

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: May 3, 2018

C. Filsfils
S. Sivabalan
K. Raza
J. Liste
F. Clad
Cisco Systems, Inc.
S. Hegde
Juniper Networks, Inc.
D. Yoyer
Bell Canada.
S. Lin
A. Bogdanov
Google, Inc.
M. Horneffer
Deutsche Telekom
D. Steinberg
Steinberg Consulting
B. Decraene
S. Litkowski
Orange Business Services
October 30, 2017

Segment Routing Policy for Traffic Engineering
draft-filsfils-spring-segment-routing-policy-03.txt

Abstract

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. The headend node steers a flow into an 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. This document details the concepts of SR Policy and steering into an SR Policy.

Requirements Language

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

Status of This Memo

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

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

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

This Internet-Draft will expire on May 3, 2018.

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 (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction	3
2.	SR Traffic Engineering Architecture	4
3.	SR Policy	5
4.	SR Policy Active Path Selection Examples	8
5.	Segment-list	10
5.1.	Explicit Null	10
6.	SR Policy Multi-Domain Database	11
7.	Operations	11
7.1.	W-ECMP	11
7.2.	Path Validation	12
7.3.	Fast Convergence	12
8.	Binding SID	13
8.1.	Benefits	13
8.2.	Allocation	14
8.2.1.	Dynamic BSID Allocation	14
8.2.2.	Explicit BSID Allocation	15
8.2.3.	Generic BSID Allocation	15
9.	Centralized Discovery	16
10.	Dynamic Path	17

10.1.	Optimization Objective	17
10.2.	Constraints	18
10.3.	SR Native Algorithm	19
10.4.	Path to SID	19
10.5.	PCE Computed Path	20
11.	Signaling Paths of an SR Policy to a Head-end	20
11.1.	BGP	20
11.2.	PCEP	21
11.3.	NETCONF	21
11.4.	CLI	21
12.	Steering into an SR Policy	21
12.1.	Incoming Active SID is a BSID	21
12.2.	Recursion on a BSID	22
12.2.1.	Multiple Colors	23
12.3.	Recursion on an on-demand dynamic BSID	23
12.3.1.	Multiple Colors	23
12.4.	An array of BSIDs associated with an IGP entry	23
12.5.	A Routing Policy on a BSID	24
13.	Optional Steering Modes for BGP Destinations	25
13.1.	Color-Only BGP Destination Steering	25
13.2.	Drop on Invalid	25
14.	Multipoint SR Policy	26
14.1.	Spray SR Policy	26
15.	Reporting SR Policy	26
16.	Work in Progress	26
17.	Acknowledgement	27
18.	Normative References	27
	Authors' Addresses	28

[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.

This document details the concepts of SR Policy and steering into an SR Policy. These apply equally to the MPLS and SRv6 instantiations of segment routing.

For reading simplicity, the illustrations are provided for the MPLS instantiations.

2. SR Traffic Engineering Architecture

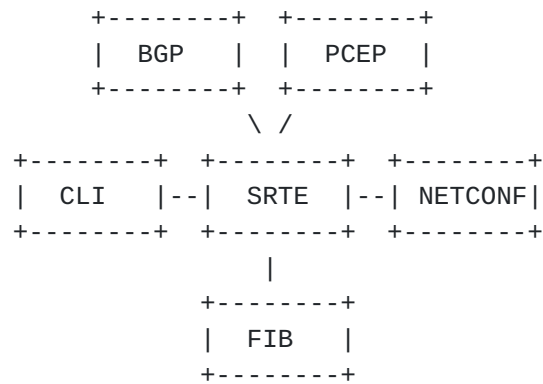


Figure 1: SR Policy architecture

The Segment Routing Traffic Engineering (SRTE) process installs a Segment Routing Policy (SR Policy) in the forwarding plane (FIB).

An SR policy is represented in FIB as a BSID-keyed entry. For traffic steering purpose, suitable SR policy is identified using either BSID or IP prefix.

For a given SR policy, the SRTE process MAY learn multiple candidate paths from different sources: NETCONF with OpenConfig or YANG model (work in progress), PCEP [[I-D.ietf-pce-pce-initiated-lsp](#)], local configuration or BGP [[I-D.previdi-idr-segment-routing-te-policy](#)].

The SRTE process selects the best candidate path and installs it in FIB.

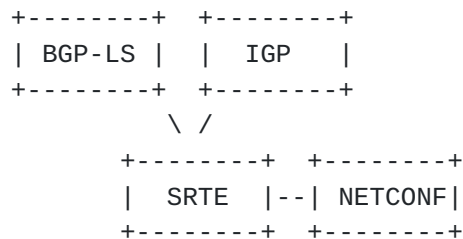


Figure 2: Topology/link-state database architecture

The SRTE process maintains an SRTE database (SRTE-DB).

The SRTE-DB includes information that a headend requires to compute path (if possible) for SR policies as well as setup and validate SR policies.

The information in the SRTE-DB includes (but not limited to):

- o Regular IGP information (topology, IGP metrics).
- o Segment Routing information (such as SRGB, Prefix-SIDs, Adj-SIDs).
- o TE Link Attributes (such as TE metric, SRLG, attribute-flag, extended admin group).

The SRTE-DB is multi-domain capable.

The attached domain topology MAY be learned via IGP, BGP-LS or NETCONF.

A non-attached (remote) domain topology MAY be learned via BGP-LS or NETCONF.

In some use-cases, the SRTE-DB may only contain the attached domain topology while in others, the SRTE-DB may contain the topology of multiple domains.

3. SR Policy

An SR Policy is identified through the following tuple:

- o The head-end where the policy is instantiated/implemented.
- o The endpoint (i.e.: the destination of the policy).
- o The color (an arbitrary numerical value).

At a given head-end, an SR Policy is fully identified by the <color, endpoint> tuple.

An endpoint can be specified as an IPv4 or IPv6 address (including null address, i.e., 0.0.0.0 and ::0 for IPv4 and IPv6 address families respectively as explained later in this document).

An SR Policy is associated with one or more candidate paths. A path refers to a list of segments (i.e. Segment-list) or a set of Segment-lists. A Segment-list represents a specific way to send traffic from the head-end to the end-point of the corresponding SR policy. If a path contains multiple Segment-lists, each list can be

associated with a weight for equal or unequal cost load-balancing (default is equal cost load-balancing).

For each SR Policy, at most one candidate path is selected, and only the selected path is used for forwarding traffic that is being steered onto that policy.

A candidate path is either dynamic or explicit.

A dynamic path expresses an optimization objective and a set of constraints. The headend computes a solution to the optimization problem as a Segment-list or a set of Segment-lists. When the headend does not have enough topological information (e.g. multi-domain problem), the headend may delegate the computation to a PCE. Whenever the network state changes, the path is recomputed.

An explicit path is a Segment-list or a set of Segment-lists.

A candidate path has a preference. If not specified, the default preference is 100.

A candidate path is associated with a single Binding SID (BSID) in the context of the corresponding SR policy. If a candidate path is selected (active path), its BSID must be the same as that of the corresponding SR policy. If more than one SR policies might happen to be associated with an identical candidate path, each candidate path MUST be associated with a unique BSID to ensure that each policy has a distinct BSID.

A candidate path is valid if it is usable. A common path validity criterion is the reachability of its constituent SIDs. The validation rules are defined in a later section.

A Path is selected (i.e. it is the best path of the policy) when it is valid and its preference is the best (highest value) among all the candidate paths of the SR Policy. The selected path is referred to as the "active path" of the SR policy in this document.

Whenever a new path is learned or the validity of an existing path changes or an existing path is changed, the selection process must be re-executed.

A headend may be informed about a candidate path for a policy <color, endpoint> by various means including: local configuration, NETCONF, PCEP or BGP. The protocol source of the path does not matter as to how an active path is chosen. For a given Policy, when comparing a candidate path learned via one means to a candidate path learned via another means (e.g., one via BGP another via PCEP), a valid path will

be regarded as preferable to the other based on the preference. 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. Refer to examples later in this document.

In the vast majority of use-cases known to date, a path is associated with a single Segment-list and each path of a policy has a different preference.

The BSID of an SR Policy refers to its selected path.

At any given time, a given BSID MUST map to a single SR policy and indirectly map to its selected path. However, the mapping from a given BSID to an SR Policy may change over the life of the SR policy, and the true identification of a policy is the tuple <headend, endpoint, color>.

An SR Policy <color, endpoint> is active at a headend as soon as this head-end knows about a valid path for this policy.

An active SR Policy installs a BSID-keyed entry in the forwarding plane with the action of steering the packets matching this entry to the selected path of the SR Policy.

If a set of Segment-lists is associated with the selected path of the policy, then the steering is flow and W-ECMP based according to the relative weight of each Segment-list.

In summary, the information model is the following:

```
SR policy F00
  Candidate-paths
    path preference 200 (selected)
      BSID1
      Weight W1, Segment-list1: SID11...SID1i
      Weight W2, Segment-list2: SID21...SID2j
    path preference 100
      BSID2
      Weight W3, Segment-list3: SID31...SID3i
      Weight W4, Segment-list4: SID41...SID4j
```

The numbers 200 and 100 are preferences of the paths associated with the policy.

In general $BSID_{in} = BSID_1 = BSID_2 \dots$. If paths of an SR policy have different BSIDs, then the BSID of the SR policy is that of the selected path.

4. SR Policy Active Path Selection Examples

Example 1:

Consider headend H where two candidate paths of the same SR Policy (endpoint, color) are signaled via BGP and whose respective NLRIs have the same route distinguishers:

NLRI A with distinguisher = RD1, color = C, endpoint = N, preference P1.

NLRI B with distinguisher = RD1, color = C, endpoint = N, preference P2.

- o Because the NLRIs are identical (same distinguisher), BGP will perform bestpath selection. Note that there are no changes to BGP best path selection algorithm.
- o H installs one advertisement as bestpath into the BGP table.
- o A single advertisement is passed to the SRTE process.
- o SRTE process does not perform any path selection.

Note that the candidate path's preference value do not have any effect on the BGP bestpath selection process.

Example 2:

Consider headend H where two candidate paths of the same SR Policy (endpoint, color) are signaled via BGP and whose respective NLRIs have different route distinguishers: NLRI A with distinguisher = RD1, color = C, endpoint = N, preference P1. NLRI B with distinguisher = RD2, color = C, endpoint = N, preference P2.

- o Because the NLRIs are different (different distinguisher), BGP will not perform bestpath selection.
- o H installs both advertisements into the BGP table.
- o Both advertisements are passed to the SRTE process.
- o SRTE process at H selects the candidate path advertised by NLRI B as the active path for the SR policy since P2 is greater than P1.

Note that the recommended approach is to use NLRIs with different distinguishers when several candidate paths for the same SR Policy (endpoint, color) are signaled via BGP to a headend.

Example 3:

Consider that a headend H learns two candidate paths of the same SR Policy (endpoint, color); one signaled via BGP and another via Local configuration.

NLRI A with distinguisher = RD1, color = C, endpoint = N, preference P1.

Local "foo" with color = C, endpoint = N, preference P2.

- o H installs NLRI A into the BGP table.
- o NLRI A and "foo" are both passed to the SRTE process.
- o SRTE process at H selects the candidate path indicated by "foo" as the active path for the SR policy since P2 is greater than P1.

When an SR Policy has multiple valid candidate paths with the the same best preference, the SRTE process at a headend MAY keep the oldest candidate path as the active path as explained in the following examples:

Example 1:

Consider headend H with two candidate paths of the same SR Policy (endpoint, color) and the same preference value.

- o NLRI A with distinguisher RD1, color C, endpoint N, preference P1 (selected as active path at time t0).
- o NLRI B with distinguisher RD2, color C, endpoint N, preference P1 (passed to SRTE process at time t1).

After t1, SRTE process at H retains candidate path associated with NLRI A as active path of the SR policy.

Example 2:

Consider headend H with two candidate paths of the same SR Policy (endpoint, color) and the same preference value.

- o Local "foo" with color C, endpoint N, preference P1 (selected as active path at time t0).
- o NLRI A with distinguisher RD1, color C, endpoint N, preference P1 (passed to SRTE process at time t1).

After t1, SRTE process at H retains candidate path associated with Local candidate path "foo" as active path of the SR policy.

Note that a headend node MUST NOT install in FIB the "merged" set of segment-lists associated with the candidate paths with the best preference.

5. Segment-list

A Segment-list includes segments of different types (1 to 8) and an optional weight value that is used for W-ECMP.

The following segment types are defined:

- 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.

The optional SID can be an MPLS label (SR applied to the MPLS dataplane) or an IPv6 SID (SRv6, SR applied to the IPv6 dataplane).

When SR is applied to MPLS, a type 1 SID may be any MPLS label including (but not limited to): segment types described in [\[I-D.ietf-spring-segment-routing\]](#), static MPLS labels and /or reserved labels.

When building the MPLS label stack or the IPv6 Segment list from the Segment List, the node instantiating the policy MUST interpret the set of Segments as follows:

- o The first Segment represents the topmost label or the first IPv6 segment. It identifies the first segment the traffic will be directed toward along the SR explicit path.
- o The last Segment represents the bottommost label or the last IPv6 segment the traffic will be directed toward along the SR explicit path.

A Segment-list is represented as <S1, S2, ... Sn> where S1 is the first SID.

5.1. Explicit Null

A Type 1 SID may be any MPLS label, including reserved labels.

For example, assuming that the desired traffic-engineered path from a headend 1 to an endpoint 4 can be expressed by the Segment-list <16002, 16003, 16004> where 16002, 16003 and 16004 respectively refer to the IPv4 Prefix SIDs bound to node 2, 3 and 4, then IPv6 traffic can be traffic-engineered from nodes 1 to 4 via the previously described path using an SR Policy with Segment-list <16002, 16003,

16004, 2> where mpls label value of 2 represents the "IPv6 Explicit NULL Label".

The penultimate node before node 4 will pop 16004 and will forward the frame on its directly connected interface to node 4.

The endpoint receives the traffic with top label "2" which indicates that the payload is an IPv6 packet.

When steering unlabeled IPv6 BGP destination traffic using an SR policy composed of Segment-list(s) based on IPv4 SIDs, the headend node SHOULD automatically impose the "IPv6 Explicit NULL Label" as bottom of stack label. Refer to "Steering" section later in this document.

6. SR Policy Multi-Domain Database

A headend can learn an attached domain topology via its IGP or a BGP-LS session. A headend can learn a non-attached domain topology via a BGP-LS session.

A headend collects all these topologies in the SR-TE database (SRTE-DB).

The SRTE-DB is multi-domain capable.

In some deployments, the SRTE-DB may only contain the attached domain topology while in others, the SRTE-DB may contain the topology of multiple domains.

7. Operations

7.1. W-ECMP

Packets steered to an SR Policy (i.e. to its BSID either via presence in the packet header as active segment or via FIB recursion) are load-balanced on a weighted basis among the Segment-lists associated with the selected path of the SR Policy.

The fraction of the flows associated with a given Segment-list is w/S_w where w is the weight of the Segment-list and S_w is the sum of the weights of the Segment-lists of the selected path of the SR Policy.

The accuracy of the weighted load-balancing depends on the platform implementation.

7.2. Path Validation

A Segment-list MUST be declared invalid when:

- o It is empty.
- o The headend is unable to resolve the first SID into one or more outgoing interface(s) and next-hop(s).
- o The headend is unable to resolve any non-first SID of type 3-to-8 into an MPLS label or an SRv6 SID.

Unreachable means that the headend has no path to the SID in its SRTE-DB.

In multi-domain deployments, it is expected that the headend be unable to verify the reachability of the SIDs in remote domains. Types 1 and 2 MUST be used for the SIDs for which the reachability cannot be verified. Note that the first SID must always be reachable whatever is type.

In addition, a Segment-list MAY be declared invalid when:

- o Its last segment is not a Prefix SID (including BGP Peer Node-SID) advertised by the node specified as the endpoint of the corresponding SR policy.
- o Its last segment is not an Adjacency SID (including BGP Peer Adjacency SID) of any of the links present on neighbor nodes and that terminate on the node specified as the endpoint of the corresponding SR policy.

A Path is invalid as soon as it has no valid Segment-list.

The headend of an SR Policy updates the validity of a Segment-list upon network topological change.

A path of an SR Policy is invalid when all its Segment-lists are invalid.

An SR Policy is invalid when all its paths are invalid.

7.3. Fast Convergence

Upon topological change, many policies could be recomputed. An implementation MAY provide a per-policy priority field. The operator MAY set this field to indicate in which order the policies should be re-computed. Such a priority may be represented by an integer in the range [0, 254] where the lowest value is the highest priority.

8. Binding SID

8.1. Benefits

The Binding SID (BSID) is fundamental to Segment Routing. It provides scaling, network opacity and service independence.

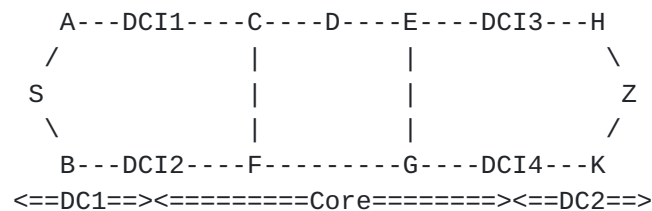


Figure 3: A Simple Datacenter Topology

A simplified illustration is provided on the basis of the previous diagram where we assume that S, A, B, Data Center Interconnect DCI1 and DCI2 share the same IGP-SR instance in the data-center 1 (DC1). DCI1, DCI2, C, D, E, F, G, DCI3 and DCI4 share the same IGP-SR domain in the core. DCI3, DCI4, H, K and Z share the same IGP-SR domain in the data-center 2 (DC2).

In this example, we assume no redistribution between the IGP's and no presence of BGP. The inter-domain communication is only provided by SR through SR Policies.

The latency from S to DCI1 equals to DCI2. The latency from Z to DCI3 equals to DCI4. All the intra-DC links have the same IGP metric 10.

The path DCI1, C, D, E, DCI3 has a lower latency and lower capacity than the path DCI2, F, G, DCI4.

The IGP metrics of all the core links are set to 10 except the links D-E which is set to 100.

A low-latency multi-domain policy from S to Z may be expressed as <DCI1, BSID, Z> where:

- o DCI1 is the prefix SID of DCI1.
- o BSID is the Binding SID bound to an SR policy <D, D2E, DCI3> instantiated at DCI1.
- o Z is the prefix SID of Z.

Without the use of an intermediate core SR Policy (efficiently summarized by a single BSID), S would need to steer its low-latency flow into the policy <DCI1, D, D2E, DCI3, Z>.

The use of a BSID (and the intermediate bound SR Policy) decreases the number of segments imposed by the source.

A BSID acts as a stable anchor point which isolates one domain from the churn of another domain. Upon topology changes within the core of the network, the low-latency path from DCI1 to DCI3 may change. While the path of an intermediate policy changes, its BSID does not change. Hence the policy used by the source does not change, hence the source is shielded from the churn in another domain.

A BSID provides opacity and independence between domains. The administrative authority of the core domain may not want to share information about its topology. The use of a BSID allows keeping the service opaque. S is not aware of the details of how the low-latency service is provided by the core domain. S is not aware of the need of the core authority to temporarily change the intermediate path.

8.2. Allocation

There are three approaches to allocate a BSID to an SR Policy: all the paths have no explicit BSID (called dynamic allocation), all the paths have the same explicit BSID (explicit allocation) and finally a mix of paths with and without explicit BSID (generic allocation).

In practice, all the use-cases seen to-date either use the explicit allocation or the dynamic allocation. The explicit allocation is most-often associated with controller-instantiated SR Policies. The dynamic allocation is most-often associated with router-based on-demand SR Policies.

8.2.1. Dynamic BSID Allocation

No path of the SR Policy have a specified BSID.

In such a case, the SR-TE implementation allocates a SID to the SR Policy and keeps it along the whole existence of the policy.

In the case of SR-MPLS, the SR-TE implementation binds a local dynamic label in the same way LDP, RSVP-TE or BGP would do.

8.2.2. Explicit BSID Allocation

All the paths of the SR Policy have the same specified BSID, with the same behavioral preference in case this specified BSID is not available.

If the specified BSID is available, then it is bound to the SR Policy and used along the existence of the policy.

If the specified BSID is not available, then a SYSLOG/NETCONF message is generated and if the preferred behavior is to fall-back on the dynamic allocation, then the dynamic allocation is performed.

If the specified BSID is not available and the operator-requested behavior is to not fall-back on the dynamic allocation, then a SYSLOG/NETCONF message is generated and the SR Policy does not install any BSID entry in the forwarding plane.

A later section will explain how controllers can discover the local SIDs available at a node N so as to pick an explicit BSID for a SR Policy to be instantiated at headend N.

8.2.2.1. Explicit BSID Allocation within a Label Block

A headend node MAY validate BSID allocation within a label block (e.g., SRLB). In this case the following applies:

- o If the specified BSID falls outside this block, then a SYSLOG/NETCONF message is generated and the SR Policy MUST NOT install any BSID entry in the forwarding plane.
- o If the specified BSID falls inside this block but conflicts with either a BSID of another policy or with another application (e.g. explicit Adjacency SID), then a SYSLOG/NETCONF message is generated and the SR Policy MUST NOT install any BSID entry in the forwarding plane.

8.2.3. Generic BSID Allocation

This section details the BSID allocation when a policy is made of paths with different BSID allocation behaviors (e.g., mix of paths with and without an explicit BSID, potentially with different explicit BSIDs).

When the selected path has a specified BSID, the SR Policy uses that BSID if this value (label in MPLS, IPv6 address in SRv6) is available (i.e. not associated with any other usage: e.g. to another MPLS client, to another SID, to another SR Policy).

If the selected path's BSID is not available, then the SR Policy keeps the previous BSID. If the SR Policy did not have a previous BSID, then the SR Policy dynamically binds a BSID to itself.

Note that a path may request that only its specified BSID be used. In that case, if that BSID is not available and that path is active, then no BSID is bound to the policy and a SYSLOG/NETCONF is triggered. In this case, the SR Policy does not install any entry indexed by a BSID in the forwarding plane.

When an SR Policy has multiple valid paths with the best preference but with different BSIDs, it is left to the implementation to decide which BSID to install. This case is unlikely in practice for two reasons. First, all known use-cases share the same BSID across all the paths of a given SR Policy. Second, all known use-cases have a different preference for each path. Hence in practice a single path will be active and with a stable BSID on a per-policy basis.

9. Centralized Discovery

This section explains how controllers can discover the local SIDs available at a node N so as to pick an explicit BSID for a SR Policy to be instantiated at headend N.

Any controller can discover the following properties of a node N (e.g. via BGP-LS, NETCONF etc.):

- o its local Segment Routing Label Block (SRLB).
- o its local topology.
- o its topology-related SIDs (Adj SID and EPE SID).
- o its SR Policies and their BSID ([\[I-D.ietf-idr-te-lsp-distribution\]](#)).

Any controller can thus infer the available SIDs in the SRLB of any node.

As an example, a controller discovers the following characteristics of N: SRLB [4000, 8000], 3 Adj SIDs (4001, 4002, 4003), 2 EPE SIDs (4004, 4005) and 3 SRTE policies (whose BSIDs are respectively 4006, 4007 and 4008). This controller can deduce that the SRLB sub-range [4009, 5000] is free for allocation.

Likely, the next question is: how do we ensure that different controllers do not pick the same available SID at the same time for different SR Policies.

Clearly, a controller is not restricted to use the next numerically available SID in the available SRLB sub-range. It can pick any label

in the subset of available labels. This random pick make the chance for a collision unlikely.

An operator could also sub-allocate the SRLB between different controllers (e.g. [4000-4499] to controller 1 and [4500-5000] to controller 2).

Inter-controller state-synchronization may be used to avoid/detect collision in BSID.

All these techniques make the likelihood of a collision between different controllers very unlikely.

In the unlikely case of a collision, the controllers will detect it through SYSLOG/NETCONF, BGP-LS reporting ([[I-D.ietf-idr-te-lsp-distribution](#)]) or PCEP notification. They then have the choice to continue the operation of their SR Policy with the dynamically allocated BSID or re-try with another explicit pick.

Note: in deployments where PCE Protocol (PCEP) is used between head-end and controller (PCE), a head-end can report BSID as well as policy attributes (e.g., type of disjointness) and operational and administrative states to controller. Similarly, a controller can also assign/update the BSID of a policy via PCEP when instantiating or updating SR Policy.

[10.](#) Dynamic Path

A dynamic path is a path that expresses an optimization objective and constraints.

The headend of the policy is responsible to compute a Segment-list ("solution Segment-list") that fits this optimization problem. The headend is responsible for computing the solution Segment-list any time the inputs to the problem change (e.g. topology changes).

[10.1.](#) Optimization Objective

We define two optimization objectives:

- o Min-Metric - requests computation of a solution Segment-list optimized for a selected metric.
- o Min-Metric with margin and maximum number of SIDs - Min-Metric with two changes: a margin of by which two paths with similar metrics would be considered equal, a constraint on the max number of SIDs in the Segment-list.

The "Min-Metric" optimization objective requests to compute a solution Segment-list such that packets flowing through the solution Segment-list use ECMP-aware paths optimized for the selected metric. The "Min-Metric" objective can be instantiated for the IGP metric xor the TE metric xor the latency extended TE metric. This metric is called the O metric (the optimized metric) to distinguish it from the IGP metric. The solution Segment-list must be computed to minimize the number of SIDs and the number of Segment-lists.

If the selected O metric is the IGP metric and the headend and tailend are in the same IGP domain, then the solution Segment-list is made of the single prefix-SID of the tailend.

When the selected O metric is not the IGP metric, then the solution Segment-list is made of prefix SIDs of intermediate nodes, Adjacency SIDs along intermediate links and potentially BSIDs of intermediate policies.

In many deployments there are insignificant metric differences between mostly equal path (e.g. a difference of 100 usec of latency between two paths from NYC to SFO would not matter in most cases). The "Min-Metric with margin" objective supports such requirement.

The "Min-Metric with margin and maximum number of SIDs" optimization objective requests to compute a solution Segment-list such that packets flowing through the solution Segment-list do not use a path whose cumulated O metric is larger than the shortest-path O metric + margin.

If this is not possible because of the number of SIDs constraint, then the solution Segment-list minimizes the O metric while meeting the maximum number of SID constraints.

10.2. Constraints

The following constraints can be defined:

- o Inclusion and/or exclusion of TE affinity.
- o Inclusion and/or exclusion of IP address.
- o Inclusion and/or exclusion of SRLG.
- o Inclusion and/or exclusion of admin-tag.
- o Maximum accumulated metric (IGP, TE and latency).
- o Maximum number of SIDs in the solution Segment-list.
- o Maximum number of weighted Segment-lists in the solution set.
- o Diversity to another service instance (e.g., link, node, or SRLG disjoint paths originating from different head-ends).

10.3. SR Native Algorithm

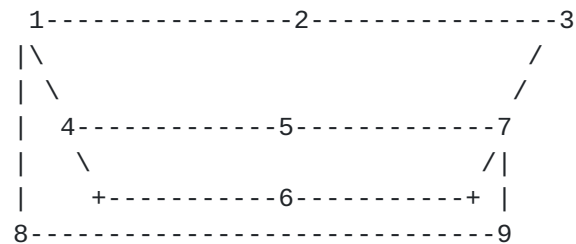


Figure 4: Illustration used to describe SR native algorithm

Let us assume that all the links have the same IGP metric of 10 and let us consider the dynamic path defined as: `Min-Metric(from 1, to 3, IGP metric, margin 0)` with constraint "avoid link 2-to-3".

A classical circuit implementation would do: prune the graph, compute the shortest-path, pick a single non-ECMP branch of the ECMP-aware shortest-path and encode it as a Segment-list. The solution Segment-list would be `<4, 5, 7, 3>`.

An SR-native algorithm would find a Segment-list that minimizes the number of SIDs and maximize the use of all the ECMP branches along the ECMP shortest path. In this illustration, the solution Segment-list would be `<7, 3>`.

In the vast majority of SR use-cases, SR-native algorithms should be preferred: they preserve the native ECMP of IP and they minimize the dataplane header overhead.

In some specific use-case (e.g. TDM migration over IP where the circuit notion prevails), one may prefer a classic circuit computation followed by an encoding into SIDs.

SR-native algorithms are a local node behavior and are thus outside the scope of this document.

10.4. Path to SID

Let us assume the below diagram where all the links have an IGP metric of 10 and a TE metric of 10 except the link AB which has an IGP metric of 20 and the link AD which has a TE metric of 100. Let us consider the `min-metric(from A, to D, TE metric, margin 0)`.

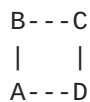


Figure 5: Illustration used to describe path to SID conversion

The solution path to this problem is ABCD.

This path can be expressed in SIDs as $\{B, D\}$; where B and D are the IGP prefix SIDs respectively associated with nodes B and D in the diagram.

Indeed, from A, the IGP path to B is AB (IGP metric 20 better than ADCB of IGP metric 30). From B, the IGP path to D is BCD (IGP metric 20 better than BAD of IGP metric 30).

While the details of the algorithm remain a local node behavior, a high-level description follows: start at the headend and find an IGP prefix SID that leads as far down the desired path as possible (without using any link not included in the desired path). If no prefix SID exists, use the Adj SID to the first neighbor along the path. Restart from the node that was reached.

10.5. PCE Computed Path

A local computation should be preferred whenever possible. When local computation is not possible (e.g., a policy's tail-end is outside the topology known to the head-end), the head-end may send path computation request to a PCE supporting PCEP extension specified in [\[I-D.ietf-pce-segment-routing\]](#).

11. Signaling Paths of an SR Policy to a Head-end

A headend H can be informed about a candidate path for an SR policy (endpoint, color) via several means: BGP, PCEP, CLI, netconf.

We remind that the selection of the best path for a policy is independent of the protocol source of the path.

11.1. BGP

Please refer to [\[I-D.previdi-idr-segment-routing-te-policy\]](#)

11.2. PCEP

Please refer to [[I-D.ietf-pce-pce-initiated-lsp](#)]

11.3. NETCONF

Operator MUST be able to install policy via NETCONF with OpenConfig/YANG models (work in progress).

11.4. CLI

Operator MUST be able to install policy via CLI.

12. Steering into an SR Policy

A headend can steer a packet flow on an SR Policy in various ways:

- o Incoming packets have an active SID matching a local BSID at the head-end.
- o Incoming packets match a BGP/Service route which recurses on the BSID of a local policy.
- o Incoming packets match a BGP/Service route which recurses on an array of paths to the BGP nhop where some of the paths in the array are local SR Policies.
- o Incoming packets match a routing policy which directs them on a local SR policy.

For simplicity of illustration, we will use the SR-MPLS example.

12.1. Incoming Active SID is a BSID

Let us assume that headend H has a local SR Policy P of Segment-list <S1, S2, S3> and BSID B.

When H receives a packet with label stack <B, L2, L3>, H pops B and pushes <S1, S2, S3>. H sends the resulting packet with label stack <S1, S2, S3, L2, L3> along the path to S1.

The previous is an example assuming that S1 is the Prefix-SID of a remote node in the network.

Actual programming of the first segment in the segment-list MUST take into account the type of segment represented by the SID.

To illustrate the previous point, let us assume that headend H has a local SR Policy Q of Segment-list <S1, S2, S3> (where S1 is an Adj-SID at H) and BSID B. When H receives a packet with label stack <B, L2, L3>, H pops B and pushes <S2, S3>. H sends the resulting packet

with label stack <S2, S3, L2, L3> along the outgoing interface(s) associated with Adj-SID S1.

H has steered the packet in the policy P.

H did not have to classify the packet. The classification was done by a node upstream of H (e.g. the source of the packet or an intermediate ingress edge node of the SR domain) and the result of this classification was efficiently encoded in the packet header as a BSID.

This is another key benefit of the segment routing in general and the binding SID in particular: the ability to encode a classification and the resulting steering in the packet header such as to better scale and simplify intermediate aggregation nodes.

12.2. Recursion on a BSID

Let us assume that headend H:

- o learns about a BGP route R/r via next-hop N, extended-color community C and label V.
- o has a local SR Policy P to (endpoint = N, color = C) of Segment-list <S1, S2, S3> and BSID B.
- o has a local BGP policy which matches on the extended-color community C and allows its usage as an SR-TE SLA steering information.

In such a case, H installs R/r in RIB/FIB with next-hop = B (instead of N).

Indeed, H's local BGP policy and the received BGP route indicate that the headend should associate R/r with an SR-TE path to N with the SLA associated with color C. The headend therefore installs the BGP route on that policy.

This can be implemented by using the BSID as a generalized nhop and installing the BGP route on that generalized next-hop.

When H receives a packet with a destination matching R/r, H pushes the label stack <S1, S2, S3, V> and sends the resulting packet along the path to S1.

Note that any label associated with the BGP route is pushed after the Segment-list of the SR Policy.

12.2.1. Multiple Colors

When a BGP route has multiple extended-color communities each with a valid SRTE policy, the BGP process installs the route on the Binding SID corresponding to the SR policy whose color is of highest numerical value.

Let us assume that headend H:

- o learns about a BGP route R/r via next-hop N, extended-color communities C1 and C2 and label V.
- o has a local SR Policy P1 to (endpoint = N, color = C1) of SID list <S1, S2, S3> and BSID B1.
- o has a local SR Policy P2 to (endpoint = N, color = C2) of SID list <S4, S5, S6> and BSID B2.
- o has a local BGP policy which matches on the extended-color communities C1 and C2 and allows their usage as an SR-TE SLA steering information.

In such a case, H installs R/r in RIB/FIB with next-hop = B2 (instead of N) because C2 > C1.

12.3. Recursion on an on-demand dynamic BSID

In the previous section, we assumed that H had a pre-established "explicit" SR Policy (endpoint N, color C).

In this section, independently to the a-priori existence of any explicit path of the SR policy (N, C), we note that the BGP process at node H triggers the SRTE process at node H to instantiate a dynamic path for the SR policy (N, C) as soon as:

- o the BGP process learns of a route R/r via N and with color C.
- o a local policy at node H authorizes the on-demand SRTE path instantiation and maps the color to a dynamic SRTE optimization template.

12.3.1. Multiple Colors

When a BGP route R/r via N has multiple extended-color communities Ci (with i=1 ... n), an individual on-demand SR-TE dynamic path request (endpoint N, color Ci) is triggered for each color Ci.

12.4. An array of BSIDs associated with an IGP entry

Let us assume that head-end H:

- o learns about a BGP route R/r via next-hop N and label V.

- o has a local SR Policy P1 to (endpoint = N, color = C1) of Segment-list <S1, S2, S3> and BSID B1.
- o has a local SR Policy P2 to (endpoint = N, color = C2) of Segment-list <S4, S5, S6> and BSID B2.
- o is configured to instantiate an array of paths to N where the entry 0 is the IGP path to N, color C1 is the first entry and Color C2 is the second entry. The index into the array is called a Forwarding Class (FC). The index can have values 0 to 7.
- o is configured to match flows in its ingress interfaces (upon any field such as Ethernet destination/source/vlan/tos or IP destination/source/DSCP or transport ports etc.) and color them with an internal per-packet forwarding-class variable (0, 1 or 2 in this example).

In such a case, H installs in RIB/FIB:

- o R/r in with next-hop N (as usual).
- o N via a recursion on an array A (instead of the immediate outgoing link associated with the IGP shortest-path to N.
- o Entry A(0) set to the immediate outgoing link of the IGP shortest-path to N.
- o Entry A(1) set to B1.
- o Entry A(2) set to B2.

H receives three packets P, P1 and P2 on its incoming interface. H colors them respectively with forwarding-class 0, 1 and 2. As a result:

- o H pushes <V> on packet P and forwards the resulting frame along the shortest-path to N (which in SR-MPLS results in the pushing of the prefix-SID of N.
- o H pushes <S1, S2, S3, V> on packet P1 and forwards the resulting frame along the shortest-path to S1.
- o H pushes <S4, S5, S6, V> on packet P2 and forwards the resulting frame along the shortest-path to S4.

If the local configuration does not specify any explicit forwarding information for an entry of the array, then this entry is filled with the same information as entry 0 (i.e. the IGP shortest-path).

This realizes per-flow steering: different flows bound to the same BGP destination R/r are steered on different SR-TE paths.

12.5. A Routing Policy on a BSID

Finally, headend H may be configured with a local routing policy which overrides any BGP/IGP path and steer a specified flow on an SR Policy.

13. Optional Steering Modes for BGP Destinations

13.1. Color-Only BGP Destination Steering

In the previous section "Recursion on a BSID", we have seen that the steering on an SR Policy is governed by the matching of the BGP route's next-hop N and the authorized color C with a local SR Policy defined by the tuple (N, C).

This is the most likely form of BGP destination steering and the one we recommend.

In this section, we define an alternative steering mechanism based only on the color.

This color-only steering variation is governed by two new flags "C" and "O" defined in the color extended community.

The Color-Only flags "CO" are set to 00 by default.

When 00, the BGP destination is preferably steered onto a valid SR Policy (N, C) where N is an IPv4/6 endpoint address and C is a color value else it is steered on the IGP path to the next-hop N. This is the classic case we described before and that we recommend.

When 01, the BGP destination is preferably steered onto a valid SR Policy (N, C) else onto a valid SR Policy (null endpoint, C) of the same address-family of N else on any valid SR Policy (any null endpoint, C) else on the IGP path to the next-hop N.

When 10, the BGP destination is preferably steered onto a valid SR Policy (N, C) else onto a valid SR Policy (null endpoint, C) of the same address-family of N else on any valid SR Policy (any null endpoint, C) else on any valid SR Policy (any endpoint, C) of the same address-family of N else on any valid SR Policy (any endpoint, C) else on the IGP path to the next-hop N.

The null endpoint is 0.0.0.0 for IPv4 and ::0 for IPv6 (all bits set to the 0 value).

When 11, it is treated like 00.

13.2. Drop on Invalid

The local BGP policy authorizing the use of an extended color community steering on an SR policy may specify that if the related SR Policy becomes invalid then the related BSID should remain in RIB/FIB and point to null0 (drop any packet recursing on that BSID).

Recall that, by default, for a BGP route R/r via next-hop N with extended-color community C, when the SR Policy (N, C) becomes invalid, then BGP re-installs R/r in RIB/FIB via N (the IGP path to N).

14. Multipoint SR Policy

14.1. Spray SR Policy

A Spray SR-TE policy is a variant of an SR-TE policy which involves packet replication.

Any traffic steered into a Spray SR Policy is replicated along the Segment-lists of its selected path.

In the context of a Spray SR Policy, the selected path SHOULD have more than one Segment-list. The weights of the Segment-lists is not applicable for a Spray SR Policy. They MUST be set to 1.

Like any SR policy, a Spray SR Policy has a BSID instantiated into the forwarding plane.

Traffic is typically steered into a Spray SR Policy in two ways:

- o local policy-based routing at the headend of the policy.
- o remote classification and steering via the BSID of the Spray SR Policy.

15. Reporting SR Policy

Stateful PCEP ([[I-D.ietf-pce-stateful-pce](#)] and [[I-D.sivabalan-pce-binding-label-sid](#)]) provides an ability for head-end to report BSID, attributes, and operational/administrative states. Using this protocol, a PCE can also update an existing SR Policy whose path computation is delegated to it as well as instantiate new SR Policy on a head-end.

BGP-LS reports an SR Policy via ([[I-D.ietf-idr-te-lsp-distribution](#)])

16. Work in Progress

- o Open configuration model.
- o Yang model.

17. Acknowledgement

18. Normative References

[GLOBECOM]

Filsfils, C., Nainar, N., Pignataro, C., Cardona, J., and P. Francois, "The Segment Routing Architecture, IEEE Global Communications Conference (GLOBECOM)", 2015.

[I-D.ietf-idr-te-lsp-distribution]

Previdi, S., Dong, J., Chen, M., Gredler, H., and j. jeffrant@gmail.com, "Distribution of Traffic Engineering (TE) Policies and State using BGP-LS", [draft-ietf-idr-te-lsp-distribution-07](#) (work in progress), July 2017.

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

Previdi, S., Filsfils, C., Bashandy, A., Gredler, H., Litkowski, S., Decraene, B., and j. jeffrant@gmail.com, "IS-IS Extensions for Segment Routing", [draft-ietf-isis-segment-routing-extensions-13](#) (work in progress), June 2017.

[I-D.ietf-pce-pce-initiated-lsp]

Crabbe, E., Minei, I., Sivabalan, S., and R. Varga, "PCEP Extensions for PCE-initiated LSP Setup in a Stateful PCE Model", [draft-ietf-pce-pce-initiated-lsp-11](#) (work in progress), October 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-10](#) (work in progress), October 2017.

[I-D.ietf-pce-stateful-pce]

Crabbe, E., Minei, I., Medved, J., and R. Varga, "PCEP Extensions for Stateful PCE", [draft-ietf-pce-stateful-pce-21](#) (work in progress), June 2017.

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

Filsfils, C., Previdi, S., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", [draft-ietf-spring-segment-routing-13](#) (work in progress), October 2017.

[I-D.previdi-idr-segment-routing-te-policy]

Previdi, S., Filsfils, C., Mattes, P., Rosen, E., and S. Lin, "Advertising Segment Routing Policies in BGP", [draft-previdi-idr-segment-routing-te-policy-07](#) (work in progress), June 2017.

[I-D.sivabalan-pce-binding-label-sid]

Sivabalan, S., Filsfils, C., Previdi, S., Tantsura, J., Hardwick, J., and D. Dhody, "Carrying Binding Label/Segment-ID in PCE-based Networks.", [draft-sivabalan-pce-binding-label-sid-03](#) (work in progress), July 2017.

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

[SIGCOMM] Hartert, R., Vissicchio, S., Schaus, P., Bonaventure, O., Filsfils, C., Telkamp, T., and P. Francois, "A Declarative and Expressive Approach to Control Forwarding Paths in Carrier-Grade Networks, ACM SIGCOMM", 2015.

Authors' Addresses

Clarence Filsfils
Cisco Systems, Inc.
Pegasus Parc
De kleetlaan 6a, DIEGEM BRABANT 1831
BELGIUM

Email: cfilsfil@cisco.com

Siva Sivabalan
Cisco Systems, Inc.
2000 Innovation Drive
Kanata, Ontario K2K 3E8
Canada

Email: msiva@cisco.com

Kamran Raza
Cisco Systems, Inc.
2000 Innovation Drive
Kanata, Ontario K2K 3E8
Canada

Email: skraza@cisco.com

Jose Liste
Cisco Systems, Inc.
821 Alder Drive
Milpitas, California 95035
USA

Email: jliste@cisco.com

Francois Clad
Cisco Systems, Inc.

Email: fclad@cisco.com

Shraddha Hegde
Juniper Networks, Inc.
Embassy Business Park
Bangalore, KA 560093
India

Email: shraddha@juniper.net

Daniel Voyer
Bell Canada.
671 de la gauchetiere W
Montreal, Quebec H3B 2M8
Canada

Email: daniel.voyer@bell.ca

Steven Lin
Google, Inc.

Email: stevenlin@google.com

Alex Bogdanov
Google, Inc.

Email: bogdanov@google.com

Martin Horneffer
Deutsche Telekom

Email: martin.horneffer@telekom.de

Dirk Steinberg
Steinberg Consulting

Email: dws@steinbergnet.net

Bruno Decraene
Orange Business Services

Email: bruno.decraene@orange.com

Stephane Litkowski
Orange Business Services

Email: stephane.litkowski@orange.com

