

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: March 12, 2020

A. Przygienda
Juniper
Y. Lee
A. Sharma
Comcast
R. White
Juniper
September 9, 2019

Flood Reflectors
draft-przygienda-flood-reflector-00

Abstract

This document provides specification of an optional ISIS extension that allows to create L2 flood reflector topologies independent of resulting forwarding within L1 areas when they are used as 'transit' to guarantee L2 connectivity between L2 "islands".

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](#)].

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 March 12, 2020.

Copyright Notice

Copyright (c) 2019 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.	Description	2
2.	Further Details	6
3.	Flood Reflection TLV	7
4.	Non-Forwarding Adjacency Sub-TLV	7
5.	Procedures	8
6.	Adjacency Forming Procedures	9
7.	Special Considerations	9
8.	IANA Considerations	10
9.	Security Considerations	10
10.	Acknowledgements	10
11.	References	10
11.1.	Informative References	10
11.2.	Normative References	11
	Authors' Addresses	11

[1.](#) Description

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 certain deployment scenarios of L2 backbones, this limit presents an obstacle.

While the standard solution to increase the scale of an IS-IS deployment is to break it up into multiple L1 flooding domains and a single L2 backbone, and alternative way is to think about "multiple" L2 flooding domains connected via L1 flooding domains. In such a solution, the L2 flooding domains are connected by "L1/L2 lanes" through the L1 areas to form a single L2 backbone again. However, in the 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 such transit areas and with that

ultimately does not help to reduce number of L2 routers and increase the scalability of L2 backbone.

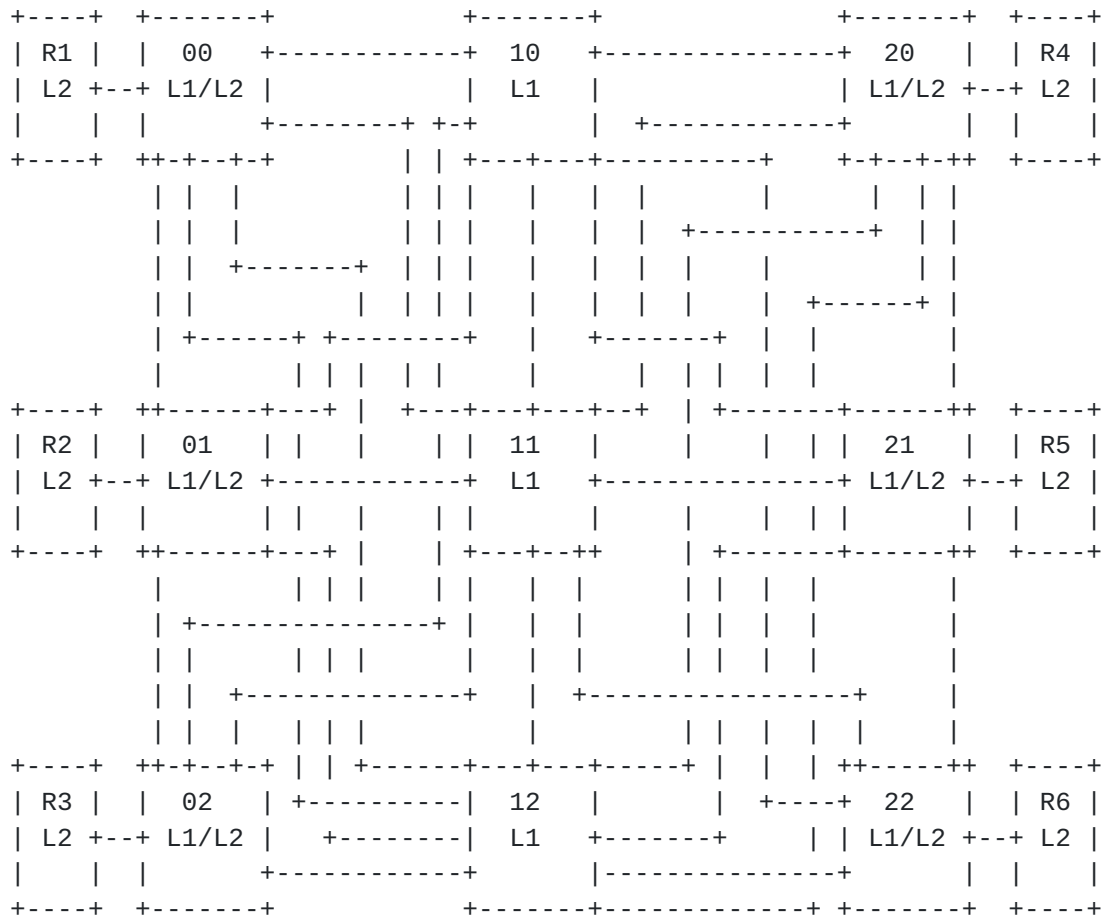


Figure 1

Figure 1 is an example of a network where a topologically rich L1 area is used to provide transit between six different routers in L2 "partitions" (R1-R6). To take advantage of the cornucopia of paths in the L1 transit, 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 maximum L2 scale limitations again.

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.

The mechanism described in [\[RFC8099\]](#) could be used in ISIS to build a full mesh of tunnels over the L1 transit, but a full mesh of tunnels

can also quickly limit the scaling. The network in Figure 2 would expose 6 L1/L2 nodes and $(5 * 6)/2 = 15$ L2 tunnels. In a slightly larger network, however, in a comparable topology containing 15 L1/L2 edge nodes the number grows very quickly to 105 tunnels.

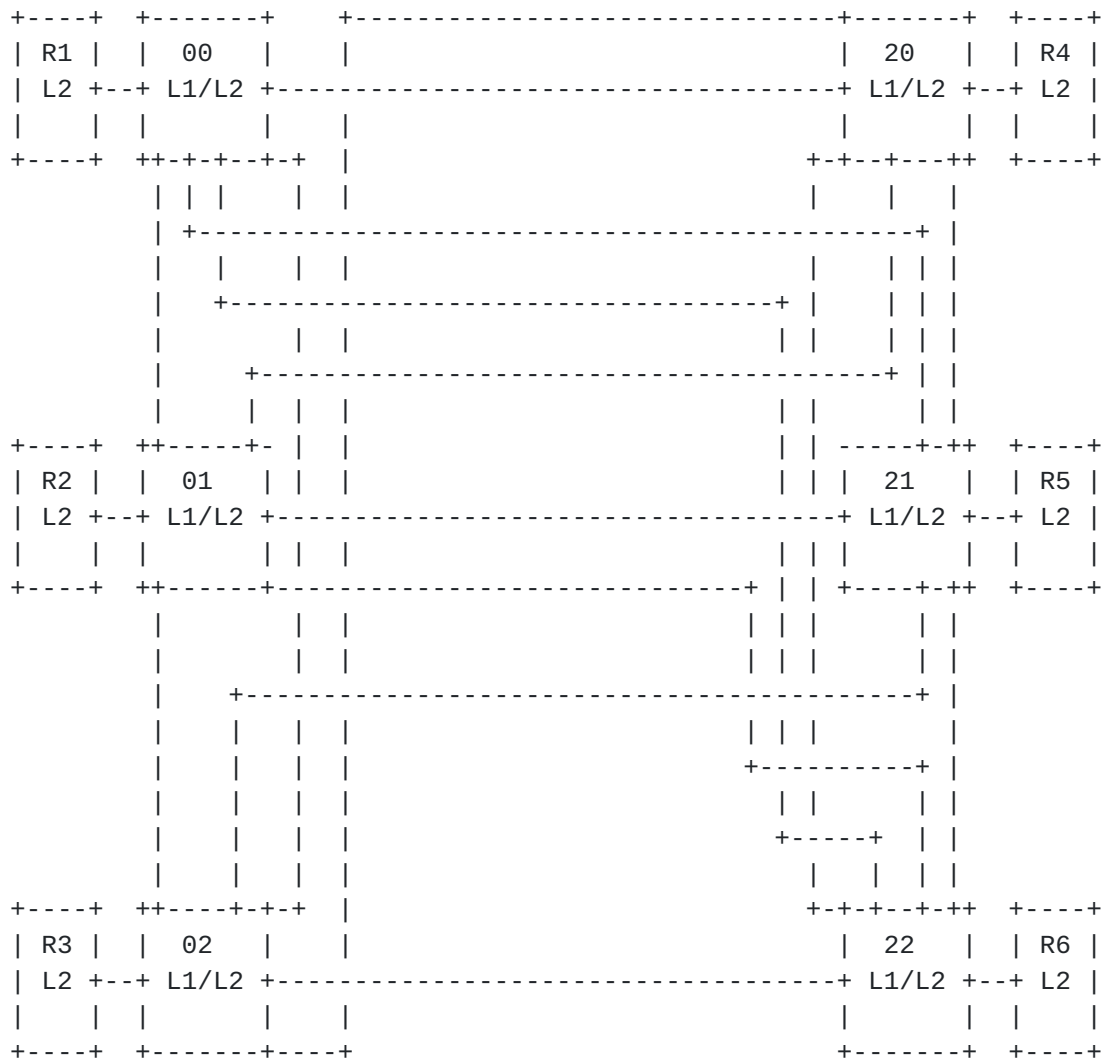


Figure 2

BGP, described in [[RFC4271](#)], faced a similar scaling problem, which has been solved in many networks by deploying BGP route reflectors, as described in [[RFC4456](#)]. And, to offer another crucial observation, BGP route reflectors do not necessarily need to be in the forwarding path.

We suggest here a similar solution for IS-IS. A good approximation of what a "flood reflector" approach would look like is shown in

Figure 3, where router 11 is used as 'reflector.' All L1/L2 routers build an L2 tunnel to such reflectors, so we end up with only 6 L2 tunnels instead of 15 of a full mesh. Multiple such reflectors can be used, of course, allowing the network operator to balance between resilience, path utilization, and state in the control plane. The resulting L2 tunnel scale is roughly $R * n$ where R is the redundancy factor or in other words, number of flood reflectors used. This compares quite favorably with $n^2 / 2$ tunnels used in a fully meshed L2 solution.

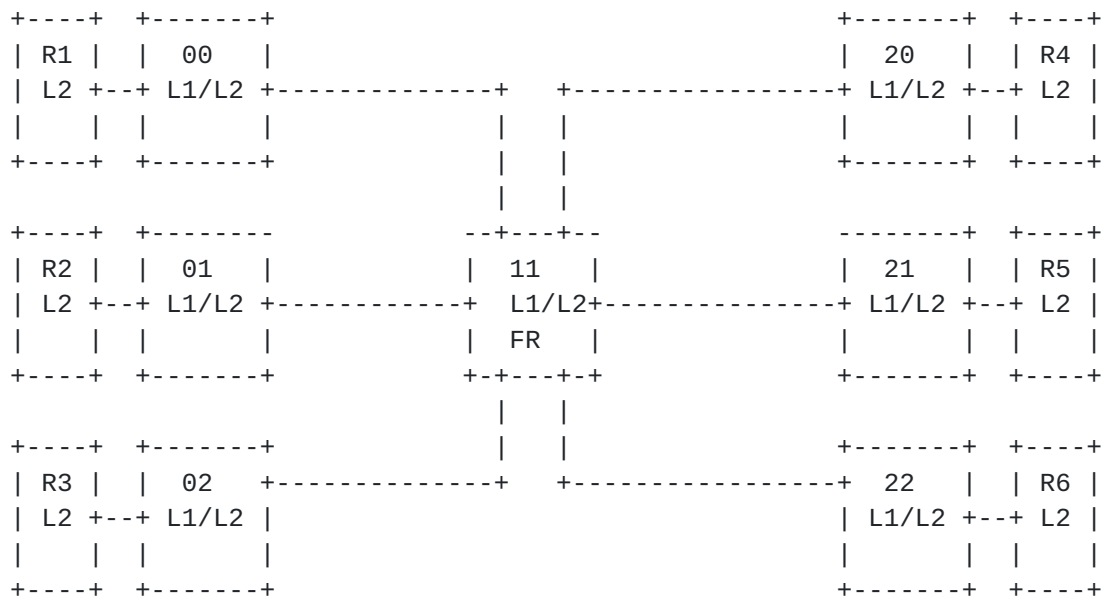


Figure 3

This proposal, however, without further qualification would concentrate forwarded traffic at router 11. It would be hence desirable to decouple the forwarding plane from the control plane, so router 11 can refllood information without being placed in the forwarding path (hence router 11 would not end up being a forwarding plane bottleneck). To achieve that goal, multiple pieces will be necessary, only one of which is a local protocol extension on the L1/L2 leafs and the 'flood reflectors'. In first approximation these extensions include:

- o A full mesh of L1 tunnels between the L1/L2 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.'

- o A 'non-forwarding adjacency' for all the adjacencies built for the purpose of reflecting flooding information. This allows these 'flood reflectors' to participate in the IS-IS control plane without being used in the forwarding plane. This is a purely local operation on the L1/L2 ingress; it does not require replacing or modifying any routers not involved in the reflection process.
- o Some system to support reflector redundancy, and potentially some way to auto-discover and advertise such adjacencies as non-forwarding. This may allow L2 nodes outside the L1 to perform optimizations in the future based on this information.

2. 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 non-forwarding adjacencies.

Second, the flooding reflectors are not required to participate in forwarding traffic through the L1 transit area. These flooding 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-19](#)]. 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 01 to router 02 can be exposed into L2 even when 01 and 02 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 L1 tunnel cost the cost of the underlying topology.

Redundancy in the solution is trivial to achieve by building multiple flood reflectors into the L1 area while all reflectors are still remaining completely stateless and do not need any kind of synchronized algorithms amongst themselves except standard ISIS flooding procedures and database.

3. Flood Reflection TLV

The Flood Reflection TLV is indicating the participation of a node as reflector and/or client. It is included in L1 area scope flooded LSPs and on L1 and L2 IIH.

```

      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       |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Priority        | FR Cluster ID |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Type: TBD

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

Reflector Priority Priority of the router to act as flood reflector in the cluster. A value of 0 indicates that the router is a client in the cluster. Any value higher than 0 indicates preference to be a flood reflector. Higher values are to be preferred by clients.

FR Cluster ID Flood Reflector Cluster Identifier to allow a node to participate in possibly multiple clusters.

4. Non-Forwarding Adjacency Sub-TLV


```

      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      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| FR Cluster ID |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Type: TBD

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

FR Cluster ID Flood Reflector Cluster Identifier to which this NFA belongs.

5. Procedures

There are a number of points to consider when implementing and deploying this solution, including:

A router participating in flood reflection **MUST** be configured as L1L2 router. It originates the Flood Reflection TLV with area flooding scope in L1 only. Normally routers on the edge of the area, i.e. with non-FR L2 adjacencies, will advertise themselves as clients. Any L1L2 non-client router in the area can act as FR.

A flood reflector can participate in a single cluster only, the clients are free to participate in multiple clusters at the same time.

Upon reception of a Flood Reflection TLV router acting as client (in case it doesn't have such L2 adjacencies already) **MUST** initialize tunnels towards all the FRs with highest priority and **MAY** initiate such tunnels to FRs with lower priority. L2 over such tunnels **MUST** be marked as non-forwarding adjacencies. If the client has a direct L2 adjacency with the flood reflector it **SHOULD** use it instead of instantiating a tunnel.

Upon reception of a Flood Reflection TLV router acting as client in case it doesn't have such direct L1 adjacencies already **SHOULD** initialize tunnels towards all the other clients in the its clusters. L1 **only** adjacencies **SHOULD** be built over such tunnels to ensure their liveliness, but other means can be used (since those adjacencies are used for L1 forwarding, it is prudent to advertise them into L1 as forwarding links).

On the reflection client, after L2 and L1 computation, all non-forwarding adjacencies used as next-hops for L2 routes MUST be examined and replaced with the correct L1 tunnel next-hop to the egress. Due to the rules in [Section 6](#) the computation in the resulting topology is relatively simple, the L2 SPF from a flood reflector client is guaranteed to reach within a hop the FR and in the following hop the L2 egress to which it has a L1 forwarding tunnel. However, if the topology has L2 paths which are not route reflected and look "shorter" than path through the FR then the computation will have to track the egress out of the L1 domain by a more advanced algorithm.

A node, when advertising the L2 NFA SHOULD include the Non-Forwarding Adjacency Sub-TLV in Extended IS reachability TLV and MT-ISN TLV.

6. Adjacency Forming Procedures

To ensure loop-free routing the ingress routers MUST follow normal L2 computation to generate L2 routes. This is because nodes outside the L1 area may not be aware that flooding reflection is performed. The resulting short cuts through the L1 area needs to be able to easily calculate the egress L1/L2 router where the tunnel tail-end is located.

To prevent complex scenarios of flood reflectors building L2 adjacencies within a cluster or across clusters or hierarchies of reflectors, a flood reflector MUST never form an L2 adjacency with a peer if the peer is not a client in the same Cluster ID. This ensures a L2 computation on an ingress link or adjacency following a non-forwarding adjacency will always traverse a client of the flood reflector to exit the flooding domain. This allows shortcuts through the L1 area to be used without any danger of forwarding loops.

Depending on pseudo-node choice in case of a broadcast domain with multiple flood reflectors attached this can lead to a partitioned LAN and hence a router discovering such a condition MUST initiate an alarm and declare misconfiguration.

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 non-forwarding adjacencies to exist. In such a case a node cannot replace a route through non-forwarding adjacency with a L1 shortcut and the client can use the L2 tunnel to the flood reflector for forwarding while 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 non-forwarding adjacency to route the packet towards them. A very, very corner case is when the flood reflector is reachable via L2 non-forwarding adjacency (due to underlying L1 partition) only 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.

Instead of modifying the computation procedures one could imagine a flood reflector solution where the FR would re-advertise the L2 prefixes with a 'third-party' next-hop but that would have less desirable convergence properties than the solution proposed and force a fork-lift of all L2 routers to make sure they disregard such prefixes unless in the same L1 domain as the FR.

8. IANA Considerations

This document will request IANA to assign new TLV type value in the ISIS TLV Codepoints registry.

This document will request IANA to assign new TLV type value in the 'Sub-TLVs for TLVs 22, 23, 25, 141, 222, and 223 (Extended IS reachability, IS Neighbor Attribute, L2 Bundle Member Attributes, inter-AS reachability information, MT-ISN, and MT IS Neighbor Attribute TLVs)' registry.

9. Security Considerations

This document introduces no new security concerns to ISIS or other specifications referenced in this document.

10. Acknowledgements

Thanks to Shraddha and Chris Bowers for thorough review.

11. References

11.1. Informative References

- [ID.[draft-ietf-idr-bgp-optimal-route-reflection-19](#)]
Raszuk et al., R., "BGP Optimal Route Reflection", July 2019.
- [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>>.

Authors' Addresses

Tony Przygienda
Juniper
1137 Innovation Way
Sunnyvale, CA
USA

Email: prz_at_juniper.net

Yiu Lee
Comcast
1800 Bishops Gate Blvd
Mount Laurel, NJ 08054
US

Email: Yiu_Lee_at_comcast.com

Alankar Sharma
Comcast
1800 Bishops Gate Blvd
Mount Laurel, NJ 08054
US

Email: Alankar_Sharma_at_comcast.com

Russ White
Juniper
1137 Innovation Way
Sunnyvale, CA
USA

Email: russw_at_juniper.net