

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: 8 June 2023

A. Przygienda, Ed.
C. Bowers
Juniper
Y. Lee
Comcast
A. Sharma
Individual
R. White
Akamai
5 December 2022

IS-IS Flood Reflection
draft-ietf-lsr-isis-flood-reflection-12

Abstract

This document describes a backward-compatible, optional IS-IS extension that allows the creation of IS-IS flood reflection topologies. Flood reflection permits topologies in which L1 areas provide transit forwarding for L2 using all available L1 nodes internally. It accomplishes this by creating L2 flood reflection adjacencies within each L1 area. Those adjacencies are used to flood L2 LSPDUs and are used in the L2 SPF computation. However, they are not ordinarily utilized for forwarding within the flood reflection cluster. This arrangement gives the L2 topology significantly better scaling properties than traditionally used flat designs. As an additional benefit, only those routers directly participating in flood reflection are required to support the feature. This allows for incremental deployment of scalable L1 transit areas in an existing, previously flat network design, without the necessity of upgrading all routers in the network.

Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

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

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

This Internet-Draft will expire on 8 June 2023.

Copyright Notice

Copyright (c) 2022 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	3
2. Glossary	8
3. Further Details	9
4. Encodings	10
4.1. Flood Reflection TLV	10
4.2. Flood Reflection Discovery Sub-TLV	11
4.3. Flood Reflection Discovery Tunnel Type Sub-Sub-TLV	12
4.4. Flood Reflection Adjacency Sub-TLV	14
4.5. Flood Reflection Discovery	15
4.6. Flood Reflection Adjacency Formation	16
5. Route Computation	16
5.1. Tunnel-Based Deployment	17
5.2. No-Tunnel Deployment	17
6. Redistribution of Prefixes	17
7. Special Considerations	18
8. IANA Considerations	18
8.1. New IS-IS TLV Codepoint	19
8.2. Sub TLVs for IS-IS Router CAPABILITY TLV	19
8.3. Sub-sub TLVs for Flood Reflection Discovery sub-TLV	19
8.4. Sub TLVs for TLVs Advertising Neighbor Information	19

9. Security Considerations	20
10. Acknowledgements	20
11. References	20
11.1. Informative References	20
11.2. Normative References	21
Authors' Addresses	22

1. Introduction

This section introduces the problem space and outlines the solution. Some of the terms may be unfamiliar to readers without extensive IS-IS background; for such readers a glossary is provided in Section 2.

Due to the inherent properties of link-state protocols the number of IS-IS routers within a flooding domain is limited by processing and flooding overhead on each node. While that number can be maximized by well-written implementations and techniques such as exponential back-offs, IS-IS will still reach a saturation point where no further routers can be added to a single flooding domain. In some L2 backbone deployment scenarios, this limit presents a significant challenge.

The traditional approach to increasing the scale of an IS-IS deployment is to break it up into multiple L1 flooding domains and a single L2 backbone. This works well for designs where an L2 backbone connects L1 access topologies, but it is limiting where a single, flat L2 domain is supposed to span large number of routers. In such scenarios, an alternative approach could be to consider multiple L2 flooding domains connected together via L1 flooding domains. In other words, L2 flooding domains are connected by "L1/L2 lanes" through the L1 areas to form a single L2 backbone again. Unfortunately, in its simplest implementation, this requires the inclusion of most, or all, of the transit L1 routers as L1/L2 to allow traffic to flow along optimal paths through those transit areas. Consequently, such an approach fails to reduce the number of L2 routers involved and with that fails to increase the scalability of the L2 backbone as well.

Figure 1 is an example of a network where a topologically rich L1 area is used to provide transit between six different L2-only routers (R1-R6). Note that the six L2-only routers do not have connectivity to one another over L2 links. To take advantage of the abundance of paths in the L1 transit area, all the intermediate systems could be placed into both L1 and L2, but this essentially combines the separate L2 flooding domains into a single one, triggering again the maximum L2 scale limitation we try to address in first place.

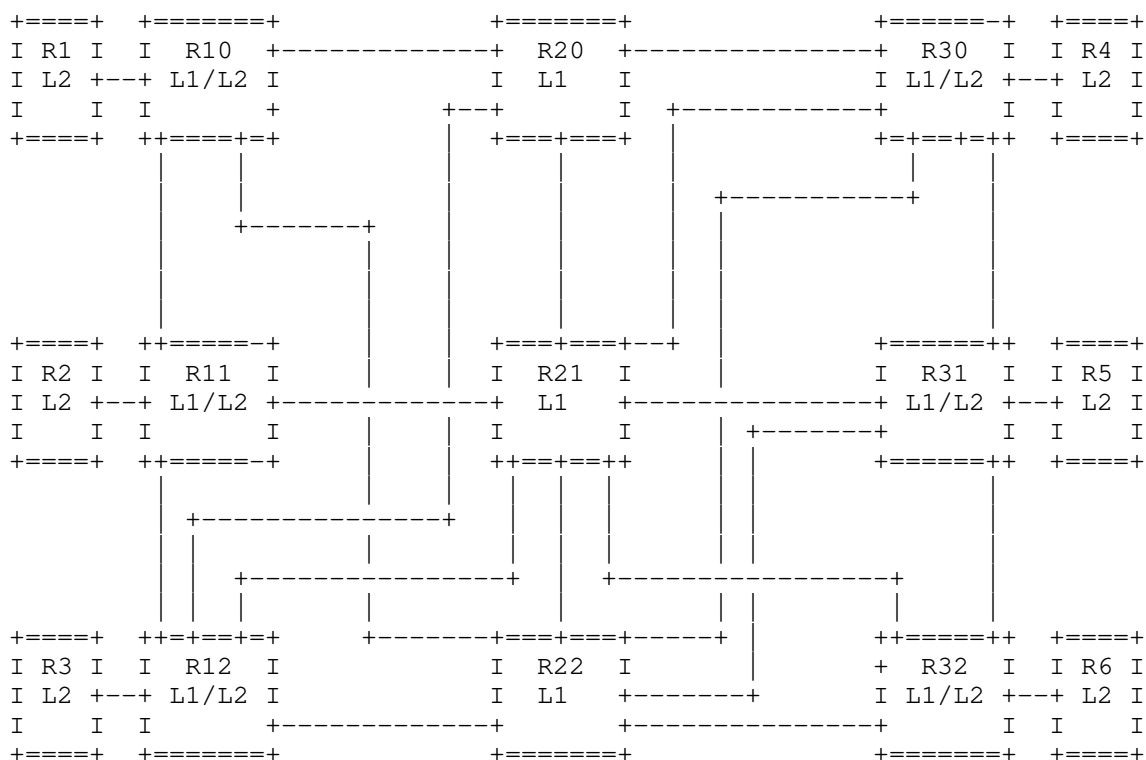


Figure 1: Example Topology of L1 with L2 Borders

A more effective solution would allow to reduce the number of links and routers exposed in L2, while still utilizing the full L1 topology when forwarding through the network.

[RFC8099] describes Topology Transparent Zones (TTZ) for OSPF. The TTZ mechanism represents a group of OSPF routers as a full mesh of adjacencies between the routers at the edge of the group. A similar mechanism could be applied to IS-IS as well. However, a full mesh of adjacencies between edge routers (or L1/L2 nodes) significantly limits the practically achievable scale of the resulting topology. The topology in Figure 1 has 6 L1/L2 nodes. Figure 2 illustrates a full mesh of L2 adjacencies between the 6 L1/L2 nodes, resulting in $(5 * 6)/2 = 15$ L2 adjacencies. In a somewhat larger topology containing 20 L1/L2 nodes, the number of L2 adjacencies in a full mesh rises to 190.

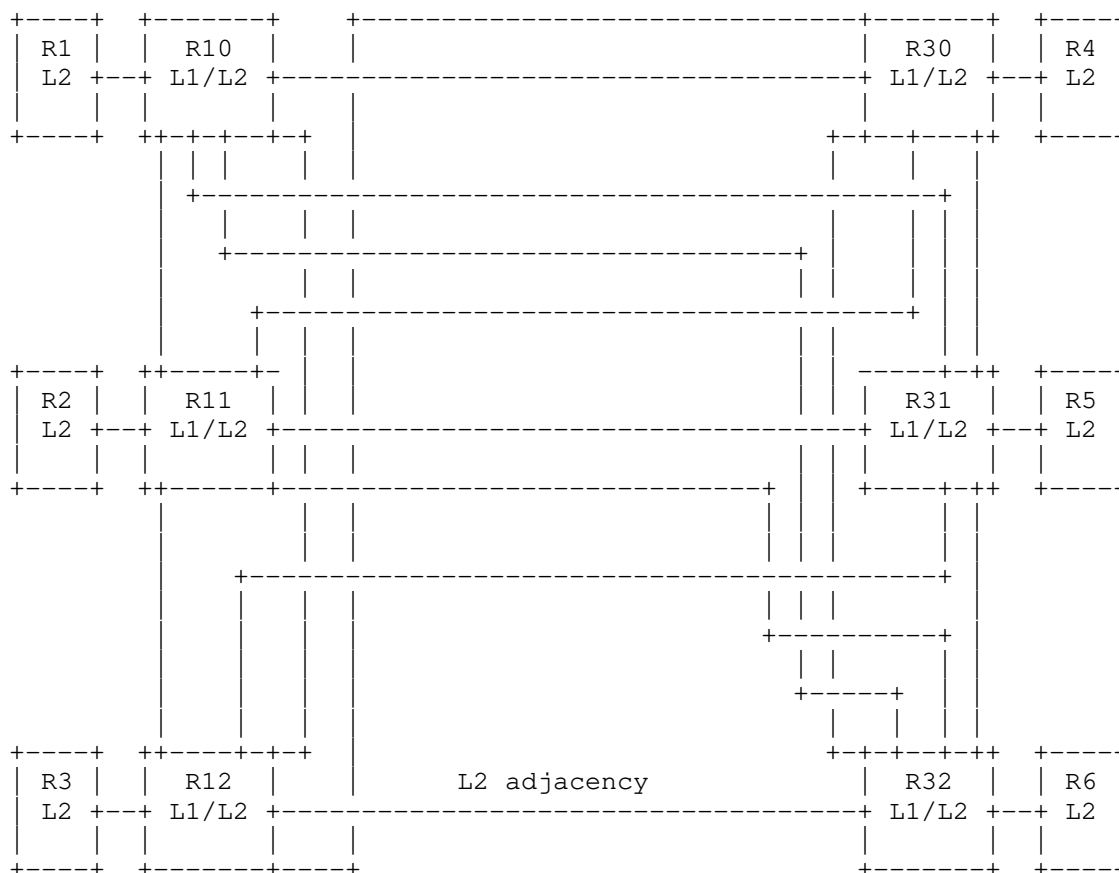


Figure 2: Example topology represented in L2 with a full mesh of L2 adjacencies between L1/L2 nodes

BGP, as specified in [RFC4271], faced a similar scaling problem, which has been solved in many networks by deploying BGP route reflectors [RFC4456]. We note that BGP route reflectors do not necessarily have to be in the forwarding path of the traffic. This non-congruity of forwarding and control path for BGP route reflectors allows the control plane to scale independently of the forwarding plane and represents an interesting degree of freedom in network architecture.

We propose in this document a similar solution for IS-IS and call it "flood reflection" in fashion analogous to "route reflection". A simple example of what a flood reflector control plane approach would look like is shown in Figure 3, where router R21 plays the role of a

flood reflector. Each L1/L2 ingress/egress router builds a tunnel to the flood reflector, and an L2 adjacency is built over each tunnel. In this solution, we need only 6 L2 adjacencies, instead of the 15 needed for a full mesh. In a somewhat larger topology containing 20 L1/L2 nodes, this solution requires only 20 L2 adjacencies, instead of the 190 needed for a full mesh. Multiple flood reflectors can be used, allowing the network operator to balance between resilience, path utilization, and state in the control plane. The resulting L2 adjacency scale is $R \cdot n$, where R is the number of flood reflectors used and n is the number of L1/L2 nodes. This compares quite favorably with $n \cdot (n-1) / 2$ L2 adjacencies required in a topologically fully meshed L2 solution.

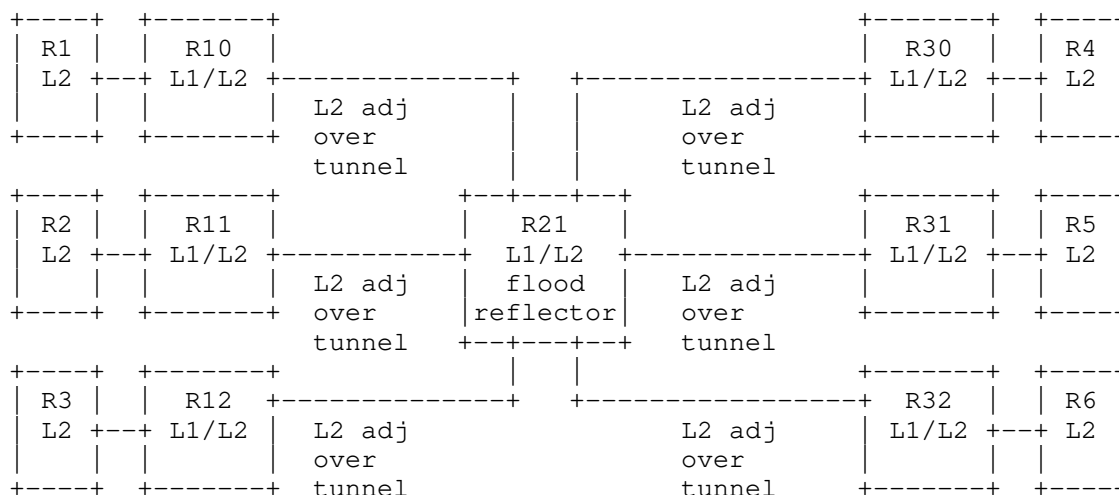


Figure 3: Example topology represented in L2 with L2 adjacencies from each L1/ L2 node to a single flood reflector

As illustrated in Figure 3, when R21 plays the role of flood reflector, it provides L2 connectivity among all of the previously disconnected L2 islands by reflooding all L2 LSPDUs. At the same time, R20 and R22 in Figure 1 remain L1-only routers. L1-only routers and L1-only links are not visible in L2. In this manner, the flood reflector allows us provide L2 control plane connectivity in a manner more scalable than a flat L2 domain.

As described so far, the solution illustrated in Figure 3 relies only on currently standardized IS-IS functionality. Without new functionality, however, the data traffic will traverse only R21. This will unnecessarily create a bottleneck at R21 since there is still available capacity in the paths crossing the L1-only routers R20 and R22 in Figure 1.

Hence, additional functionality is compulsory to allow the L1/L2 edge nodes (R10-12 and R30-32 in Figure 3) to recognize that the L2 adjacency to R21 should not be used for forwarding. The L1/L2 edge nodes should forward traffic that would normally be forwarded over the L2 adjacency to R21 over L1 links instead. This would allow the forwarding within the L1 area to use the L1-only nodes and links shown in Figure 1 as well. It allows networks to be built that use the entire forwarding capacity of the L1 areas, while at the same time introducing control plane scaling benefits provided by L2 flood reflectors.

It is expected that deployment at scale, and suitable time in operation, will provide sufficient evidence to either make this extension a standard, or suggest necessary modifications to accomplish this.

The remainder of this document defines the remaining extensions necessary for a complete flood reflection solution:

- * It defines a special 'flood reflector adjacency' built for the purpose of reflecting flooding information. These adjacencies allow 'flood reflectors' to participate in the IS-IS control plane without necessarily being used in the forwarding plane. Maintenance of such adjacencies is a purely local operation on the L1/L2 ingress and flood reflectors; it does not require replacing or modifying any routers not involved in the reflection process. In practical deployments, it is far less tricky to just upgrade the routers involved in flood reflection rather than have a flag day for the whole IS-IS domain.
- * It specifies an (optional) full mesh of tunnels between the L1/L2 ingress routers, ideally load-balancing across all available L1 links. This harnesses all forwarding paths between the L1/L2 edge nodes without injecting unneeded state into the L2 flooding domain or creating 'choke points' at the 'flood reflectors' themselves. The specification is agnostic as to the tunneling technology used but provides enough information for automatic establishment of such tunnels if desired. The discussion of IS-IS adjacency formation and/or liveness discovery on such tunnels is outside the scope of this specification and is largely a choice of the underlying implementation. A solution without tunnels is also

possible by introducing the correct scoping of reachability information between the levels. This is described in more detail later.

- * Finally, the document defines support of reflector redundancy and an (optional) way to auto-discover and annotate flood reflector adjacencies on advertisements. Such additional information in link advertisements allows L2 nodes outside the L1 area to recognize a flood reflection cluster and its adjacencies.

2. Glossary

The following terms are used in this document.

ISIS Level-1 and Level-2 areas, mostly abbreviated as L1 and L2: Traditional ISIS concepts whereas a routing domain has two "levels" with a single L2 area being the "backbone" that connects multiple L1 areas for scaling and reliability purposes. In traditional ISIS L2 can be used as transit for L1-L1 traffic but L1 areas cannot be used for that purpose since L2 level must be "connected" and all traffic flows along L2 routers until it arrives at the destination L1 area.

Flood Reflector:

Node configured to connect in L2 only to flood reflector clients and reflect (reflood) IS-IS L2 LSPs amongst them.

Flood Reflector Client:

Node configured to build Flood Reflector Adjacencies to Flood Reflectors, and normal adjacencies to other clients and L2 nodes not participating in flood reflection.

Flood Reflector Adjacency:

IS-IS L2 adjacency where one end is a Flood Reflector Client and the other a Flood Reflector, and the two have the same Flood Reflector Cluster ID.

Flood Reflector Cluster:

Collection of clients and flood reflectors configured with the same cluster identifier.

Tunnel-Based Deployment:

Deployment where Flood Reflector Clients build a partial or full mesh of tunnels in L1 to "shortcut" forwarding of L2 traffic through the cluster.

No-Tunnel Deployment:

Deployment where Flood Reflector Clients redistribute L2 reachability into L1 to allow forwarding through the cluster without underlying tunnels.

Tunnel Endpoint:

An endpoint that allows the establishment of a bi-directional tunnel carrying IS-IS control traffic or alternately, serves as the origin of such a tunnel.

L1 shortcut:

A tunnel between two clients visible in L1 only that is used as a next-hop, i.e. to carry data traffic in tunnel-based deployment mode.

Hot-Potato Routing:

In context of this document, a routing paradigm where L2->L1 routes are less preferred than L2 routes [RFC5302].

3. Further Details

Several considerations should be noted in relation to such a flood reflection mechanism.

First, this allows multi-area IS-IS deployments to scale without any major modifications in the IS-IS implementation on most of the nodes deployed in the network. Unmodified (traditional) L2 routers will compute reachability across the transit L1 area using the flood reflector adjacencies.

Second, the flood reflectors are not required to participate in forwarding traffic through the L1 transit area. These flood reflectors can be hosted on virtual devices outside the forwarding topology.

Third, astute readers will realize that flooding reflection may cause the use of suboptimal paths. This is similar to the BGP route reflection suboptimal routing problem described in [ID.draft-ietf-idr-bgp-optimal-route-reflection-28]. The L2 computation determines the egress L1/L2 and with that can create illusions of ECMP where there is none, and in certain scenarios lead to an L1/L2 egress which is not globally optimal. This represents a straightforward instance of the trade-off between the amount of control plane state and the optimal use of paths through the network often encountered when aggregating routing information.

One possible solution to this problem is to expose additional topology information into the L2 flooding domains. In the example network given, links from router R10 to router R11 can be exposed into L2 even when R10 and R11 are participating in flood reflection. This information would allow the L2 nodes to build 'shortcuts' when the L2 flood reflected part of the topology looks more expensive to cross distance wise.

Another possible variation is for an implementation to approximate with the tunnel cost the cost of the underlying topology.

Redundancy can be achieved by configuring multiple flood reflectors in a L1 area. Multiple flood reflectors do not need any synchronization mechanisms amongst themselves, except standard IS-IS flooding and database maintenance procedures.

4. Encodings

4.1. Flood Reflection TLV

The Flood Reflection TLV is a top-level TLV that MUST appear in L2 IIHs on all Flood Reflection Adjacencies. The Flood Reflection TLV indicates the flood reflector cluster (based on Flood Reflection Cluster ID) that a given router is configured to participate in. It also indicates whether the router is configured to play the role of either flood reflector or flood reflector client. The Flood Reflection Cluster ID and flood reflector roles advertised in the IIHs are used to ensure that flood reflector adjacencies are only formed between a flood reflector and flood reflector client, and that the Flood Reflection Cluster IDs match. The Flood Reflection 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      | C |  RESERVED  |
+-----+-----+-----+-----+-----+-----+-----+-----+
|          Flood Reflection Cluster ID          |
+-----+-----+-----+-----+-----+-----+-----+-----+
|      Sub-TLVs ...      |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Type: 161

Length: The length, in octets, of the following fields.

C (Client): This bit is set to indicate that the router acts as a

flood reflector client. When this bit is NOT set, the router acts as a flood reflector. On a given router, the same value of the C-bit MUST be advertised across all interfaces advertising the Flood Reflection TLV in IIHs.

RESERVED: This field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

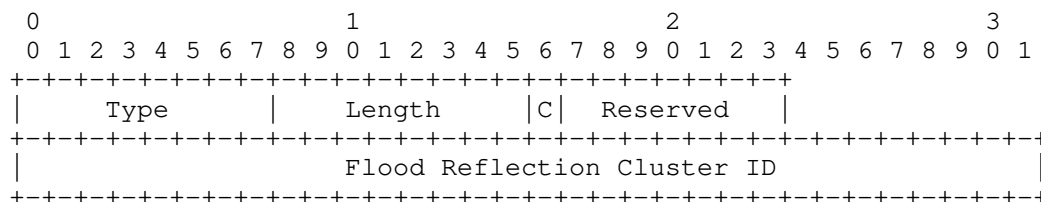
Flood Reflection Cluster ID: Flood Reflection Cluster Identifier. The same arbitrary 32-bit value MUST be assigned to all of the flood reflectors and flood reflector clients in the same L1 area. The value MUST be unique across different L1 areas within the IGP domain. In case of violation of those rules multiple L1 areas may become a single cluster or a single area may split in flood reflection sense and several mechanisms such as auto-discovery of tunnels may not work correctly. On a given router, the same value of the Flood Reflection Cluster ID MUST be advertised across all interfaces advertising the Flood Reflection TLV in IIHs. When a router discovers that a node is using more than a single Cluster IDs based on its advertised TLVs and IIHs, the node MAY log such violations subject to rate limiting. This implies that a flood reflector MUST NOT participate in more than a single L1 area. In case of Cluster ID value of 0, the TLV containing it MUST be ignored.

Sub-TLVs: Optional sub-TLVs. For future extensibility, the format of the Flood Reflection TLV allows for the possibility of including optional sub-TLVs. No sub-TLVs of the Flood Reflection TLV are defined in this document.

The Flood Reflection TLV SHOULD NOT appear more than once in an IIH. A router receiving one or more Flood Reflection TLVs in the same IIH MUST use the values in the first TLV and it SHOULD log such violations subject to rate limiting.

4.2. Flood Reflection Discovery Sub-TLV

The Flood Reflection Discovery sub-TLV is advertised as a sub-TLV of the IS-IS Router Capability TLV-242, defined in [RFC7981]. The Flood Reflection Discovery sub-TLV is advertised in L1 and L2 LSPs with area flooding scope in order to enable the auto-discovery of flood reflection capabilities. The Flood Reflection Discovery sub-TLV has the following format:



Type: 161

Length: The length, in octets, of the following fields.

C (Client): This bit is set to indicate that the router acts as a flood reflector client. When this bit is NOT set, the router acts as a flood reflector.

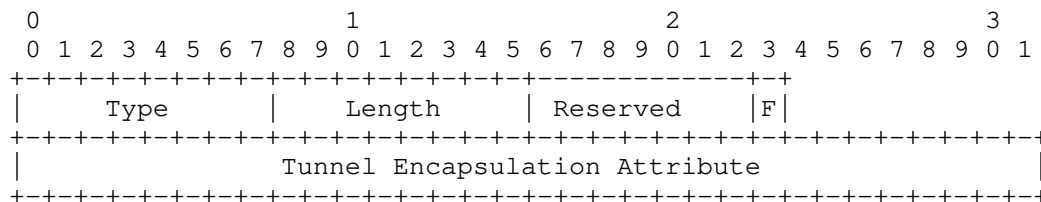
RESERVED: This field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

Flood Reflection Cluster ID: The Flood Reflection Cluster Identifier is the same as that defined in the Flood Reflection TLV and obeys the same rules.

The Flood Reflection Discovery sub-TLV SHOULD NOT appear more than once in TLV 242. A router receiving one or more Flood Reflection Discovery sub-TLVs in TLV 242 MUST use the values in the first sub-TLV of the lowest numbered fragment and it SHOULD log such violations subject to rate limiting.

4.3. Flood Reflection Discovery Tunnel Type Sub-Sub-TLV

Flood Reflection Discovery Tunnel Type sub-sub-TLV is advertised optionally as a sub-sub-TLV of the Flood Reflection Discovery Sub-TLV, defined in Section 4.2. It allows the automatic creation of L2 tunnels to be used as flood reflector adjacencies and L1 shortcut tunnels. The Flood Reflection Tunnel Type sub-sub-TLV has the following format:



Type: 161

Length: The length, in octets, of zero or more of the following fields.

Reserved: SHOULD be 0 on transmission and MUST be ignored on reception.

F Flag: When set indicates flood reflection tunnel endpoint, when clear, indicates possible L1 shortcut tunnel endpoint.

Tunnel Encapsulation Attribute: Carries encapsulation type and further attributes necessary for tunnel establishment as defined in [RFC9012]. In context of this attribute the protocol Type sub-TLV as defined in [RFC9012] MAY be included to ensure proper encapsulation of IS-IS frames. In case such a sub-TLV is included and the F flag is set (i.e. the resulting tunnel is a flood reflector adjacency) this sub-TLV MUST include a type that allows to carry encapsulated IS-IS frames. Furthermore, such tunnel type MUST be able to transport IS-IS frames of size up to 'originatingL2LSPBufferSize'.

A flood reflector receiving Flood Reflection Discovery Tunnel Type sub-sub-TLVs in Flood Reflection Discovery sub-TLV with F flag set (i.e. the resulting tunnel is a flood reflector adjacency) SHOULD use one or more of the specified tunnel endpoints to automatically establish one or more tunnels that will serve as flood reflection adjacency(-ies) to the clients advertising the endpoints.

A flood reflection client receiving one or more Flood Reflection Discovery Tunnel Type sub-sub-TLVs in Flood Reflection Discovery sub-TLV with F flag clear (i.e. the resulting tunnel is used to support tunnel-based mode) from other leaves MAY use one or more of the specified tunnel endpoints to automatically establish one or more tunnels that will serve as L1 tunnel shortcuts to the clients advertising the endpoints.

In case of automatic flood reflection adjacency tunnels and in case IS-IS adjacencies are being formed across L1 shortcuts all the aforementioned rules in Section 4.5 apply as well.

Optional address validation procedures as defined in [RFC9012] MUST be disregarded.

It remains to be observed that automatic tunnel discovery is an optional part of the specification and can be replaced or mixed with statically configured tunnels for either flood reflector adjacencies and/or tunnel-based shortcuts. Specific implementation details how both mechanisms interact are specific to an implementation and mode of operation and are outside the scope of this document.

Flood reflector adjacencies rely on IS-IS L2 liveness procedures. In case of L1 shortcuts the mechanism used to ensure liveness and tunnel integrity are outside the scope of this document.

4.4. Flood Reflection Adjacency Sub-TLV

The Flood Reflection Adjacency sub-TLV is advertised as a sub-TLV of TLVs 22, 23, 25, 141, 222, and 223 (the "TLVs Advertising Neighbor Information"). Its presence indicates that a given adjacency is a flood reflector adjacency. It is included in L2 area scope flooded LSPs. The Flood Reflection Adjacency 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	2	3	4	5	6	7	8	9	0	1																						
Type																Length																C	Reserved																														
Flood Reflection Cluster ID																																																															

Type: 161

Length: The length, in octets, of the following fields.

C (Client): This bit is set to indicate that the router advertising this adjacency is a flood reflector client. When this bit is NOT set, the router advertising this adjacency is a flood reflector.

RESERVED: This field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

Flood Reflection Cluster ID: The Flood Reflection Cluster Identifier is the same as that defined in the Flood Reflection TLV and obeys the same rules.

The Flood Reflection Adjacency sub-TLV SHOULD NOT appear more than once in a given TLV. A router receiving one or more Flood Reflection Adjacency sub-TLVs in a TLV MUST use the values in the first sub-TLV of the lowest numbered fragment and it SHOULD log such violations subject to rate limiting.

4.5. Flood Reflection Discovery

A router participating in flood reflection as client or reflector MUST be configured as an L1/L2 router. It MAY originate the Flood Reflection Discovery sub-TLV with area flooding scope in L1 and L2. Normally, all routers on the edge of the L1 area (those having traditional L2 adjacencies) will advertise themselves as flood reflector clients. Therefore, a flood reflector client will have both traditional L2 adjacencies and flood reflector L2 adjacencies.

A router acting as a flood reflector MUST NOT form any traditional L2 adjacencies except with flood reflector clients. It will be an L1/L2 router only by virtue of having flood reflector L2 adjacencies. A router desiring to act as a flood reflector MAY advertise itself as such using the Flood Reflection Discovery sub-TLV in L1 and L2.

A given flood reflector or flood reflector client can only participate in a single cluster, as determined by the value of its Flood Reflection Cluster ID and should disregard other routers' TLVs for flood reflection purposes if the cluster ID is not matching.

Upon reception of Flood Reflection Discovery sub-TLVs, a router acting as flood reflector SHOULD initiate a tunnel towards each flood reflector client with which it shares a Flood Reflection Cluster ID using one or more of the tunnel encapsulations provided with F flag is set. The L2 adjacencies formed over such tunnels MUST be marked as flood reflector adjacencies. If the client or reflector has a direct L2 adjacency with the according remote side it SHOULD use it instead of instantiating a tunnel.

In case the optional auto-discovery mechanism is not implemented or enabled a deployment MAY use statically configured tunnels to create flood reflection adjacencies.

The IS-IS metrics for all flood reflection adjacencies in a cluster SHOULD be identical.

Upon reception of Flood Reflection Discovery TLVs, a router acting as a flood reflector client MAY initiate tunnels with L1-only adjacencies towards any of the other flood reflector clients with lower router IDs in its cluster using encapsulations with F flag clear. These tunnels MAY be used for forwarding to improve the load-balancing characteristics of the L1 area. If the clients have a direct L2 adjacency they SHOULD use it instead of instantiating a new tunnel.

4.6. Flood Reflection Adjacency Formation

In order to simplify implementation complexity, this document does not allow the formation of complex hierarchies of flood reflectors and clients or allow multiple clusters in a single L1 area. Consequently, all flood reflectors and flood reflector clients in the same L1 area MUST share the same Flood Reflector Cluster ID. Deployment of multiple cluster IDs in the same L1 area are outside the scope of this document.

A flood reflector MUST NOT form flood reflection adjacencies with flood reflector clients with a different Cluster ID. A flood reflector MUST NOT form any traditional L2 adjacencies.

Flood reflector clients MUST NOT form flood reflection adjacencies with flood reflectors with a different Cluster ID.

Flood reflector clients MAY form traditional L2 adjacencies with flood reflector clients or nodes not participating in flood reflection. When two flood reflector clients form a traditional L2 adjacency the Cluster IDs are disregarded.

The Flood Reflector Cluster ID and flood reflector roles advertised in the Flood Reflection TLVs in IIHs are used to ensure that flood reflection adjacencies that are established meet the above criteria.

On change in either flood reflection role or cluster ID on IIH on the local or remote side the adjacency has to be reset. It is then re-established if possible.

Once a flood reflection adjacency is established, the flood reflector and the flood reflector client MUST advertise the adjacency by including the Flood Reflection Adjacency Sub-TLV in the Extended IS reachability TLV or MT-ISN TLV.

5. Route Computation

To ensure loop-free routing, the flood reflection client MUST follow the normal L2 computation to determine L2 routes. This is because nodes outside the L1 area will generally not be aware that flood reflection is being performed. The flood reflection clients need to produce the same result for the L2 route computation as a router not participating in flood reflection.

5.1. Tunnel-Based Deployment

In the tunnel-based option the reflection client, after L2 and L1 computation, MUST examine all L2 routes with flood reflector next-hop adjacencies. Such next-hops must be replaced by the corresponding tunnel next-hops to the correct egress nodes of the flood reflection cluster.

5.2. No-Tunnel Deployment

In case of deployment without underlying tunnels, the necessary L2 routes are distributed into the area, normally as L2->L1 routes. Due to the rules in Section 4.6 the computation in the resulting topology is relatively simple, the L2 SPF from a flood reflector client is guaranteed to reach the Flood Reflector within a single hop, and in the following hop the L2 egress to which it has a forwarding tunnel. All the flood reflector tunnel nexthops in the according L2 route can hence be removed and if the L2 route has no other ECMP L2 nexthops, the L2 route MUST be suppressed in the RIB by some means to allow the less preferred L2->L1 route to be used to forward traffic towards the advertising egress.

In the particular case the client has L2 routes which are not flood reflected, those will be naturally preferred (such routes normally "hot-potato" packets out of the L1 area). However in the case the L2 route through the flood reflector egress is "shorter" than such present non flood reflected L2 routes, the node SHOULD ensure that such routes are suppressed so the L2->L1 towards the egress still takes preference. Observe that operationally this can be resolved in a relatively simple way by configuring flood reflector adjacencies to have a high metric, i.e. the flood reflector topology becomes "last resort" and the leaves will try to "hot-potato" out the area as fast as possible which is normally the desirable behavior.

In No-tunnel deployment all L1/L2 edge nodes MUST be flood reflection clients.

6. Redistribution of Prefixes

In case of no-tunnel deployment per Section 5.2 a client that does not have any L2 flood reflector adjacencies MUST NOT redistribute L2 routes into the cluster.

The L2 prefix advertisements redistributed into an L1 that contains flood reflectors SHOULD be normally limited to L2 intra-area routes (as defined in [RFC7775]), if the information exists to distinguish them from other L2 prefix advertisements.

On the other hand, in topologies that make use of flood reflection to hide the structure of L1 areas while still providing transit forwarding across them using tunnels, we generally do not need to redistribute L1 prefix advertisements into L2.

7. Special Considerations

In pathological cases setting the overload bit in L1 (but not in L2) can partition L1 forwarding, while allowing L2 reachability through flood reflector adjacencies to exist. In such a case a node cannot replace a route through a flood reflector adjacency with a L1 shortcut and the client MAY use the L2 tunnel to the flood reflector for forwarding but in any case it MUST initiate an alarm and declare misconfiguration.

A flood reflector with directly L2 attached prefixes should advertise those in L1 as well since based on preference of L1 routes the clients will not try to use the L2 flood reflector adjacency to route the packet towards them. A very unlikely corner case can occur when the flood reflector is reachable via L2 flood reflector adjacency (due to underlying L1 partition) exclusively, in which case the client can use the L2 tunnel to the flood reflector for forwarding towards those prefixes while it MUST initiate an alarm and declare misconfiguration.

A flood reflector MUST NOT set the attached bit on its LSPs.

In certain cases where reflectors are attached to same broadcast medium, and where some other L2 router, which is neither a flood reflector nor a flood reflector client (a "non-FR router"), attaches to the same broadcast medium, flooding between the reflectors in question might not succeed, potentially partitioning the flood reflection domain. This could happen specifically in the event that the non-FR router is chosen as the designated intermediate system ("DIS", the designated router). Since, per Section 4.6, a flood reflector MUST NOT form an adjacency with a non-FR router, the flood reflector(s) will not be represented in the pseudo-node.

To avoid this situation, it is RECOMMENDED that flood reflectors not be deployed on the same broadcast medium as non-FR routers.

A router discovering such condition MUST initiate an alarm and declare misconfiguration.

8. IANA Considerations

This document requests allocation for the following IS-IS TLVs and Sub-TLVs, and requests creation of a new registry.

8.1. New IS-IS TLV Codepoint

This document requests the following IS-IS TLV under the IS-IS Top-Level TLV Codepoints registry::

Value	Name	IIH	LSP	SNP	Purge
161	Flood Reflection	y	n	n	n

8.2. Sub TLVs for IS-IS Router CAPABILITY TLV

This document request the following registration in the "sub-TLVs for IS-IS Router CAPABILITY TLV" registry.

Type	Description
161	Flood Reflection Discovery

8.3. Sub-sub TLVs for Flood Reflection Discovery sub-TLV

This document requests creation of a new registry named "Sub-sub TLVs for Flood Reflection Discovery sub-TLV" under the "IS-IS TLV Codepoints" grouping. The Registration Procedures for this registry are Expert Review, following the common expert review guidance given for the grouping.

The range of values in this registry is 0-255. The registry should be seeded with the following initial registration:

Type	Description
161	Flood Reflection Discovery Tunnel Encapsulation Attribute

8.4. Sub TLVs for TLVs Advertising Neighbor Information

This document requests the following registration in the "IS-IS Sub-TLVs for TLVs Advertising Neighbor Information" registry.

Type	Description	22	23	25	141	222	223
161	Flood Reflector Adjacency	y	y	n	y	y	y

9. Security Considerations

This document uses flood reflection tunnels to carry IS-IS control traffic. If an attacker can inject traffic into such a tunnel, the consequences could be in the most extreme case the complete subversion of the IS-IS level 2 information. Therefore, a mechanism inherent to the tunnel technology should be taken to prevent such injection. Since the available security procedures will vary by deployment and tunnel type, the details of securing tunnels are beyond the scope of this document.

This document specifies information used to form dynamically discovered shortcut tunnels. If an attacker were able to hijack the endpoint of such a tunnel and form an adjacency, it could divert short-cut traffic to itself, placing itself on-path and facilitating on-path attacks or could even completely subvert the IS-IS level 2 routing. Therefore, steps should be taken to prevent such capture by using mechanism inherent to the tunnel type used. Since the available security procedures will vary by deployment and tunnel type, the details of securing tunnels are beyond the scope of this document.

Additionally, the usual IS-IS security mechanisms [RFC5304] SHOULD be deployed to prevent misrepresentation of routing information by an attacker in case a tunnel is compromised if the tunnel itself does not provide mechanisms strong enough guaranteeing the integrity of the messages exchanged.

10. Acknowledgements

The authors thank Shraddha Hegde, Peter Psenak, Acee Lindem, Robert Raszuk and Les Ginsberg for their thorough review and detailed discussions. Thanks are also extended to Michael Richardson for an excellent routing directorate review. John Scudder ultimately spent significant time helping to make the document more comprehensible and coherent.

11. References

11.1. Informative References

[ID.draft-ietf-idr-bgp-optimal-route-reflection-28]
Raszuk et al., R., "BGP Optimal Route Reflection", July 2019, <<https://www.ietf.org/id/draft-ietf-idr-bgp-optimal-route-reflection-28.txt>>.

- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", RFC 4456, DOI 10.17487/RFC4456, April 2006, <<https://www.rfc-editor.org/info/rfc4456>>.
- [RFC8099] Chen, H., Li, R., Retana, A., Yang, Y., and Z. Liu, "OSPF Topology-Transparent Zone", RFC 8099, DOI 10.17487/RFC8099, February 2017, <<https://www.rfc-editor.org/info/rfc8099>>.

11.2. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC5302] Li, T., Smit, H., and T. Przygienda, "Domain-Wide Prefix Distribution with Two-Level IS-IS", RFC 5302, DOI 10.17487/RFC5302, October 2008, <<https://www.rfc-editor.org/info/rfc5302>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<https://www.rfc-editor.org/info/rfc5304>>.
- [RFC7775] Ginsberg, L., Litkowski, S., and S. Previdi, "IS-IS Route Preference for Extended IP and IPv6 Reachability", RFC 7775, DOI 10.17487/RFC7775, February 2016, <<https://www.rfc-editor.org/info/rfc7775>>.
- [RFC7981] Ginsberg, L., Previdi, S., and M. Chen, "IS-IS Extensions for Advertising Router Information", RFC 7981, DOI 10.17487/RFC7981, October 2016, <<https://www.rfc-editor.org/info/rfc7981>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

[RFC9012] Patel, K., Van de Velde, G., Sangli, S., and J. Scudder,
"The BGP Tunnel Encapsulation Attribute", RFC 9012,
DOI 10.17487/RFC9012, April 2021,
<<https://www.rfc-editor.org/info/rfc9012>>.

Authors' Addresses

Tony Przygienda (editor)
Juniper
1137 Innovation Way
Sunnyvale, CA
United States of America
Email: prz@juniper.net

Chris Bowers
Juniper
1137 Innovation Way
Sunnyvale, CA
United States of America
Email: cbowers@juniper.net

Yiu Lee
Comcast
1800 Bishops Gate Blvd
Mount Laurel, NJ 08054
United States of America
Email: Yiu_Lee@comcast.com

Alankar Sharma
Individual
Email: as3957@gmail.com

Russ White
Akamai
Email: russ@riw.us