Zandi

LinkedIn

222 2nd Street

San Francisco, CA 94105

United States of America

Email: szandi@linkedin.com

Wim Henderickx

Nokia

copernicuslaan 50

2018 Antwerp

Belgium

Email: wim.henderickx@nokia.com