

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: December 24, 2017

S. Previdi, Ed.
C. Filsfils
Cisco Systems, Inc.
P. Mattes
Microsoft
E. Rosen
Juniper Networks
S. Lin
Google
June 22, 2017

Advertising Segment Routing Policies in BGP
draft-previdi-idr-segment-routing-te-policy-07

Abstract

This document defines a new BGP SAFI with a new NLRI in order to advertise a candidate path of a Segment Routing Policy (SR Policy). An SR Policy is a set of candidate paths consisting of one or more segment lists. The headend of an SR Policy may learn multiple candidate paths for an SR Policy. Candidate paths may be learned via a number of different mechanisms, e.g., CLI, NetConf, PCEP, or BGP. This document specifies the way in which BGP may be used to distribute candidate paths. New sub-TLVs for the Tunnel Encapsulation Attribute are defined.

Status of This Memo

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

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

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

This Internet-Draft will expire on December 24, 2017.

Copyright Notice

Copyright (c) 2017 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1.	Introduction	3
1.1.	Requirements Language	5
2.	SR TE Policy Encoding	5
2.1.	SR TE Policy SAFI and NLRI	5
2.2.	SR TE Policy and Tunnel Encapsulation Attribute	7
2.3.	Remote Endpoint and Color	8
2.4.	SR TE Policy Sub-TLVs	8
2.4.1.	Preference sub-TLV	8
2.4.2.	SR TE Binding SID Sub-TLV	9
2.4.3.	Segment List Sub-TLV	10
3.	Extended Color Community	21
4.	SR Policy Operations	21
4.1.	Configuration and Advertisement of SR TE Policies	22
4.2.	Reception of an SR Policy NLRI	22
4.2.1.	Acceptance of an SR Policy NLRI	22
4.2.2.	Usable SR Policy NLRI	23
4.2.3.	Passing a usable SR Policy NLRI to the SRTE Process	24
4.2.4.	Propagation of an SR Policy	24
4.3.	Flowspec and SR Policies	24
5.	Contributors	24
6.	Acknowledgments	25
7.	Implementation Status	25
8.	IANA Considerations	26
8.1.	Existing Registry: Subsequent Address Family Identifiers (SAFI) Parameters	26
8.2.	Existing Registry: BGP Tunnel Encapsulation Attribute Tunnel Types	26
8.3.	Existing Registry: BGP Tunnel Encapsulation Attribute sub-TLVs	27
8.4.	New Registry: SR Policy List Sub-TLVs	27
9.	Security Considerations	27

10.	References	27
10.1.	Normative References	27
10.2.	Informational References	28
	Authors' Addresses	29

[1.](#) Introduction

Segment Routing (SR) allows a headend node to steer a packet flow along any path. Intermediate per-flow states are eliminated thanks to source routing [[I-D.ietf-spring-segment-routing](#)].

The headend node is said to steer a flow into an Segment Routing Policy (SR Policy).

The header of a packet steered in an SR Policy is augmented with the ordered list of segments associated with that SR Policy.

[[I-D.filsfils-spring-segment-routing-policy](#)] details the concepts of SR Policy and steering into an SR Policy. These apply equally to the MPLS and SRv6 instantiations of segment routing.

As highlighted in section 2 of [[I-D.filsfils-spring-segment-routing-policy](#)]:

- o an SR policy may have multiple candidate paths learned via various mechanisms (CLI, NetConf, PCEP or BGP);
- o the SRTE process selects the best candidate path for a Policy;
- o the SRTE process binds a BSID to the selected path of the Policy;
- o the SRTE process installs the selected path and its BSID in the forwarding plane.

This document specifies the way to use BGP to distribute one or more of the candidate paths of an SR policy to the headend of that policy. The SRTE process ([[I-D.filsfils-spring-segment-routing-policy](#)]) of the headend receives candidate paths from BGP, and possibly other sources as well, and the SRTE process then determines the selected path of the policy.

This document specifies a way of representing SR policies and their candidate paths in BGP UPDATE messages. BGP can then be used to propagate the SR policies and candidate paths. The usual BGP rules for BGP propagation and "bestpath selection" are used. At the headend of a specific policy, this will result in one or more candidate paths being installed into the "BGP table". These paths are then passed to the SRTE process. The SRTE process may compare

them to candidate paths learned via other mechanisms, and will choose one or more paths to be installed in the data plane. BGP itself does not install SRTE candidate paths into the data plane.

This document defines a new BGP address family (SAFI). In UPDATE messages of that address family, the NLRI identifies an SR policy, and the attributes specify candidate paths of that policy.

While for simplicity we may write that BGP advertises an SR Policy, it has to be understood that BGP advertises a candidate path of an SR policy and that this SR Policy might have several other candidate paths provided via BGP (via an NLRI with a different distinguisher as defined in this document), PCEP, NETCONF or local policy configuration.

Typically, a controller defines the set of policies and advertise them to policy head-end routers (typically ingress routers). The policy advertisement uses BGP extensions defined in this document. The policy advertisement is, in most but not all of the cases, tailored for a specific policy head-end. In this case the advertisement may sent on a BGP session to that head-end and not propagated any further.

Alternatively, a router (i.e.: an BGP egress router) advertises SR Policies representing paths to itself. In this case, it is possible to send the policy to each head-end over a BGP session to that head-end, without requiring any further propagation of the policy.

An SR Policy intended only for the receiver will, in most cases, not traverse any Route Reflector (RR, [[RFC4456](#)]).

In some situations, it is undesirable for a controller or BGP egress router to have a BGP session to each policy head-end. In these situations, BGP Route Reflectors may be used to propagate the advertisements, or it may be necessary for the advertisement to propagate through a sequence of one or more ASes. To make this possible, an attribute needs to be attached to the advertisement that enables a BGP speaker to determine whether it is intended to be a head-end for the advertised policy. This is done by attaching one or more Route Target Extended Communities to the advertisement ([[RFC4360](#)]).

The BGP extensions for the advertisement of SR Policies include following components:

- o A new Subsequent Address Family Identifier (SAFI) whose NLRI identifies an SR Policy.

- o A set of new TLVs to be inserted into the Tunnel Encapsulation Attribute (as defined in [[I-D.ietf-idr-tunnel-encaps](#)]) specifying candidate paths of the SR policy, as well as other information about the SR policy.
- o One or more IPv4 address format route-target extended community ([[RFC4360](#)]) attached to the SR Policy advertisement and that indicates the intended head-end of such SR Policy advertisement.
- o The Color Extended Community (as defined in [[I-D.ietf-idr-tunnel-encaps](#)]) and used in order to steer traffic into an SR Policy, as described in [[I-D.filsfils-spring-segment-routing-policy](#)]. This document ([Section 3](#)) modifies the format of the Color Extended Community by using the two leftmost bits of the RESERVED field.

1.1. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [[RFC2119](#)].

2. SR TE Policy Encoding

2.1. SR TE Policy SAFI and NLRI

A new SAFI is defined: the SR Policy SAFI, (codepoint 73 assigned by IANA (see [Section 8](#)) from the "Subsequent Address Family Identifiers (SAFI) Parameters" registry).

The SR Policy SAFI uses a new NLRI defined as follows:

```
+-----+
| NLRI Length      | 1 octet
+-----+
| Distinguisher    | 4 octets
+-----+
| Policy Color     | 4 octets
+-----+
| Endpoint         | 4 or 16 octets
+-----+
```

where:

- o NLRI Length: 1 octet of length expressed in bits as defined in [[RFC4760](#)].

- o Distinguisher: 4-octet value uniquely identifying the policy in the context of <color, endpoint> tuple. The distinguisher has no semantic value and is solely used by the SR Policy originator to make unique (from an NLRI perspective) multiple occurrences of the same SR Policy.
- o Policy Color: 4-octet value identifying (with the endpoint) the policy. The color is used to match the color of the destination prefixes to steer traffic into the SR Policy [[I-D.filsfils-spring-segment-routing-policy](#)].
- o Endpoint: identifies the endpoint of a policy. The Endpoint may represent a single node or a set of nodes (e.g., an anycast address or a summary address). The Endpoint is an IPv4 (4-octet) address or an IPv6 (16-octet) address according to the AFI of the NLRI.

The color and endpoint are used to automate the steering of BGP Payload prefixes on SR policy ([\[I-D.filsfils-spring-segment-routing-policy\]](#)).

The NLRI containing the SR Policy is carried in a BGP UPDATE message [[RFC4271](#)] using BGP multiprotocol extensions [[RFC4760](#)] with an AFI of 1 or 2 (IPv4 or IPv6) and with a SAFI of 73 (assigned by IANA from the "Subsequent Address Family Identifiers (SAFI) Parameters" registry).

An update message that carries the MP_REACH_NLRI or MP_UNREACH_NLRI attribute with the SR Policy SAFI MUST also carry the BGP mandatory attributes. In addition, the BGP update message MAY also contain any of the BGP optional attributes.

The next-hop of the SR Policy SAFI NLRI is set based on the AFI. For example, if the AFI is set to IPv4 (1), then the next-hop is encoded as a 4-byte IPv4 address. If the AFI is set to IPv6 (2), then the next-hop is encoded as a 16-byte IPv6 address of the router.

It is important to note that any BGP speaker receiving a BGP message with an SR Policy NLRI, will process it only if the NLRI is among the best paths as per the BGP best path selection algorithm. In other words, this document does not modify the BGP propagation or bestpath selection rules.

It has to be noted that if several candidate paths of the same SR Policy (endpoint, color) are signaled via BGP to a head-end, it is recommended that each NLRI use a different distinguisher. If BGP has installed into the BGP table two advertisements whose respective

NLRIs have the same color and endpoint, but different distinguishers, both advertisements are passed to the SRTE process.

2.2. SR TE Policy and Tunnel Encapsulation Attribute

The content of the SR Policy is encoded in the Tunnel Encapsulation Attribute originally defined in [[I-D.ietf-idr-tunnel-encaps](#)] using a new Tunnel-Type TLV (codepoint is 15, assigned by IANA (see [Section 8](#)) from the "BGP Tunnel Encapsulation Attribute Tunnel Types" registry).

The SR Policy Encoding structure is as follows:

SR Policy SAFI NLRI: <Distinguisher, Policy-Color, Endpoint>

Attributes:

 Tunnel Encaps Attribute (23)

 Tunnel Type: SR Policy

 Binding SID

 Preference

 Segment List

 Weight

 Segment

 Segment

 ...

 ...

where:

- o SR Policy SAFI NLRI is defined in [Section 2.1](#).
- o Tunnel Encapsulation Attribute is defined in [[I-D.ietf-idr-tunnel-encaps](#)].
- o Tunnel-Type is set to 15 (assigned by IANA from the "BGP Tunnel Encapsulation Attribute Tunnel Types" registry).
- o Preference, Binding SID, Segment-List, Weight and Segment are defined in this document.
- o Additional sub-TLVs may be defined in the future.

A Tunnel Encapsulation Attribute MUST NOT contain more than one TLV of type "SR Policy".

Multiple occurrences of "Segment List" MAY be encoded within the same SR Policy.

Multiple occurrences of "Segment" MAY be encoded within the same Segment List.

2.3. Remote Endpoint and Color

The Remote Endpoint and Color sub-TLVs, as defined in [\[I-D.ietf-idr-tunnel-encaps\]](#), MAY also be present in the SR Policy encodings.

If present, the Remote Endpoint sub-TLV MUST match the Endpoint of the SR Policy SAFI NLRI.

If present, the Color sub-TLV MUST match the Policy Color of the SR Policy SAFI NLRI.

2.4. SR TE Policy Sub-TLVs

This section defines the SR Policy sub-TLVs.

Preference, Binding SID, Segment-List are assigned from the "BGP Tunnel Encapsulation Attribute sub-TLVs" registry.

Weight and Segment Sub-TLVs are assigned from a new registry defined in this document and called: "SR Policy List Sub-TLVs". See [Section 8](#) for the details of the registry.

2.4.1. Preference sub-TLV

The Preference sub-TLV does not have any effect on the BGP bestpath selection or propagation procedures. The contents of this sub-TLV are used by the SRTE process ([\[I-D.ietf-spring-segment-routing-policy\]](#)).

The Preference sub-TLV is optional, MUST NOT appear more than once in the SR Policy and has following format:

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   | Length |   Flags   | RESERVED |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     Preference (4 octets)                                     |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

- o Type: TBD3 (to be assigned by IANA from the "BGP Tunnel Encapsulation Attribute sub-TLVs" registry).
- o Length: 6.

- o Type: TBD4 (to be assigned by IANA from the "BGP Tunnel Encapsulation Attribute sub-TLVs" registry).
- o Length: specifies the length of the value field not including Type and Length fields. Can be 2 or 6 or 18.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Binding SID: if length is 2, then no Binding SID is present. If length is 6 then the Binding SID contains a 4-octet SID. If length is 18 then the Binding SID contains a 16-octet IPv6 SID.

2.4.3. Segment List Sub-TLV

The Segment List TLV encodes a single explicit path towards the endpoint. The Segment List sub-TLV includes the elements of the paths (i.e.: segments) as well as an optional Weight TLV.

The Segment List sub-TLV may exceed 255 bytes length due to large number of segments. Therefore a 2-octet length is required. According to [I-D.ietf-idr-tunnel-encaps], the first bit of the sub-TLV codepoint defines the size of the length field. Therefore, for the Segment List sub-TLV a code point of 128 (or higher) is used. See [Section 8](#) for details of codepoints allocation.

The Segment List sub-TLV is mandatory and MAY appear multiple times in the SR Policy.

The Segment-List Sub-TLV MUST contain at least one Segment Sub-TLV and MAY contain a Weight Sub-TLV.

The Segment List sub-TLV has the following format:

```

      0               1               2               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   |             Length             |  RESERVED  |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                               sub-TLVs                               //
```

where:

- o Type: TBD5 (to be assigned by IANA from the "BGP Tunnel Encapsulation Attribute sub-TLVs" registry).
- o Length: the total length (not including the Type and Length fields) of the sub-TLVs encoded within the Segment List sub-TLV.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o sub-TLVs:
 - * An optional single Weight sub-TLV.
 - * One or more Segment sub-TLVs.

[I-D.filsfils-spring-segment-routing-policy] defines several types of Segment Sub-TLVs:


```
Type 1: SID only, in the form of MPLS Label
Type 2: SID only, in the form of IPv6 address
Type 3: IPv4 Node Address with optional SID
Type 4: IPv6 Node Address with optional SID
Type 5: IPv4 Address + index with optional SID
Type 6: IPv4 Local and Remote addresses with optional SID
Type 7: IPv6 Address + index with optional SID
Type 8: IPv6 Local and Remote addresses with optional SID
```

2.4.3.2.1. Type 1: SID only, in the form of MPLS Label

The Type-1 Segment Sub-TLV encodes a single SID in the form of an MPLS label. The format is as follows:

									1									2									3												
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+												
Type									Length									Flags									RESERVED												
+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+												
Label									TC S									TTL																					
+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+									+--+--+--+--+--+--+--+--+--+												

where:

- o Type: 1 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 6.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Label: 20 bits of label value.
- o TC: 3 bits of traffic class.
- o S: 1 bit of bottom-of-stack.
- o TTL: 1 octet of TTL.

The following applies to the Type-1 Segment sub-TLV:

- o The S bit SHOULD be zero upon transmission, and MUST be ignored upon reception.

- o If the originator wants the receiver to choose the TC value, it sets the TC field to zero.
- o If the originator wants the receiver to choose the TTL value, it sets the TTL field to 255.
- o If the originator wants to recommend a value for these fields, it puts those values in the TC and/or TTL fields.
- o The receiver MAY override the originator's values for these fields. This would be determined by local policy at the receiver. One possible policy would be to override the fields only if the fields have the default values specified above.

2.4.3.2.2. Type 2: SID only, in the form of IPv6 address

The Type-2 Segment Sub-TLV encodes a single SID in the form of an IPv6 SID. The format is as follows:

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   | Length |   Flags   | RESERVED |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                               IPv6 SID (16 octets)                               //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

- o Type: 2 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 18.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o IPv6 SID: 16 octets of IPv6 address.

The IPv6 Segment Identifier (IPv6 SID) is defined in [\[I-D.ietf-6man-segment-routing-header\]](#).

- o The IPv4 Node Address MUST be present.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).
 - * IPV6 SID: a 16 octet IPv6 address.
- o If length is 6, then only the IPv4 Node Address is present.

- o If length is 10, then the IPv4 Node Address and the MPLS SID are present.
- o If length is 22, then the IPv4 Node Address and the IPv6 SID are present.

2.4.3.2.4. Type 4: IPv6 Node Address with optional SID

The Type-4 Segment Sub-TLV encodes an IPv6 node address and an optional SID in the form of either an MPLS label or an IPv6 address. The format is as follows:

```

      0               1               2               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   | Length |   Flags   | RESERVED |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//          IPv6 Node Address (16 octets)          //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//          SID (optional, 4 or 16 octets)          //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

- o Type: 4 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 18 or 22 or 34.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o IPv6 Node Address: a 16 octet IPv6 address representing a node.
- o SID: either 4 octet MPLS SID or a 16 octet IPv6 SID.

The following applies to the Type-4 Segment sub-TLV:

- o The IPv6 Node Address MUST be present.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).

- * IPv6 SID: a 16 octet IPv6 address.
- o If length is 18, then only the IPv6 Node Address is present.
- o If length is 22, then the IPv6 Node Address and the MPLS SID are present.
- o If length is 34, then the IPv6 Node Address and the IPv6 SID are present.

2.4.3.2.5. Type 5: IPv4 Address + Local Interface ID with optional SID

The Type-5 Segment Sub-TLV encodes an IPv4 node address, a local interface Identifier (Local Interface ID) and an optional SID in the form of either an MPLS label or an IPv6 address. The format is as follows:

0										1										2										3									
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1								
Type										Length										Flags										RESERVED									
										Local Interface ID (4 octets)																													
										IPv4 Node Address (4 octets)																													
//										SID (optional, 4 or 16 octets)																				//									

where:

- o Type: 5 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 10 or 14 or 26.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Local Interface ID: 4 octets as defined in [\[I-D.ietf-pce-segment-routing\]](#).
- o IPv4 Node Address: a 4 octet IPv4 address representing a node.

- o SID: either 4 octet MPLS SID or a 16 octet IPv6 SID.

The following applies to the Type-5 Segment sub-TLV:

- o The IPv4 Node Address MUST be present.
- o The Local Interface ID MUST be present.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).
 - * IPV6 SID: a 16 octet IPv6 SID.
- o If length is 10, then the IPv4 Node Address and Local Interface ID are present.
- o If length is 14, then the IPv4 Node Address, the Local Interface ID and the MPLS SID are present.
- o If length is 26, then the IPv4 Node Address, the Local Interface ID and the IPv6 SID are present.

[2.4.3.2.6](#). Type 6: IPv4 Local and Remote addresses with optional SID

The Type-6 Segment Sub-TLV encodes an adjacency local address, an adjacency remote address and an optional SID in the form of either an MPLS label or an IPv6 address. The format is as follows:

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|      Type      | Length      |   Flags   |  RESERVED  |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                                     Local IPv4 Address (4 octets) |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                                     Remote IPv4 Address (4 octets) |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
//                                     SID (4 or 16 octets)          //
```

where:

- o Type: 6 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).

- o Length is 10 or 14 or 26.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Local IPv4 Address: a 4 octet IPv4 address.
- o Remote IPv4 Address: a 4 octet IPv4 address.
- o SID: either 4 octet MPLS SID or a 16 octet IPv6 SID.

The following applies to the Type-6 Segment sub-TLV:

- o The Local IPv4 Address MUST be present and represents an adjacency local address.
- o The Remote IPv4 Address MUST be present and represents the remote end of the adjacency.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).
 - * IPv6 SID: a 16 octet IPv6 address.
- o If length is 10, then only the IPv4 Local and Remote addresses are present.
- o If length is 14, then the IPv4 Local address, IPv4 Remote address and the MPLS SID are present.
- o If length is 26, then the IPv4 Local address, IPv4 Remote address and the IPv6 SID are present.

[2.4.3.2.7](#). Type 7: IPv6 Address + Local Interface ID with optional SID

The Type-7 Segment Sub-TLV encodes an IPv6 node address, a local interface identifier (Local Interface ID) and an optional SID in the form of either an MPLS label or an IPv6 address. The format is as follows:


```

      0               1               2               3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   | Length |   Flags   | RESERVED |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     Local Interface ID (4 octets) |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                     IPv6 Node Address (16 octets) //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                     SID (optional, 4 or 16 octets) //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

- o Type: 7 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 22 or 26 or 38.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Local Interface ID: 4 octets of interface index.
- o IPv6 Node Address: a 16 octet IPv6 address representing a node.
- o SID: either 4 octet MPLS SID or a 16 octet IPv6 SID.

The following applies to the Type-7 Segment sub-TLV:

- o The IPv6 Node Address MUST be present.
- o The Local Interface ID MUST be present.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).
 - * IPV6 SID: a 16 octet IPv6 address.
- o If length is 22, then the IPv6 Node Address and Local Interface ID are present.

- o If length is 26, then the IPv6 Node Address, the Local Interface ID and the MPLS SID are present.
- o If length is 38, then the IPv6 Node Address, the Local Interface ID and the IPv6 SID are present.

2.4.3.2.8. Type 8: IPv6 Local and Remote addresses with optional SID

The Type-8 Segment Sub-TLV encodes an adjacency local address, an adjacency remote address and an optional SID in the form of either an MPLS label or an IPv6 address. The format is as follows:

```

      0               1               2               3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Type   | Length |   Flags   | RESERVED |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//          Local IPv6 Address (16 octets)          //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//          Remote IPv6 Address (16 octets)          //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//          SID (4 or 16 octets)                     //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

where:

- o Type: 8 (to be assigned by IANA from the registry "SR Policy List Sub-TLVs" defined in this document).
- o Length is 34 or 38 or 50.
- o Flags: 1 octet of flags. None are defined at this stage. Flags SHOULD be set to zero on transmission and MUST be ignored on receipt.
- o RESERVED: 1 octet of reserved bits. SHOULD be unset on transmission and MUST be ignored on receipt.
- o Local IPv6 Address: a 16 octet IPv6 address.
- o Remote IPv6 Address: a 16 octet IPv6 address.
- o SID: either 4 octet MPLS SID or a 16 octet IPv6 SID.

The following applies to the Type-8 Segment sub-TLV:

- o The Local IPv6 Address MUST be present and represents an adjacency local address.
- o The Remote IPv6 Address MUST be present and represents the remote end of the adjacency.
- o The SID is optional and MAY be of one of the following formats:
 - * MPLS SID: a 4 octet label containing label, TC, S and TTL as defined in [Section 2.4.3.2.1](#).
 - * IPv6 SID: a 16 octet IPv6 address.
- o If length is 34, then only the IPv6 Local and Remote addresses are present.
- o If length is 38, then the IPv6 Local address, IPv4 Remote address and the MPLS SID are present.
- o If length is 50, then the IPv6 Local address, IPv4 Remote address and the IPv6 SID are present.

3. Extended Color Community

The Color Extended Community as defined in [\[I-D.ietf-idr-tunnel-encaps\]](#) is used to steer traffic into a policy.

When the Color Extended Community is used for the purpose of steering the traffic into an SRTE policy, the RESERVED field (as defined in [\[I-D.ietf-idr-tunnel-encaps\]](#)) is changed as follows:

```

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

```

where C0 bits are defined as the "Color-Only" bits.

[\[I-D.filsfils-spring-segment-routing-policy\]](#) defines the influence of these bits on the automated steering of BGP Payload traffic onto SRTE policies.

4. SR Policy Operations

As described in this document, the consumer of a SR Policy NLRI is not the BGP process. The BGP process is in charge of the origination and propagation of the SR Policy NLRI but its installation and use is

outside the scope of BGP
([\[I-D.filsfils-spring-segment-routing-policy\]](#)).

4.1. Configuration and Advertisement of SR TE Policies

Typically, but not limited to, an SR Policy is configured into a controller.

Multiple SR Policy NLRIs may be present with the same <color, endpoint> tuple but with different content when these SR policies are intended to different head-ends.

The distinguisher of each SR Policy NLRI prevents undesired BGP route selection among these SR Policy NLRIs and allow their propagation across route reflectors [[RFC4456](#)].

Moreover, one or more route-target SHOULD be attached to the advertisement, where each route-target identifies one or more intended head-ends for the advertised SR policy.

If no route-target is attached to the SR Policy NLRI, then it is assumed that the originator sends the SR Policy update directly (e.g., through a BGP session) to the intended receiver. In such case, the NO_ADVERTISE community MUST be attached to the SR Policy update.

4.2. Reception of an SR Policy NLRI

On reception of an SR Policy NLRI, a BGP speaker MUST determine if it's first acceptable, then it determines if it is usable.

4.2.1. Acceptance of an SR Policy NLRI

When a BGP speaker receives an SR Policy NLRI from a neighbor it has to determine if it's acceptable. The following applies:

- o The SR Policy NLRI MUST include a distinguisher, color and endpoint field which implies that the length of the NLRI MUST be either 12 or 24 octets (depending on the address family of the endpoint). If the NLRI is not one of the legal lengths, a router supporting this document and that imports the route MUST consider it to be malformed and MUST apply the "treat-as-withdraw" strategy of [[RFC7606](#)].
- o The SR Policy update MUST have either the NO_ADVERTISE community or at least one route-target extended community in IPv4-address format. If a router supporting this document receives an SR policy update with no route-target extended communities and no

NO_ADVERTISE community, the update MUST NOT be sent to the SRTE process. Furthermore, it SHOULD be considered to be malformed, and the "treat-as-withdraw" strategy of [[RFC7606](#)] applied.

- o The Tunnel Encapsulation Attribute MUST be attached to the BGP Update and MUST have the Tunnel Type set to SR Policy (value to be assigned by IANA).
- o Within the SR Policy NLRI, at least one Segment List sub-TLV MUST be present.
- o Within the Segment List sub-TLV at least one Segment sub-TLV MUST be present.

A router that receives an SR Policy update that is not valid according to these criteria MUST treat the update as malformed. The route MUST NOT be passed to the SRTE process, and the "treat-as-withdraw" strategy of [[RFC7606](#)].

The Remote Endpoint and Color sub-TLVs, as defined in [[I-D.ietf-idr-tunnel-encaps](#)], MAY also be present in the SR Policy NLRI encodings. If present, the Remote Endpoint sub-TLV MUST match the Endpoint of the SR Policy SAFI NLRI. If they don't match, the SR Policy advertisement MUST be considered as unacceptable. If present, the Color sub-TLV MUST match the Policy Color of the SR Policy SAFI NLRI. If they don't match, the SR Policy advertisement MUST be considered as unacceptable.

A unacceptable SR Policy update that has a valid NLRI portion with invalid attribute portion MUST be considered as a withdraw of the SR Policy.

A unacceptable SR Policy update that has an invalid NLRI portion MUST trigger a reset of the BGP session.

[4.2.2](#). Usable SR Policy NLRI

If one or more route-targets are present, then at least one route-target MUST match one of the BGP Identifiers of the receiver in order for the update to be considered usable. The BGP Identifier is defined in [[RFC4271](#)] as a 4 octet IPv4 address. Therefore the route-target extended community MUST be of the same format.

If one or more route-targets are present and no one matches any of the local BGP Identifiers, then, while the SR Policy NLRI is acceptable, it is not usable. It has to be noted that if the receiver has been explicitly configured to do so, it MAY propagate the SR Policy NLRI to its neighbors as defined in [Section 4.2.4](#).

Usable SR Policy NLRIs are sent to the Segment Routing Traffic Engineering (SRTE) process. The description of the SRTE process is outside the scope of this document and it's described in [\[I-D.filsfils-spring-segment-routing-policy\]](#).

4.2.3. Passing a usable SR Policy NLRI to the SRTE Process

Once BGP has determined that the SR Policy NLRI is usable, BGP passes the path to the SRTE process ([\[I-D.filsfils-spring-segment-routing-policy\]](#)).

The SRTE process applies the rules defined in [\[I-D.filsfils-spring-segment-routing-policy\]](#) to determine whether a path is valid and to select the best path among the valid paths.

4.2.4. Propagation of an SR Policy

By default, a BGP node receiving an SR Policy NLRI MUST NOT propagate it to any EBGP neighbor.

However, a node MAY be explicitly configured to advertise a received SR Policy NLRI to neighbors according to normal BGP rules (i.e., EBGP propagation by an ASBR or iBGP propagation by a Route-Reflector).

SR Policy NLRIs that have been determined acceptable and valid can be propagated, even the ones that are not usable.

Only SR Policy NLRIs that do not have the NO_ADVERTISE community attached to them can be propagated.

4.3. Flowspec and SR Policies

The SR Policy can be carried in context of a Flowspec NLRI ([\[RFC5575\]](#)). In this case, when the redirect to IP next-hop is specified as in [\[I-D.ietf-idr-flowspec-redirect-ip\]](#), the tunnel to the next-hop is specified by the segment list in the Segment List sub-TLVs. The Segment List (e.g., label stack or IPv6 segment list) is imposed to flows matching the criteria in the Flowspec route to steer them towards the next-hop as specified in the SR Policy SAFI NLRI.

5. Contributors

Arjun Sreekantiah
Cisco Systems
US

Email: asreekan@cisco.com

Dhanendra Jain
Cisco Systems
US

Email: dhjain@cisco.com

Acee Lindem
Cisco Systems
US

Email: acee@cisco.com

Siva Sivabalan
Cisco Systems
US

Email: msiva@cisco.com

Imtiyaz Mohammad
Arista Networks
India

Email: imtiyaz@arista.com

6. Acknowledgments

The authors of this document would like to thank Shyam Sethuram and John Scudder for their comments and review of this document.

7. Implementation Status

Note to RFC Editor: Please remove this section prior to publication, as well as the reference to [RFC 7942](#).

This section records the status of known implementations of the protocol defined by this specification at the time of posting of this Internet-Draft, and is based on a proposal described in [\[RFC7942\]](#). The description of implementations in this section is intended to assist the IETF in its decision processes in progressing drafts to RFCs. Please note that the listing of any individual implementation here does not imply endorsement by the IETF. Furthermore, no effort has been spent to verify the information presented here that was supplied by IETF contributors. This is not intended as, and must not be construed to be, a catalog of available implementations or their features. Readers are advised to note that other implementations may exist.

According to [[RFC7942](#)], "this will allow reviewers and working groups to assign due consideration to documents that have the benefit of running code, which may serve as evidence of valuable experimentation and feedback that have made the implemented protocols more mature. It is up to the individual working groups to use this information as they see fit".

Several early implementations exist and will be reported in detail in a forthcoming version of this document. For purposes of early interoperability testing, when no FCFS code point was available, implementations have made use of the following values:

- o Preference sub-TLV: 6
- o Binding SID sub-TLV: 7
- o Segment List sub-TLV: 128

When IANA-assigned values are available, implementations will be updated to use them.

8. IANA Considerations

This document defines new Sub-TLVs in following existing registries:

- o Subsequent Address Family Identifiers (SAFI) Parameters
- o BGP Tunnel Encapsulation Attribute Tunnel Types
- o BGP Tunnel Encapsulation Attribute sub-TLVs

This document also defines a new registry: "SR Policy List Sub-TLVs".

8.1. Existing Registry: Subsequent Address Family Identifiers (SAFI) Parameters

This document defines a new SAFI in the registry "Subsequent Address Family Identifiers (SAFI) Parameters" that has been assigned by IANA:

Codepoint	Description	Reference

73	SR Policy SAFI	This document

8.2. Existing Registry: BGP Tunnel Encapsulation Attribute Tunnel Types

This document defines a new Tunnel-Type in the registry "BGP Tunnel Encapsulation Attribute Tunnel Types" that has been assigned by IANA:

Codepoint	Description	Reference

15	SR Policy Type	This document

8.3. Existing Registry: BGP Tunnel Encapsulation Attribute sub-TLVs

This document defines new sub-TLVs in the registry "BGP Tunnel Encapsulation Attribute sub-TLVs" to be assigned by IANA:

Codepoint	Description	Reference

TBD3	Preference sub-TLV	This document
TBD4	Binding SID sub-TLV	This document
TBD5	Segment List sub-TLV	This document

8.4. New Registry: SR Policy List Sub-TLVs

This document defines a new registry called "SR Policy List Sub-TLVs". The allocation policy of this registry is "First Come First Served (FCFS)" according to [[RFC5226](#)].

Following Sub-TLV codepoints are defined:

Value	Description	Reference

1	MPLS SID sub-TLV	This document
2	IPv6 SID sub-TLV	This document
3	IPv4 Node and SID sub-TLV	This document
4	IPv6 Node and SID sub-TLV	This document
5	IPv4 Node, index and SID sub-TLV	This document
6	IPv4 Local/Remote addresses and SID sub-TLV	This document
7	IPv6 Node, index and SID sub-TLV	This document
8	IPv6 Local/Remote addresses and SID sub-TLV	This document
9	Weight sub-TLV	This document

9. Security Considerations

TBD.

10. References

10.1. Normative References

[I-D.ietf-idr-tunnel-encaps]
 Rosen, E., Patel, K., and G. Velde, "The BGP Tunnel Encapsulation Attribute", [draft-ietf-idr-tunnel-encaps-06](#) (work in progress), June 2017.

- [I-D.ietf-pce-segment-routing]
Sivabalan, S., Filsfils, C., Tantsura, J., Henderickx, W.,
and J. Hardwick, "PCEP Extensions for Segment Routing",
[draft-ietf-pce-segment-routing-09](#) (work in progress),
April 2017.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate
Requirement Levels", [BCP 14](#), [RFC 2119](#),
DOI 10.17487/RFC2119, March 1997,
<<http://www.rfc-editor.org/info/rfc2119>>.
- [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,
<<http://www.rfc-editor.org/info/rfc4271>>.
- [RFC4360] Sangli, S., Tappan, D., and Y. Rekhter, "BGP Extended
Communities Attribute", [RFC 4360](#), DOI 10.17487/RFC4360,
February 2006, <<http://www.rfc-editor.org/info/rfc4360>>.
- [RFC4760] Bates, T., Chandra, R., Katz, D., and Y. Rekhter,
"Multiprotocol Extensions for BGP-4", [RFC 4760](#),
DOI 10.17487/RFC4760, January 2007,
<<http://www.rfc-editor.org/info/rfc4760>>.
- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an
IANA Considerations Section in RFCs", [RFC 5226](#),
DOI 10.17487/RFC5226, May 2008,
<<http://www.rfc-editor.org/info/rfc5226>>.
- [RFC5575] Marques, P., Sheth, N., Raszuk, R., Greene, B., Mauch, J.,
and D. McPherson, "Dissemination of Flow Specification
Rules", [RFC 5575](#), DOI 10.17487/RFC5575, August 2009,
<<http://www.rfc-editor.org/info/rfc5575>>.
- [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,
<<http://www.rfc-editor.org/info/rfc7606>>.

[10.2.](#) Informational References

[I-D.filsfils-spring-segment-routing-policy]

Filsfils, C., Sivabalan, S., Yoyer, D., Nanduri, M., Lin, S., bogdanov@google.com, b., Horneffer, M., Clad, F., Steinberg, D., Decraene, B., and S. Litkowski, "Segment Routing Policy for Traffic Engineering", [draft-filsfils-spring-segment-routing-policy-00](#) (work in progress), February 2017.

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

Previdi, S., Filsfils, C., Raza, K., Leddy, J., Field, B., daniel.voyer@bell.ca, d., daniel.bernier@bell.ca, d., Matsushima, S., Leung, I., Linkova, J., Aries, E., Kosugi, T., Vyncke, E., Lebrun, D., Steinberg, D., and R. Raszuk, "IPv6 Segment Routing Header (SRH)", [draft-ietf-6man-segment-routing-header-06](#) (work in progress), March 2017.

[I-D.ietf-idr-flowspec-redirect-ip]

Uttaro, J., Haas, J., Texier, M., Andy, A., Ray, S., Simpson, A., and W. Henderickx, "BGP Flow-Spec Redirect to IP Action", [draft-ietf-idr-flowspec-redirect-ip-02](#) (work in progress), February 2015.

[I-D.ietf-spring-segment-routing]

Filsfils, C., Previdi, S., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", [draft-ietf-spring-segment-routing-12](#) (work in progress), June 2017.

[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, <<http://www.rfc-editor.org/info/rfc4456>>.

[RFC7942] Sheffer, Y. and A. Farrel, "Improving Awareness of Running Code: The Implementation Status Section", [BCP 205](#), [RFC 7942](#), DOI 10.17487/RFC7942, July 2016, <<http://www.rfc-editor.org/info/rfc7942>>.

Authors' Addresses

Stefano Previdi (editor)
Cisco Systems, Inc.
IT

Email: stefano@previdi.net

Clarence Filsfils
Cisco Systems, Inc.
Brussels
BE

Email: cfilsfil@cisco.com

Paul Mattes
Microsoft
One Microsoft Way
Redmond, WA 98052
USA

Email: pamattes@microsoft.com

Eric Rosen
Juniper Networks
10 Technology Park Drive
Westford, MA 01886
US

Email: erosen@juniper.net

Steven Lin
Google

Email: stevenlin@google.com

