

Layer 2 Virtual Private Networks
Internet-Draft
Intended status: Informational
Expires: January 17, 2013

O. Dornon
J. Kotalwar
Alcatel-Lucent
J. Zhang
Juniper Networks, Inc.
V. Hemige
Alcatel-Lucent
July 16, 2012

PIM Snooping over VPLS
draft-ietf-l2vpn-vpls-pim-snooping-02

Abstract

In Virtual Private LAN Service (VPLS), as also in IEEE Bridged Networks, the switches simply flood multicast traffic on all ports in the LAN by default. IGMP Snooping is commonly deployed to ensure multicast traffic is not forwarded on ports without IGMP receivers. The procedures and recommendations for IGMP Snooping are defined in [IGMP-SNOOP]. But when any protocol other than IGMP is used, the common practice is to simply flood multicast traffic to all ports. PIM-SM, PIM-SSM, PIM-BIDIR are widely deployed routing protocols. PIM Snooping procedures are important to restrict multicast traffic to only the routers interested in receiving such traffic.

While most of the PIM Snooping procedures defined here also apply to IEEE Bridged Networks, VPLS demands certain special procedures due to the split-horizon rules that require the Provider Edge (PE) devices to co-operate. This document describes the procedures and recommendations for PIM-Snooping in VPLS to facilitate replication to only those ports behind which there are interested PIM routers and/or IGMP hosts. This document also describes procedures for PIM Proxy. PIM Proxy is required on PEs for VPLS Multicast to work correctly when Join suppression is enabled in the VPLS. PIM Proxy also helps scale VPLS Multicast much better than just PIM Snooping.

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 <http://datatracker.ietf.org/drafts/current/>.

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

This Internet-Draft will expire on January 17, 2013.

Copyright Notice

Copyright (c) 2012 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	5
1.1. Assumptions	6
1.2. PIM Snooping and PIM Proxy Complexity	6
1.3. Definitions	6
2. Multicast Traffic over VPLS	7
2.1. Constraining of IP Multicast in a VPLS	8
2.2. IPv6 Considerations	9
2.3. PIM-SM (*,*,RP) Considerations	9
2.4. PIM Packet Types to Snoop	9
2.5. PIM Snooping vs PIM Proxy	9
2.5.1. Differences between PIM Snooping and PIM Proxy	10
2.5.2. PIM Control Message Latency	11
2.5.3. When to Snoop and When to Proxy	11
3. PIM Snooping for VPLS	12
3.1. General Rules for PIM Snooping in VPLS	13
3.1.1. Snooping PIM Packets	13
3.1.2. Preserving Assert Trigger	13
3.2. Discovering PIM Routers	14
3.3. PIM-SM and PIM-SSM	15
3.3.1. Building PIM-SM Snooping States	15
3.3.2. Explanation for per (S,G,N) states	17
3.3.3. Receiving (*,G) PIM-SM Join/Prune Messages	18
3.3.4. Receiving (S,G) PIM-SM Join/Prune Messages	20
3.3.5. Receiving (S,G,rpt) Join/Prune Messages	22
3.3.6. Sending Join/Prune Messages Upstream	22
3.4. Bidirectional-PIM (PIM-BIDIR)	23
3.5. Interaction with IGMP Snooping	24
3.6. PIM-DM	24
3.6.1. Building PIM-DM Snooping States	24
3.6.2. PIM-DM Downstream Per-Port PIM(S,G,N) State Machine	25
3.6.3. Triggering ASSERT election in PIM-DM	25
3.7. PIM Proxy	25
3.7.1. Downstream PIM Proxy behavior	26
3.7.2. Upstream PIM Proxy behavior	26
3.7.3. Source IP Address in Proxy PIM Join/Prune Packets	26
3.8. Directly Connected Multicast Source	27
3.9. Data Forwarding Rules	27
3.9.1. PIM-SM Data Forwarding Rules	28
3.9.2. PIM-BIDIR Data Forwarding Rules	29
3.9.3. PIM-DM Data Forwarding Rules	30
4. IANA Considerations	31
5. Security Considerations	31
6. Contributors	31
7. Acknowledgements	32
8. References	32
8.1. Normative References	32

8.2. Informative References	32
Appendix A. PIM-BIDIR Thoughts	33
Appendix B. Example Network Scenario	33
B.1. Pim Snooping Example	34
B.2. PIM Proxy Example with (S,G) / (*,G) interaction	36
Authors' Addresses	39

1. Introduction

In Virtual Private LAN Service (VPLS), the Provider Edge (PE) devices provide a logical interconnect such that Customer Edge (CE) devices belonging to a specific VPLS instance appear to be connected by a single LAN. Forwarding information base for particular VPLS instance is populated dynamically by source MAC address learning. This is a straightforward solution to support unicast traffic, with reasonable flooding for unicast unknown traffic. Since a VPLS provides LAN emulation for IEEE bridges as well as for routers, the unicast and multicast traffic need to follow the same path for layer-2 protocols to work properly. As such, multicast traffic is treated as broadcast traffic and is flooded to every site in the VPLS instance. VPLS solutions (i.e., [VPLS-LDP] and [VPLS-BGP]) perform replication for multicast traffic at the ingress PE devices. As stated in the VPLS Multicast Requirements draft [VPLS-MCAST-REQ], there are two issues with VPLS Multicast today: A. Multicast traffic is replicated to non-member sites. B. Replication of PWS on shared physical path.

This document solves Issue A of [VPLS-MCAST-REQ] by ensuring that IP multicast traffic is not forwarded to non-member sites. Issue B is outside the scope of this document. The different mechanisms to tunnel IP multicast traffic in a VPLS from the ingress PE to the egress PEs are discussed in [VPLS-MCAST-TREES]. The solution in this document when combined with the solutions proposed in the working group to solve Issue B will provide a complete VPLS Multicast solution set.

Using IGMP/PIM Snooping in VPLS has the following advantages:

- o It improves IP Multicast bandwidth usage in the VPLS core by ensuring traffic is replicated only to PEs with member sites. Note that this is not necessarily optimum, as there can still be bandwidth waste if traffic from a PE to other PE(s) is not forwarded along a minimum cost spanning tree.
- o It prevents sending multicast traffic to non-member sites.

Procedures for IGMP Snooping are specified in [IGMP-SNOOP]. This document describes the procedures for Protocol Independent Multicast (PIM) snooping over VPLS for efficient distribution of IP multicast traffic. It also describes the rules when both IGMP and PIM are active in a VPLS instance.

This document also describes procedures for PIM Proxy. PIM Proxy is required on PEs for VPLS Multicast to work correctly when Join suppression is enabled in the VPLS. PIM Proxy also helps scale VPLS Multicast much better than just PIM Snooping.

1.1. Assumptions

Since this draft describes the procedures for PIM Snooping and PIM Proxy, the draft assumes that the reader has a good understanding of the PIM protocols. The text in this draft is written in the same style as the PIM RFCs to help correlate the concepts and to make it easier to follow. In order to avoid replicating the text relating to PIM protocol handling here, this draft assumes that the user will infer such detail from the PIM RFC referenced in this document. Deviations in protocol handling specific to PIM Snooping and PIM Proxy are specified in this draft. There could be cross references into definitions of macros and procedures from the PIM RFCs.

1.2. PIM Snooping and PIM Proxy Complexity

The PIM Snooping and PIM Proxy solutions described here requires a switch to examine and operate on only PIM Hello and PIM Join/Prune packets. The switch does not need to examine any other PIM packets.

The switch does not need to have any routing tables like is required in PIM Multicast Routing. It knows how to forward Join/Prunes by looking at the Upstream Neighbor field in the Join/Prune packets.

The switch does not need to know about Rendezvous Points (RP) and does not have to maintain any RP Set. All that is transparent to a PIM Snooping switch.

Most of the procedures in PIM Snooping and PIM Proxy in the handling of PIM Hellos and PIM Join/Prune packets are very similar to that of a PIM Router.

The solutions described here provide complete separation of control and data planes.

A PIM Proxy solution minimizes the control plane messages received at CE routers by proxying one message upstream on behalf of a large number of downstream CEs. As such control plane messaging is very similar to that of a PIM Router.

1.3. Definitions

There are several definitions referenced in this document that are well described in the PIM RFCs [PIM-SM], [PIM-BIDIR], [PIM-DM]. The following definitions and abbreviations are used throughout this document:

- o A port is defined as either an attachment circuit (AC) or a Pseudo-Wire (PW).
- o When we say a PIM message is 'received' on a port, it means that a PIM Snooping switch snooped the PIM message.

Abbreviations used in the document:

- o S: IP Address of the Multicast Source.
- o G: IP Address of the Multicast Group.
- o N: Upstream Neighbor field in a Join/Prune/Graft message.
- o Rport(N): Port on which neighbor N is learnt

Other definitions are explained in the sections where they are introduced.

2. Multicast Traffic over VPLS

In VPLS, if a PE receives a frame from an Attachment Circuit (AC) with no matching entry in the forwarding information base for that particular VPLS instance, it floods the frame to all other PEs (which are part of this VPLS instance) and to directly connected ACs (other than the one that the frame is received from). The flooding of a frame occurs when:

- o The destination MAC address has not been learned,
- o The destination MAC address is a broadcast address,
- o The destination MAC address is a multicast address.

Malicious attacks (e.g., receiving unknown frames constantly) aside, the first situation is handled by VPLS solutions as long as destination MAC address can be learned. After that point on, the frames will not be flooded. A PE is REQUIRED to have safeguards, such as unknown unicast limiting and MAC table limiting, against malicious unknown unicast attacks.

There is no way around flooding broadcast frames. To prevent runaway broadcast traffic from adversely affecting the VPLS service and the SP network, a PE is REQUIRED to have tools to rate limit the broadcast traffic as well.

Similar to broadcast frames, multicast frames are flooded as well, as

a PE cannot know where multicast members reside. Rate limiting multicast traffic, while possible, should be done carefully since several network control protocols relies on multicast. For one thing, layer-2 and layer-3 protocols utilize multicast for their operation. For instance, Bridge Protocol Data Units (BPDUs) use an IEEE assigned all bridges multicast MAC address, and OSPF is multicast to all OSPF routers multicast MAC address. If the rate-limiting of multicast traffic is not done properly, the customer network will experience instability and poor performance. For the other, it is not straightforward to determine the right rate limiting parameters for multicast.

A VPLS solution MUST NOT affect the operation of customer layer-2 protocols (e.g., BPDUs). Additionally, a VPLS solution MUST NOT affect the operation of layer-3 protocols.

In the following section, we describe procedures to constrain the flooding of IP multicast traffic in a VPLS.

2.1. Constraining of IP Multicast in a VPLS

For a PE in a VPLS (a layer-2 device) to constrain IP multicast traffic, it needs to be able to learn which CEs are interested in receiving multicast traffic for what flows.

The most obvious solution is to snoop IP multicast control traffic at the PEs. Snooping as a solution to constrain multicast traffic makes sense under the following circumstances:

- o The CE-CE protocol the PEs snoop is a popular and widely deployed protocol.
- o It does not require any changes on the CEs and it should be completely transparent to the CEs.

IGMP/MLD and PIM are the popular IP Multicast Routing protocols today. Other routing protocols such as DVMRP or MOSPF are outside the scope of this document.

This document describes the guidelines for PIM Snooping and PIM Proxy in VPLS. The specifications in this document could be used for either PIM Snooping or PIM Proxy. The PIM Proxy solution is described in section Section 3.7. Differences that need to be observed while implementing one or the other and recommendations on which method to employ in different scenarios are noted in section Section 2.5. We will largely refer to PIM "Snooping" in this document. Unless specifically specified, the same procedures should apply to a Proxy solution as well.

In the following sub-sections, we provide some guidelines for the implementation of PIM snooping in VPLS. Snooping techniques need to be employed on ACs at the downstream PEs. Snooping techniques can also be employed on PWs at the upstream PEs. This may work well for small to medium scale deployments. However, if there are a large number of VPLS instances with a large number of PEs per instances, then the amount of snooping required at the upstream PEs can overwhelm the upstream PEs.

2.2. IPv6 Considerations

In VPLS, PEs forward Ethernet frames received from CEs and as such are agnostic of the layer-3 protocol used by the CEs. However, as an IGMP and PIM snooping switch, the PE would have to look deeper into the IP and IGMP/PIM packets and build snooping state based on that. The PIM Protocol specifications handle both IPv4 and IPv6. The specification for PIM Snooping in this draft can be applied to both IPv4 and IPv6 payloads.

2.3. PIM-SM (*,*,RP) Considerations

This draft does not address (*,*,RP) states in the VPLS network. Although [PIM-SM] specifies that routers MUST support (*,*,RP) states, there are very few implementations that actually support it in actual deployments. Given the complexity of supporting (*,*,RP) states and knowing that there is little to no use to supporting it, this draft omits the specification relating to (*,*,RP) support.

2.4. PIM Packet Types to Snoop

A PIM Snooping switch need only snoop on PIM Hellos and PIM Join/Prune packets. All other PIM packets can be transparently flooded unexamined.

2.5. PIM Snooping vs PIM Proxy

PIM Snooping switches simply snoop on PIM packets as they are being forwarded in the VPLS. As such it truly provides transparent LAN services since no customer packets are modified or consumed or new packets introduced in the VPLS. It is also slightly simpler to implement than PIM Proxy. However for PIM Snooping to work correctly, it is a requirement that CE routers MUST disable Join suppression in the VPLS.

Given that a large number of existing CE deployments do not support disabling of Join suppression and given the operational complexity for a provider to manage disabling of Join suppression in the VPLS, it becomes a difficult solution to deploy. Another disadvantage of

PIM Snooping as a solution is that it does not scale as well as PIM Proxy. If there are a large number of CEs in a VPLS, then every CE will see every other CE's Join/Prune messages.

PIM Proxy on the PEs has the advantage that it does not require Join suppression to be disabled in the VPLS. Multicast as a VPLS service can be very easily be provided without requiring any changes on the CE routers. It also helps scale VPLS Multicast very well since the PEs intelligently forward only one Join/Prune message for a given flow and only to the upstream CE.

PIM Proxy as a solution however loses the transparency argument since Join/Prunes could get modified or even consumed at a PE. Also, new packets could get introduced in the VPLS. However, this loss of transparency is limited to PIM Join/Prune packets. It is in the interest of optimizing multicast in the VPLS and helping a VPLS network scale much better. Data traffic will still be completely transparent.

2.5.1. Differences between PIM Snooping and PIM Proxy

For PIM-SM and PIM-BIDIR, a PIM Snooping/Proxy Switch only needs to examine PIM Hello and Join/Prune messages. PIM Proxy for PIM-DM is for future study and is not currently specified in this draft.

A proxy switch performs proxy only for the Join/Prune messages. PIM Hello messages are snooped by both PIM Snooping and PIM Proxy switches.

Details on the PIM Proxy solution are discussed in section Section 3.7. This section is presented here to say that most of the procedures to follow (unless explicitly specified) are common to both PIM Snooping and PIM Proxy. Differences between a PIM Snooping switch and a PIM Proxy switch can be summarized as the following:

PIM Snooping	PIM Proxy
1. PIM Snooping switches snoop Hello and Join/Prune messages while they are transparently flooded in the VPLS.	1. PIM Proxy switches also snoop PIM Hello messages while they are transparently flooded in the VPLS. But they consume PIM Join/Prune messages and do not flood them as is in the VPLS.
2. PIM Snooping switches do not originate any PIM packets.	2. PIM Proxy switches may originate new or modified Join/Prune packets.

Other than the above simple differences, most of the procedures are common to PIM Snooping and PIM Proxy. In the text to follow, we describe the procedures for PIM "Snooping". Unless specifically stated otherwise, such procedures apply to PIM Proxy as well.

2.5.2. PIM Control Message Latency

A PIM Snooping or PIM Proxy switch snoops on PIM Hello packets while transparently flooding it in the VPLS. As such there is no latency introduced by the VPLS in the delivery of PIM Hello packets to remote CEs in the VPLS.

A PIM Proxy switch consumes PIM Join/Prune packets and generates proxy Join/Prune packets to be sent upstream. This can result in additional latency for a downstream CE to receive multicast traffic after it has sent a Join. When a downstream CE prunes a multicast stream, the traffic should stop flowing to the CE with no additional latency introduced by the VPLS.

A PIM Snooping switch snoops on PIM Join/Prune packets while transparently flooding them in the VPLS. There is no latency introduced by the VPLS in the delivery of PIM Join/Prune packets when PIM Snooping is employed.

2.5.3. When to Snoop and When to Proxy

Explicit Tracking (ET) is enabled in a VPLS when all PIM CE Routers in the VPLS advertise Tracking Support in their PIM Hello messages. If even one does not advertise Tracking Support, then all PIM CE routers disable ET in the VPLS. When ET is enabled, it implies that Join Suppression is disabled and vice versa.

PIM Snooping PEs can determine if ET is enabled or disabled in a VPLS by examining PIM Hellos. If ET is disabled, PIM Proxy MUST be used. If ET is enabled, either PIM Snooping or PIM Proxy can be used, unless the PIM control message latency due to proxy is a concern, in which case PIM Snooping SHOULD be used.

3. PIM Snooping for VPLS

IGMP snooping procedures described in [IGMP-SNOOP] provide efficient delivery of IP multicast traffic in a given VPLS service when end stations are connected to the VPLS. However, when VPLS is offered as a WAN service it is likely that the CE devices are routers and would run PIM between them. To provide efficient IP multicasting in such cases, it is necessary that the PE routers offering the VPLS service do PIM snooping.

PIM is a multicast routing protocol, which runs exclusively between routers. PIM shares many of the common characteristics of a routing protocol, such as discovery messages (e.g., neighbor discovery using Hello messages), topology information (e.g., multicast tree), and error detection and notification (e.g., dead timer and designated router election). On the other hand, PIM does not participate in any kind of exchange of databases, as it uses the unicast routing table to provide reverse path information for building multicast trees. There are a few variants of PIM. In [PIM-DM], multicast data is pushed towards the members similar to broadcast mechanism. PIM-DM constructs a separate delivery tree for each multicast group. As opposed to PIM-DM, other PIM flavors (PIM-SM [PIM-SM], PIM-SSM [PIM-SSM], and PIM-BIDIR [PIM-BIDIR]) invoke a pull methodology instead of push technique.

PIM routers periodically exchange Hello messages to discover and maintain stateful sessions with neighbors. After neighbors are discovered, PIM routers can signal their intentions to join or prune specific multicast groups. This is accomplished by having downstream routers send an explicit Join/Prune message (for the sake of generalization, consider Graft messages for PIM-DM as Join messages) to the upstream routers. The Join/Prune message can be group specific (*,G) or group and source specific (S,G).

In PIM snooping, a PE snoops on the PIM message exchanged between routers, and builds its multicast states.

Based on the multicast states, it forwards IP multicast traffic accordingly to avoid unnecessary flooding.

In the following sub-sections, snooping mechanisms for each variety

of PIM are specified.

3.1. General Rules for PIM Snooping in VPLS

The following rules for the correct operation of PIM snooping MUST be followed.

- o PIM messages and multicast data traffic forwarded by PEs MUST follow the split-horizon rule for mesh PWs.
- o PIM snooping states in a PE MUST be per VPLS instance.
- o PIM assert triggers MUST be preserved to the extent necessary to avoid sending duplicate traffic to the same PE (see Section 3.1.2).

3.1.1. Snooping PIM Packets

PIM-SM snooping PEs need to snoop on just the PIM Hello and PIM Join/Prune messages to build its multicast states.

- o PIM-DM snooping PEs have to also snoop on PIM Graft and PIM State Refresh messages.

3.1.2. Preserving Assert Trigger

In PIM-SM/DM, there are scenarios where multiple routers could be forwarding the same multicast traffic on a LAN. When this happens, using PIM Assert Election process by sending PIM Assert Messages, routers ensure that only the Assert Winner forwards traffic on the LAN. The Assert Election is a data driven event and happens only if a router sees traffic on the interface to which it should be forwarding the traffic. In the case of VPLS with snooping, two routers may forward the same flow at the same time but each copy may reach different set of PEs, and that is acceptable from the point of view of avoiding duplicate traffic. If the two copies may reach the same PE then the sending routers must be able to see each other's traffic, in order to trigger Assert Election and stop duplicate traffic.

To achieve that, PIM-SM Snooping MUST not only forward multicast traffic for an (S,G) on the ports on which they snooped Joins(S,G)/Joins(*,G), but also towards the upstream neighbor(s)). In other words, the ports on which the upstream neighbors are learnt must be added to the outgoing port list along with the ports on which Joins are snooped.

Similarly, PIM-DM Snooping SHOULD make sure that asserts can be

triggered (Section 3.6.3).

The above logic needs to be facilitated without breaking VPLS Split Horizon Rules. i.e. traffic should not be forwarded on the port on which it was received, and traffic arriving on a PW MUST NOT be forwarded onto other PW(s).

3.2. Discovering PIM Routers

A PIM Snooping PE MUST snoop on PIM Hellos received on ACs and PWs. i.e. the PE transparently floods the PIM Hello while snooping on it. PIM Hellos are used by the snooping switch to discover PIM routers and their characteristics.

For each neighbor discovered by a PE, it includes an entry in the PIM Neighbor Database with the following fields:

- o Layer 2 encapsulation for the Router sending the PIM Hello.
- o IP Address and address family of the Router sending the PIM Hello.
- o Port (AC / PW) on which the PIM Hello was received.
- o Hello TLVs

The PE should be able to interpret and act on Hello TLVs currently defined in the PIM RFCs. The TLVs of particular interest in this document are:

- o Hello-Hold-Time
- o Tracking Support
- o DR Priority

Please refer to [PIM-SM] for a list of the Hello TLVs. When a PIM Hello is received, the PE MUST reset the neighbor-expiry- timer to Hello-Hold-Time. If a PE does not receive a Hello message from a router within Hello-Hold-Time, the PE MUST remove that neighbor from its PIM Neighbor Database. If a PE receives a Hello message from a router with Hello-Hold-Time value set to zero, the PE MUST remove that router from the PIM snooping state immediately.

From the PIM Neighbor Database, a PE MUST be able to use the procedures defined in [PIM-SM] to identify the PIM Designated Router in the VPLS instance. It should also be able to determine if Tracking Support is active in the VPLS instance.

3.3. PIM-SM and PIM-SSM

The key characteristic of PIM-SM and PIM-SSM is explicit join behavior. In this model, multicast traffic is only forwarded to locations that specifically request it. The root node of a tree is the Rendezvous Point (RP) in case of a shared tree (PIM-SM only) or the first hop router that is directly connected to the multicast source in the case of a shortest path tree. All the procedures described in this section apply to both PIM-SM and PIM-SSM, except for the fact that there is no (*,G) state in PIM-SSM.

3.3.1. Building PIM-SM Snooping States

PIM-SM and PIM-SSM Snooping states are built by snooping on the PIM-SM Join/Prune messages received on AC/PWs.

The downstream state machine of a PIM-SM snooping switch very closely resembles the downstream state machine of PIM-SM routers. The downstream state consists of:

Per downstream (Port, *, G):

- o DownstreamJPState: One of { "NoInfo" (NI), "Join" (J), "Prune Pending" (PP) }

Per downstream (Port, *, G, N):

- o Prune Pending Timer (PPT(N))
- o Join Expiry Timer (ET(N))

Per downstream (Port, S, G):

- o DownstreamJPState: One of { "NoInfo" (NI), "Join" (J), "Prune Pending" (PP) }

Per downstream (Port, S, G, N):

- o Prune Pending Timer (PPT(N))
- o Join Expiry Timer (ET(N))

Per downstream (Port, S, G, rpt):

- o DownstreamJPRptState: One of { "NoInfo" (NI), "Pruned" (P), "Prune Pending" (PP) }

Per downstream (Port, S, G, rpt, N):

- o Prune Pending Timer (PPT(N))
- o Join Expiry Timer (ET(N))

Where S is the address of the multicast source, G is the Group address and N is the upstream neighbor field in the Join/Prune message. Notice that unlike on PIM-SM routers where PPT and ET are per (Interface, S, G), PIM Snooping switches have to maintain PPT and ET per (Port, S, G, N). The reasons for this are explained in section Section 3.3.2

Apart from the above states, we define the following state summarization macros.

UpstreamNeighbors(*,G): If there is one or more Join(*,G) received on any port with upstream neighbor N and ET(N) is active, then N is added to UpstreamNeighbors(*,G). This set is used to determine if a Join(*,G) or a Prune(*,G) with upstream neighbor N needs to be sent upstream.

UpstreamNeighbors(S,G): If there is one or more Join(S,G) received on any port with upstream neighbor N and ET(N) is active, then N is added to UpstreamNeighbors(S,G). This set is used to determine if a Join(S,G) or a Prune(S,G) with upstream neighbor N needs to be sent upstream.

UpstreamPorts(*,G): This is the set of all Rport(N) ports where N is in the set UpstreamNeighbors(*,G). Multicast Streams forwarded using a (*,G) match MUST be forwarded to these ports in addition to downstream ports. So UpstreamPorts(*,G) MUST be added to OutgoingPortList(*,G).

UpstreamPorts(S,G): This is the set of all Rport(N) ports where N is in the set UpstreamNeighbors(S,G). UpstreamPorts(S,G) MUST be added to OutgoingPortList(S,G).

InheritedUpstreamPorts(S,G): This is the union of UpstreamPorts(S,G) and UpstreamPorts(*,G).

UpstreamPorts(S,G,rpt): If PruneDesired(S,G,rpt) becomes true, then this set is set to UpstreamPorts(*,G). Otherwise, this set is empty. UpstreamPorts(*,G) (-) UpstreamPorts(S,G,rpt) MUST be added to OutgoingPortList(S,G).

UpstreamPorts(G): This set is the union of all the UpstreamPorts(S,G) and UpstreamPorts(*,G) for a given G. Proxy (S,G) Join/Prune and (*,G) Join/Prune messages MUST be sent to a subset of UpstreamPorts(G) as specified in section Section 3.3.6.1.

PWPorts: This is the set of all PWs.

OutgoingPortList(*,G): This is the set of all ports to which traffic needs to be forwarded on a (*,G) match.

OutgoingPortList(S,G): This is the set of all ports to which traffic needs to be forwarded on an (S,G) match.

See section Section 3.9 on Data Forwarding Rules for the specification on how OutgoingPortList is calculated.

NumETsActive(Port,*,G): Number of (Port,*,G,N) entries that have Expiry Timer running. This macro keeps track of the number of Join(*,G)s that are received on this Port with different upstream neighbors.

NumETsActive(Port,S,G): Number of (Port,S,G,N) entries that have Expiry Timer running. This macro keeps track of the number of Join(*,G)s that are received on this Port with different upstream neighbors.

RpfVectorTlvs(*,G): RPF Vectors [RPF-VECTOR] are TLVs that may be present in received Join(*,G) messages. If present, they must be copied to RpfVectorTlvs(*,G).

RpfVectorTlvs(S,G): RPF Vectors [RPF-VECTOR] are TLVs that may be present in received Join(S,G) messages. If present, they must be copied to RpfVectorTlvs(S,G).

Since there are a few differences between the downstream state machines of PIM-SM Routers and PIM-SM snooping switches, we specify the details of the downstream state machine of PIM-SM snooping switches at the risk of repeating most of the text documented in [PIM-SM].

3.3.2. Explanation for per (S,G,N) states

In PIM Routing protocols, states are built per (S,G). On a router, an (S,G) has only one RPF-Neighbor. However, a PIM Snooping switch does not have the Layer 3 routing information available to the routers in order to determine the RPF-Neighbor for a multicast flow. It merely discovers it by snooping the Join/Prune message. A PE could have snooped on two or more different Join/Prune messages for the same (S,G) that could have carried different Upstream-Neighbor fields. This could happen during transient network conditions or due to dual-homed sources. A PE cannot make assumptions on which one to pick, but instead must facilitate the CE routers decide which Upstream Neighbor gets elected the RPF-Neighbor. And for this

purpose, the PE will have to track downstream and upstream Join/Prune states per (S,G,N).

3.3.3. Receiving (*,G) PIM-SM Join/Prune Messages

A Join(*,G) or Prune(*,G) is considered "received" if the following conditions are met:

- o The port on which it arrived is not RPort(N) where N is the upstream-neighbor N of the Join/Prune(*,G), or,
- o if both RPort(N) and the arrival port are PWs, then there exists at least one other (*,G,Nx) or (Sx,G,Nx) state with an AC UpstreamPort.

For simplicity, the case where both RPort(N) and the arrival port are PWs is referred to as PW-only Join/Prune in this document. The PW-only Join/Prune handling is so that the RPort(N) PW can be added to the related forwarding entries' OutgoingPortList to trigger Assert, but that is only needed for those states with AC UpstreamPort. Note that in PW-only case, it is ok for the arrival port and RPort(N) to be the same. See Appendix Appendix B for examples.

When a router receives a Join(*,G) or a Prune(*,G) with upstream neighbor N, it must process the message as defined in the state machine below. Note that the macro computations of the various macros resulting from this state machine transition is exactly as specified in the PIM-SM RFC [PIM-SM].

We define the following per-port (*,G,N) macro to help with the state machine below.

Figure 1 : Downstream per-port (*,G) state machine in tabular form

Event	Previous State		
	NoInfo (NI)	Join (J)	Prune-Pend
Receive Join(*,G)	-> J state Action RxJoin(N)	-> J state Action RxJoin(N)	-> J state Action RxJoin(N)
Receive Prune(*,G) and NumETsActive<=1	-	-> PP state Start PPT(N)	-> PP state
Receive Prune(*,G) and NumETsActive>1	-	-> J state Start PPT(N)	-
PPT(N) expires	-	-> J state Action PPTExpiry(N)	-> NI state Action PPTExpiry(N)
ET(N) expires and NumETsActive<=1	-	-> NI state Action ETExpiry(N)	-> NI state Action ETExpiry(N)
ET(N) expires and NumETsActive>1	-	-> J state Action ETExpiry(N)	-> NI state Action ETExpiry(N)

Action RxJoin(N):

If ET(N) is not already running, then start ET(N). Otherwise restart ET(N). If N is not already in UpstreamNeighbors(*,G), then add N to UpstreamNeighbors(*,G) and trigger a Join(*,G) with upstream neighbor N to be forwarded upstream. If there are RPF Vector TLVs in the received (*,G) message and if they are different from the recorded RpfVectorTlvs(*,G), then copy them into RpfVectorTlvs(*,G).

Action PPTExpiry(N):

Same as Action ETExpiry(N) below, plus Send a Prune-Echo(*,G) with upstream-neighbor N on the downstream port.

Action ETExpiry(N):

Disable timers ET(N) and PPT(N). Delete neighbor state (Port,*,G,N). If there are no other (Port,*,G) states with NumETsActive(Port,*,G) > 0, transition DownstreamJPState to NoInfo. If there are no other (Port,*,G,N) state (different ports but for the same N), remove N from UpstreamPorts(*,G) - this also serves as a trigger for US FSM (JoinDesired(*,G,N) becomes FALSE).

3.3.4. Receiving (S,G) PIM-SM Join/Prune Messages

A Join(S,G) or Prune(S,G) is considered "received" if the following conditions are met:

- o The port on which it arrived is not Rport(N) where N is the upstream-neighbor N of the Join/Prune(S,G), or,
- o if both RPort(N) and the arrival port are PWs, then there exists at least one other (*,G,Nx) or (S,G,Nx) state with an AC UpstreamPort.

For simplicity, the case where both RPort(N) and the arrival port are PWs is referred to as PW-only Join/Prune in this document. The PW-only Join/Prune handling is so that the RPort(N) PW can be added to the related forwarding entries' OutgoingPortList to trigger Assert, but that is only needed for those states with AC UpstreamPort. See Appendix Appendix B for examples.

When a router receives a Join(S,G) or a Prune(S,G) with upstream neighbor N, it must process the message as defined in the state machine below. Note that the macro computations of the various macros resulting from this state machine transition is exactly as specified in the PIM-SM RFC [PIM-SM].

Figure 2: Downstream per-port (S,G) state machine in tabular form

Event	Previous State		
	NoInfo (NI)	Join (J)	Prune-Pend
Receive Join(S,G)	-> J state Action RxJoin(N)	-> J state Action RxJoin(N)	-> J state Action RxJoin(N)
Receive Prune (S,G) and NumETsActive<=1	-	-> PP state Start PPT(N)	-> PP state
Receive Prune(S,G) and NumETsActive>1	-	-> J state Start PPT(N)	-
PPT(N) expires	-	-> J state Action PPTEpiry(N)	-> NI state Action PPTEpiry(N)
ET(N) expires and NumETsActive<=1	-	-> NI state Action ETExpiry(N)	-> NI state Action ETExpiry(N)
ET(N) expires and NumETsActive>1	-	-> J state Action ETExpiry(N)	-> NI state Action ETExpiry(N)

Action RxJoin(N):

If ET(N) is not already running, then start ET(N). Otherwise, restart ET(N).

If N is not already in UpstreamNeighbors(S,G), then add N to UpstreamNeighbors(S,G) and trigger a Join(S,G) with upstream neighbor N to be forwarded upstream. If there are RPF Vector TLVs in the received (S,G) message and if they are different from the recorded RpfVectorTlvs(S,G), then copy them into RpfVectorTlvs(S,G).

Action PPTEpiry(N):

Same as Action ETExpiry(N) below, plus Send a Prune-Echo(S,G) with upstream-neighbor N on the downstream port.

Action ETEpiry(N):

Disable timers ET(N) and PPT(N). Delete neighbor state (Port,S,G,N). If there are no other (Port,S,G) states with NumETsActive(Port,S,G) > 0, transition DownstreamJPState to NoInfo. If there are no other (Port,S,G,N) state (different ports but for the same N), remove N from UpstreamPorts(S,G) - this also serves as a trigger for US FSM (JoinDesired(S,G,N) becomes FALSE).

3.3.5. Receiving (S,G,rpt) Join/Prune Messages

A Join(S,G,rpt) or Prune(S,G,rpt) is "received" when the port on which it was received is not also the port on which the upstream-neighbor N of the Join/Prune(S,G,rpt) was learnt.

While it is important to ensure that the (S,G) and (*,G) state machines allow for handling per (S,G,N) states, it is not as important for (S,G,rpt) states. It suffices to say that the downstream (S,G,rpt) state machine is the same as what is defined in section 4.5.4 of the PIM-SM RFC [PIM-SM].

3.3.6. Sending Join/Prune Messages Upstream

This section applies only to a PIM Proxy Switch and not to a PIM Snooping Switch.

A PIM Proxy PE MUST implement the Upstream FSM for which the procedures are similar to what is defined in section 4.5.6 of [PIM-SM]. Similar to Downstream FSM described above, the Upstream FSM is also per Upstream Neighbor.

For the purposes of the Upstream FSM, a Join or Prune message with upstream neighbor N is "seen" on a PIM Snooping switch if the port on which the message was received is also Rport(N), and the port is an AC. The AC requirement is needed because a Join received on the Rport(N) PW must not suppress this PE's Join on that PW.

In order to correctly facilitate assert among the CE routers, such Join/Prunes need to sent not only towards the upstream neighbor, but also on certain PWs as described below.

If RpfVectorTlvs(*,G) is not empty, then it must be encoded in a Join(*,G) message sent upstream.

If RpfVectorTlvs(S,G) is not empty, then it must be encoded in a Join(S,G) message sent upstream.

3.3.6.1. Where to send Join/Prune messages

The following rules apply, to both refresh and triggered (S,G)/(*,G) Join/Prune messages.

- o The upstream neighbor field N in the Join/Prune to be sent is set to the N in the corresponding Upstream FSM.
- o if Rport(N) is an AC, send the message to Rport(N).
- o Additionally, if OutgoingPortList(X,G,N) contains at least one AC, then the message MUST be sent to at least all the PWs in UpstreamPorts(G) (for (*,G)) or InheritedUpstreamPorts(S,G) (for (S,G)). Alternatively, the message MAY be sent to all PWs.

Sending to a subset of PWs as described above guarantees that if traffic (of the same flow) from two upstream routers were to reach this PE, then the two routers will receive from each other, triggering assert.

Sending to all PWs guarantees that if two upstream routers both send traffic for the same flow (even if it is to different set of downstream PEs), then they'll receive from each other, triggering assert.

3.4. Bidirectional-PIM (PIM-BIDIR)

PIM-BIDIR is a variation of PIM-SM. The main differences between PIM-SM and Bidirectional-PIM are as follows:

- o There are no source-based trees, and source-specific multicast is not supported (i.e., no (S,G) states) in PIM-BIDIR.
- o Multicast traffic can flow up the shared tree in PIM-BIDIR.
- o To avoid forwarding loops, one router on each link is elected as the Designated Forwarder (DF) for each RP in PIM-BIDIR.

The main advantage of PIM-BIDIR is that it scales well for many-to-many applications. However, the lack of source-based trees means that multicast traffic is forced to remain on the shared tree.

As described in [PIM-BIDIR], parts of a PIM-BIDIR enabled network may forward traffic without exchanging Join/Prune messages, for instance between DF's and the RPL.

As the described procedures for Pim snooping rely on the presence of Join/Prune messages, enabling Pim snooping on PIM-BIDIR networks

could break the PIM-BIDIR functionality. Deploying Pim snooping on PIM-BIDIR enabled networks will require some further study, some thoughts are gathered in Appendix A.

3.5. Interaction with IGMP Snooping

Whenever IGMP Snooping is enabled in conjunction with PIM Snooping in the same VPLS instance the switch SHOULD follow these rules:

- o To maintain the list of multicast routers and ports on which they are attached, the switch SHOULD NOT use the rules as described in section 2.1.1.(1) of RFC4541 [IGMP-SNOOP] but SHOULD rely on the neighbors discovered by PIM Snooping . This list SHOULD then be used to apply the forwarding rule as described in 2.1.1.(1) of RFC4541 [IGMP-SNOOP].
- o If the switch supports proxy-reporting, as described in section 2.1.1.(2) of RFC4541 [IGMP-SNOOP], all IGMP membership information learned on a port to which a PIM neighbor is attached SHOULD NOT be included in the summarized upstream report.

3.6. PIM-DM

The characteristics of PIM-DM is flood and prune behavior. Shortest path trees are built as a multicast source starts transmitting.

3.6.1. Building PIM-DM Snooping States

PIM-DM Snooping states are built by snooping on the PIM-DM Join, Prune, Graft and State Refresh messages received on AC/PWs and State-Refresh Messages sent on AC/PWs. By snooping on these PIM-DM messages, a PE builds the following states per (S,G,N) where S is the address of the multicast source, G is the Group address and N is the upstream neighbor to which Prunes/Grafts are sent by downstream CEs:

Per PIM (S,G,N):

Port PIM (S,G,N) Prune State:

- * DownstreamPState(S,G,N,Port): One of {"NoInfo" (NI), "Pruned" (P), "PrunePending" (PP)}
- * Prune Pending Timer (PPT)
- * Prune Timer (PT)

- * Upstream Port (valid if the PIM(S,G,N) Prune State is "Pruned").

3.6.2. PIM-DM Downstream Per-Port PIM(S,G,N) State Machine

The downstream per-port PIM(S,G,N) state machine is as defined in section 4.4.2 of [PIM-DM] with a few changes relevant to PIM Snooping. When reading section 4.4.2 of [PIM-DM] for the purposes of PIM-Snooping please be aware that the downstream states are built per (S, G, N, Downstream-Port} in PIM-Snooping and not per {Downstream-Interface, S, G} as in a PIM-DM router. As noted in the previous section Section 3.6.1, the states (DownstreamPState) and timers (PPT and PT) are per (S,G,N,P).

3.6.3. Triggering ASSERT election in PIM-DM

Since PIM-DM is a flood-and-prune protocol, traffic is flooded to all routers unless explicitly pruned. Since PIM-DM routers do not prune on non-RPF interfaces, PEs should typically not receive Prunes on Rport(RPF-neighbor). So the asserting routers should typically be in pim_oiflist(S,G). In most cases, assert election should occur naturally without any special handling since data traffic will be forwarded to the asserting routers.

However, there are some scenarios where a prune might be received on a port which is also an upstream port (UP). If we prune the port from pim_oiflist(S,G), then it would not be possible for the asserting routers to determine if traffic arrived on their downstream port. This can be fixed by adding pim_iifs(S,G) to pim_oiflist(S,G) so that data traffic flows to the UP ports.

3.7. PIM Proxy

As noted earlier, PIM Snooping will work correctly only if Join Suppression is disabled in the VPLS. If Join Suppression is enabled in the VPLS, then PEs MUST do PIM Proxy for VPLS Multicast to work correctly.

A PIM Proxy switch behaves like a PIM Router by doing most of the functionality of a PIM Router. The complexity however is much lesser on a switch since many of the issues that a PIM Router has to deal with are not relevant on a switch. A PIM Router needs to be able to build and maintain RP-Sets. They also have to deal with the Register and Assert State Machines. There are other complexities for a PIM Router resulting from inter-domain multicast. A PIM Snooping or PIM Proxy switch can be agnostic of all of this. All that a PIM Proxy switch cares about is building multicast states using PIM Hellos and PIM Join/Prune message. As such it's complexity is greatly reduced.

Other than the procedures defined here, the rest of the procedures that apply to PIM Snooping apply to PIM Proxy as well.

3.7.1. Downstream PIM Proxy behavior

Only PIM Join/Prune messages are proxied. Hellos MUST be snooped while being flooded in the VPLS. i.e. PIM Hellos MUST NOT be consumed at a PE and regenerated.

All other PIM packet types are flooded in the VPLS without any processing.

Performing only proxy of Join/Prune messages keeps the switch behavior very similar to that of a PIM router without introducing too much additional complexity. It keeps the PIM Proxy solution fairly simple. Since Join/Prunes are forwarded by a PE along the slow-path and all other PIM packet types are forwarded along the fast-path, it is very likely that packets forwarded along the fast-path will arrive "ahead" of Join/Prune packets at a CE router (note the stress on the fact that fast-path messages will never arrive after Join/Prunes). Of particular importance are Hello packets sent along the fast-path. We can construct a variety of scenarios resulting in out of order delivery of Hellos and Join/Prune messages. However, there should be no deviation from normal expected behavior observed at the CE router receiving these messages out of order.

The other option for a PIM Proxy solution is to proxy both Hello and Join/Prune messages that a PE is interested in building states for. If Hellos are being proxied, then it becomes necessary that the PE proxy all other PIM packet types also. Because if Hellos are received after other packet types are received at a CE router, then bad things will happen. That means every PIM packet has to be sent along the slow-path. This greatly increases the complexity on the CE router, it is very compute intensive and does not scale well. Also, proxying Hellos will result in added latency to delivery of Hello messages to a CE and that affects multicast convergence in the VPLS.

3.7.2. Upstream PIM Proxy behavior

Since a PIM Proxy switch consumes Join/Prune messages, it must also originate PIM Join/Prune messages to be sent upstream. On ACs, both triggered and refresh Join/Prunes are forwarded as PIM packets.

3.7.3. Source IP Address in Proxy PIM Join/Prune Packets

The source IP address in PIM packets sent upstream SHOULD be the address of a PIM downstream neighbor in the corresponding join/prune state. The address picked MUST NOT be the upstream neighbor field to

be encoded in the packet. The layer 2 encapsulation for the selected source IP address MUST be the encapsulation recorded in the PIM Neighbor database for that IP address.

If Explicit Tracking (ET) is disabled in the VPLS, then it does not matter what Source IP Address is picked in the packets sent upstream as long as we adhere to the rule in the previous paragraph.

If ET is enabled, it means that a CE router is interested in tracking every CE that wishes to join a stream. If a PE determines that ET is enabled, then it SHOULD use PIM Snooping procedures instead of PIM Proxy.

3.8. Directly Connected Multicast Source

If there is a source in the CE network that connects directly into the VPLS instance, then multicast traffic from that source MUST be sent to all PIM routers on the VPLS instance apart from the igmp receivers in the VPLS. If there is already (S,G) or (*,G) snooping state that is formed on any PE, this will not happen per the current forwarding rules and guidelines. So, in order to determine if traffic needs to be flooded to all routers, a PE must be able to determine if the traffic came from a host on that LAN. There are three ways to address this problem:

- o The PE would have to do ARP snooping to determine if a source is directly connected.
- o Another option is to have configuration on all PEs to say there are CE sources that are directly connected to the VPLS instance and disallow snooping for the groups for which the source is going to send traffic. This way traffic from that source to those groups will always be flooded within the provider network.
- o A third option is to require that sources of CE multicast routers must appear behind a router.

3.9. Data Forwarding Rules

First we define the rules that are common to PIM-SM, PIM-BIDIR and PIM-DM PEs. Forwarding rules for each protocol type is specified in the sub-sections.

If there is no matching forwarding state, then the PE MAY either discard the packet or send it towards all the snooped PIM CE routers or to a configured set of ports. How this is determined is outside the scope of this document.

The following general rules MUST be followed when forwarding multicast traffic in a VPLS:

- o Traffic arriving on a port MUST NOT be forwarded back onto the same port.
- o Due to VPLS Split-Horizon rules, traffic ingressing on a PW MUST NOT be forwarded to any other PW.

3.9.1. PIM-SM Data Forwarding Rules

Per the rules in [PIM-SM] and per the additional rules specified in this document,

```
OutgoingPortList(*,G) = immediate_olist(*,G) (+)
                        UpstreamPorts(*,G) (+)
                        Rport(PimDR)
```

```
OutgoingPortList(S,G) = inherited_olist(S,G) (+)
                        UpstreamPorts(S,G) (+)
                        (UpstreamPorts(*,G) (-)
                        UpstreamPorts(S,G,rpt)) (+)
                        Rport(PimDR)
```

[PIM-SM] specifies how `immediate_olist(*,G)` and `inherited_olist(S,G)` are built. `PimDR` is the IP address of the PIM DR in the VPLS.

The PIM-SM Snooping forwarding rules are defined below in pseudocode:

```
BEGIN
  iif is the incoming port of the multicast packet.
  S is the Source IP Address of the multicast packet.
  G is the Destination IP Address of the multicast packet.

  If there is (S,G) state on the PE
  Then
    OutgoingPortList = OutgoingPortList(S,G)
  Else if there is (*,G) state on the PE
  Then
    OutgoingPortList = OutgoingPortList(*,G)
  Else
    OutgoingPortList = UserDefinedPortList
  Endif

  If iif is an AC
  Then
    OutgoingPortList = OutgoingPortList (-) iif
  Else
    ## iif is a PW
    OutgoingPortList = OutgoingPortList (-) PWPorts
  Endif

  Forward the packet to OutgoingPortList.
END
```

First if there is (S,G) state on the PE, then the set of outgoing ports is OutgoingPortList(S,G).

Otherwise if there is (*,G) state on the PE, the set of outgoing ports is OutgoingPortList(*,G).

The packet is forwarded to the selected set of outgoing ports while observing the general rules above in section Section 3.9

3.9.2. PIM-BIDIR Data Forwarding Rules

The PIM-BIDIR Snooping forwarding rules are defined below in pseudocode:

```
BEGIN
  iif is the incoming port of the multicast packet.
  G is the Destination IP Address of the multicast packet.

  If there is forwarding state for G
  Then
    OutgoingPortList = olist(G)
  Else
    OutgoingPortList = UserDefinedPortList
  Endif

  If iif is an AC
  Then
    OutgoingPortList = OutgoingPortList (-) iif
  Else
    ## iif is a PW
    OutgoingPortList = OutgoingPortList (-) PWPorts
  Endif

  Forward the packet to OutgoingPortList.
END
```

If there is forwarding state for G, then forward the packet to olist(G) while observing the general rules above in section Section 3.9

[PIM-BIDIR] specifies how olist(G) is constructed.

3.9.3. PIM-DM Data Forwarding Rules

The PIM-DM Snooping data forwarding rules are defined below in pseudocode:

```
BEGIN
    iif is the incoming port of the multicast packet.
    S is the Source IP Address of the multicast packet.
    G is the Destination IP Address of the multicast packet.

    If there is (S,G) state on the PE
    Then
        OutgoingPortList = olist(S,G)
    Else
        OutgoingPortList = UserDefinedPortList
    Endif

    If iif is an AC
    Then
        OutgoingPortList = OutgoingPortList (-) iif
    Else
        ## iif is a PW
        OutgoingPortList = OutgoingPortList (-) PWPorts
    Endif

    Forward the packet to OutgoingPortList.
END
```

If there is forwarding state for (S,G), then forward the packet to olist(S,G) while observing the general rules above in section Section 3.9

[PIM-DM] specifies how olist(S,G) is constructed.

4. IANA Considerations

This document makes no request of IANA.

Note to RFC Editor: this section may be removed on publication as an RFC.

5. Security Considerations

Security considerations provided in VPLS solution documents (i.e., [VPLS-LDP] and [VPLS-BGP]) apply to this document as well.

6. Contributors

Karl (Xiangrong) Cai and Princy Elizabeth made significant contributions to bring the specification to its current state,

especially in the area of Join forwarding rules.

Yetik Serbest, Ray Qiu, Suresh Boddapati co-authored earlier versions.

7. Acknowledgements

Many members of the L2VPN and PIM working groups have contributed to and provided valuable comments and feedback to this draft, including Vach Kompella, Shane Amante, Sunil Khandekar, Rob Nath, Marc Lassere, Yuji Kamite, Yiqun Cai, Ali Sajassi, Jozef Raets, Himanshu Shah (Ciena), Himanshu Shah (Alcatel-Lucent).

8. References

8.1. Normative References

- [PIM-BIDIR] Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano, "Bidirectional Protocol Independent Multicast (BIDIR-PIM)", RFC 5015, 2007.
- [PIM-DM] Adams, A., Nicholas, J., and W. Siadak, "Protocol Independent Multicast Version 2 - Dense Mode Specification", RFC 3973, 2005.
- [PIM-SM] Fenner, B., Handley, M., Holbrook, H., and I. Kouvelas, "Protocol Independent Multicast- Sparse Mode (PIM-SM): Protocol Specification (Revised)", RFC 4601, 2006.
- [PIM-SSM] Holbrook, H. and B. Cain, "Source-Specific Multicast for IP", RFC 4607, 2006.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, 1997.
- [RPF-VECTOR] Wijnands, I., Boers, A., and E. Rosen, "The Reverse Path Forwarding (RPF) Vector TLV", RFC 5496, 2009.

8.2. Informative References

- [IGMP-SNOOP] Christensen, M., Kimball, K., and F. Solensky, "Considerations for IGMP and MLD Snooping Switches", RFC 4541, 2006.

[VPLS-BGP]

Kompella, K. and Y. Rekhter, "Virtual Private LAN Service using BGP for Auto-Discovery and Signaling", RFC 4761, 2007.

[VPLS-LDP]

Lasserre, M. and V. Kompella, "Virtual Private LAN Services using LDP Signaling", RFC 4762, 2007.

[VPLS-MCAST-REQ]

Kamite, Y., Wada, Y., Serbest, Y., Morin, T., and L. Fang, "Requirements for Multicast Support in Virtual Private LAN Services", RFC 5501, 2009.

[VPLS-MCAST-TREES]

Aggarwal, R., Kamite, Y., Fang, L., and Y. Rekhter, "Multicast in VPLS", draft-ietf-l2vpn-vpls-mcast-10, Work in Progress.

Appendix A. PIM-BIDIR Thoughts

This section describes some guidelines that may be used to preserve PIM-BIDIR functionality in combination with Pim Snooping.

In order to preserve PIM-BIDIR Pim snooping routers need to set up forwarding states so that :

- o on the RPL all traffic is forwarded to all Rport(N)
- o on any other interface traffic is always forwarded to the DF

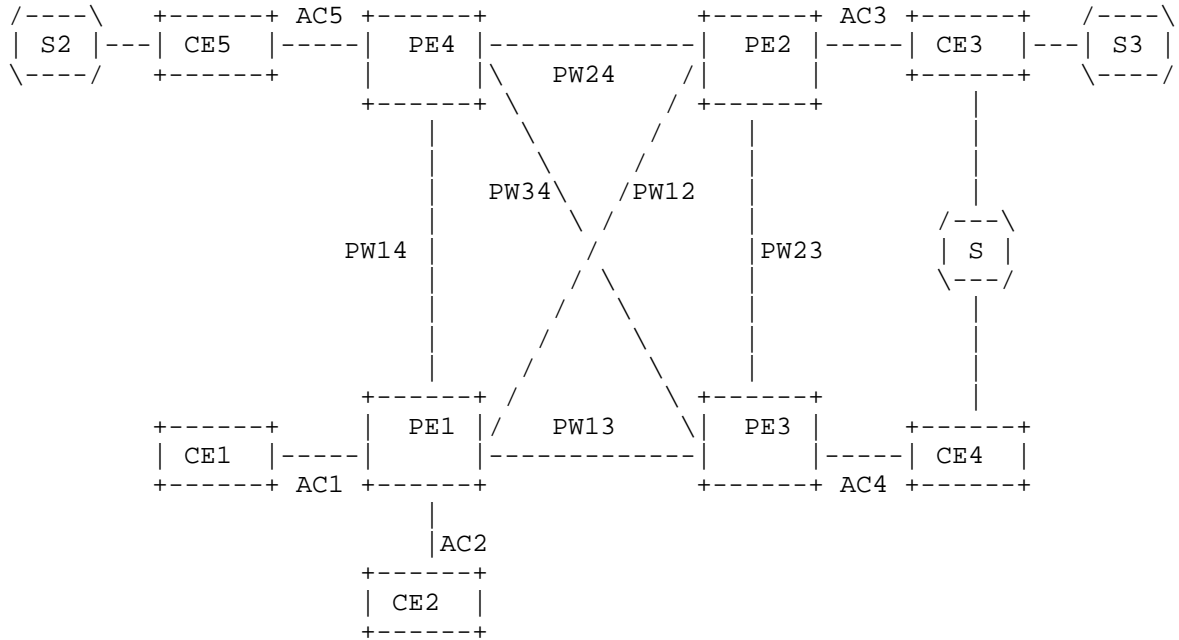
The information needed to setup these states may be obtained by :

- o determining the mapping between group(range) and RP
- o snooping and storing DF election information
- o determining where the RPL is, this could be achieved by static configuration, or by combining the information mentioned in previous bullets.

Appendix B. Example Network Scenario

Let us consider the scenario in Figure 3.

An Example Network for Triggering Assert



In the examples below, $JT(\text{Port}, S, G, N)$ is the downstream Join Expiry Timer on the specified Port for the (S, G) with upstream neighbor N .

B.1. Pim Snooping Example

In the network depicted in Figure 3, S is the source of a multicast stream (S, G) . $CE1$ and $CE2$ both have two ECMP routes to reach the source.

1. $CE1$ Sends a $\text{Join}(S, G)$ with $\text{UpstreamNeighbor}(S, G) = CE3$.
2. $PE1$ snoops on the $\text{Join}(S, G)$ while flooding it in the VPLS. $PE2$ and $PE3$ also snoop on the $\text{Join}(S, G)$ while flooding it in the VPLS.

The resulting states at the PEs is as follows:

At $PE1$:

```

JT(AC1, S, G, CE3)      = JP_HoldTime
UpstreamNeighbors(S, G) = { CE3 }
UpstreamPorts(S, G)     = { PW12 }
OutgoingPortList(S, G)  = { AC1, PW12 }

```

At $PE2$:

```

JT(PW12, S, G, CE3)     = JP_HoldTime
UpstreamNeighbors(S, G) = { CE3 }

```

```
UpstreamPorts(S,G)      = { AC3 }
OutgoingPortList(S,G)   = { PW12, AC3 }
```

At PE3:

PE3 doesn't create a forwarding state for (S,G) because the Join(S,G) was received on a PW and the Upstream RPort is a PW too. <<<<<

3. The multicast stream (S,G) flows along CE3 -> PE2 -> PE1 -> CE1
4. Now CE2 sends a Join(S,G) with Upstream Neighbor(S,G) = CE4.
5. All PEs snoop on the Join(S,G).

The resulting states at the PEs:

At PE1:

```
JT(AC1,S,G,CE3)         = active
JT(AC2,S,G,CE4)         = JP_HoldTime.
UpstreamNeighbors(S,G)  = { CE3, CE4 }
UpstreamPorts(S,G)      = { PW12, PW13 }
OutgoingPortList(S,G)   = { AC1, PW12, AC2, PW13 }
```

At PE2: Note: Since PE2 already has (S,G) state, it does not ignore the Join(S,G) even though it received the Join(S,G) on a PW and the Upstream Rport is a PW. <<<<<<<

```
JT(PW12,S,G,CE4)        = JP_HoldTime
JT(PW12,S,G,CE3)        = active
UpstreamNeighbors(S,G)  = { CE3, CE4 }
UpstreamPorts(S,G)      = { AC3, PW23 }
OutgoingPortList(S,G)   = { PW12, AC3, PW23 }
```

At PE3:

```
JT(PW13,S,G,CE4)        = JP_HoldTime
UpstreamNeighbors(S,G)  = { CE4 }
UpstreamPorts(S,G)      = { AC4 }
OutgoingPortList(S,G)   = { PW13, AC4 }
```

6. The multicast stream (S,G) flows into the VPLS from the two CEs CE3 and CE4. PE2 forwards the stream received from CE3 to PW23 and PE3 forwards the stream to AC4. This facilitates the CE routers to trigger assert election. Let us say CE3 becomes the assert winner.
7. CE3 sends an Assert message to the VPLS. The PEs flood the Assert message without examining it.
8. CE4 stops sending the multicast stream to the VPLS.
9. CE2 notices an RPF change due to Assert and sends a Prune(S,G) with Upstream Neighbor = CE4. CE2 also sends a Join(S,G) with Upstream Neighbor = CE3.
10. All the PEs start a prune-pend timer on the ports on which

they received the Prune(S,G). When the prune-pend timer expires, all PEs will remove the downstream (S,G,CE4) states.

Resulting states at the PEs:

At PE1:

```

JT(AC1,S,G,CE3)      = active
UpstreamNeighbors(S,G) = { CE3 }
UpstreamPorts(S,G)   = { PW12 }
OutgoingPortList(S,G) = { AC1, AC2, PW12 }

```

At PE2:

```

JT(PW12,S,G,CE3)     = active
UpstreamNeighbors(S,G) = { CE3 }
UpstreamPorts(S,G)   = { AC3 }
OutgoingPortList(S,G) = { PW12, AC3 }

```

At PE3: no (S,G) state.

Note that at the end of the assert election, there should be no duplicate traffic forwarded downstream and traffic should flow only on the desired path. Also note that there are no unnecessary (S,G) states on PE3 after the assert election.

B.2. PIM Proxy Example with (S,G) / (*,G) interaction

In the same network, let us assume CE4 is the Upstream Neighbor towards the RP for G.

1. CE1 Sends a Join(S,G) with Upstream Neighbor(S,G) = CE3.
2. PE1 consumes the Join(S,G). PE1 looks up the neighbor database and determines CE3 was learnt on PW12. PE1 sends a Proxy Join(S,G) to the resulting UpstreamPorts(G). i.e. it sends the proxy Join(S,G) on PW12.
3. Likewise, PE2 consumes the Join(S,G) and sends a proxy Join(S,G) on AC3 with Upstream Neighbor = CE3.

The resulting states at the PEs is as follows:

At PE1:

```

JT(AC1,S,G,CE3)      = JP_HoldTime
UpstreamNeighbors(S,G) = { CE3 }
UpstreamPorts(S,G)   = { PW12 }
OutgoingPortList(S,G) = { AC1, PW12 }

```

At PE2:

```

JT(PW12,S,G,CE3)     = JP_HoldTime
UpstreamNeighbors(S,G) = { CE3 }

```

```
UpstreamPorts(S,G)      = { AC3 }
OutgoingPortList(S,G)   = { PW12, AC3 }
```

At PE3: PE3 did not receive any PIM Join(S,G). So it has no (S,G) state.

4. The multicast stream (S,G) flows along CE3 -> PE2 -> PE1 -> CE1.

5. Now let us say CE1 sends a Join(*,G) towards CE4.

6. PE1 consumes the Join(*,G). PE1 sends a Proxy Join(*,G) to the resulting UpstreamPorts(G). Since UpstreamPorts(G) now has both PW12 and PW13, the Join(*,G) gets sent on both PW12 and PW13.

Note that the UpstreamPorts(S,G) and OutgoingPortList(S,G) inherit the corresponding (*,G) sets, but not vice versa.

remove "but not vice versa"

COMMENT : > Original "but not vice versa" applies to OutgoingPortList(S,G) only, I assume, because of the earlier definition:

UpstreamPorts(G): This set is the union of all the UpstreamPorts(S,G) and UpstreamPorts(*,G) for a given G

7. PE2 and PE3 perform a similar function. PE2 received the Join(*,G) on a PW and the Upstream Neighbor is also on a PW. Hence PE2 only adds UpstreamPorts(*,G) to OutgoingPortList(*,G) and not the downstream port PW12.

At PE1:

```
JT(AC1,S,G,CE3)         = active
UpstreamNeighbors(S,G)   = { CE3 }
UpstreamPorts(S,G)       = { PW12, PW13 }
OutgoingPortList(S,G)    = { AC1, PW12, PW13 }
```

```
JT(AC1*,G,CE4)          = JP_HoldTime.
UpstreamNeighbors(*,G)   = { CE4 }
UpstreamPorts(*,G)       = { PW13 }
OutgoingPortList(*,G)    = { AC1, PW13 }
```

```
UpstreamPorts(G)         = { PW12, PW13 }
```

At PE2:

```
JT(PW12,S,G,CE3)        = active
UpstreamNeighbors(S,G)   = { CE3 }
UpstreamPorts(S,G)       = { AC3, PW23 }
OutgoingPortList(S,G)    = { PW12, AC3, PW23 }
```

```
JT(PW12*,G,CE4)         = JP_HoldTime
UpstreamNeighbors(*,G)   = { CE4 }
UpstreamPorts(G)         = { PW23 }
```

```
OutgoingPortList(*,G)      = { PW23 }
```

At PE3:

```
JT(PW13,*,G,CE4) = JP_HoldTime
UpstreamNeighbors(*,G) = { CE4 }
UpstreamPorts(*,G) = { AC4 }
OutgoingPortList(*,G) = { PW13, AC4 }
```

8. The above state results in both (S,G) and (*,G) streams to be forwarded to AC1. The above state also results in the (S,G) stream to be forwarded from CE3 to CE4 resulting in an (S,G) assert election. Following the assert election, CE3 becomes the (S,G) assert winner. CE4 stops sending (S,G) stream down the RPT.
9. CE1 notices an RPF change due to assert. It sends a Prune(S,G,rpt) with Upstream Neighbor = CE4.
10. PE1 consumes the Prune(S,G,rpt) and forwards the Prune(S,G,rpt) to both PW12 and PW13. PE2 consumes the Prune(S,G,rpt) and updates its states. PE3 updates its states and forwards the Prune(S,G,rpt) on AC4.

At PE1:

```
JT(AC1,S,G,CE3) = active
UpstreamNeighbors(S,G) = { CE3 }
UpstreamPorts(S,G) = { PW12 }
OutgoingPortList(S,G) = { AC1, PW12 }

JT(AC1,*,G,CE4) = active.
UpstreamNeighbors(*,G) = { CE4 }
UpstreamPorts(*,G) = { PW13 }
OutgoingPortList(*,G) = { AC1, PW13 }
```

At PE2:

```
JT(PW12,S,G,CE3) = active
UpstreamNeighbors(S,G) = { CE3 }
UpstreamPorts(*,G) = { AC3 }
OutgoingPortList(S,G) = { PW12, AC3 }

JT(PW12,*,G,CE4) = JP_HoldTime
UpstreamNeighbors(*,G) = { CE4 }
UpstreamPorts(*,G) = { PW23 }
OutgoingPortList(*,G) = { PW23 }
```

At PE3:

```
JT(PW13,*,G,CE4) = JP_HoldTime
UpstreamNeighbors(*,G) = { CE4 }
UpstreamPorts(G) = { AC4 }
OutgoingPortList(*,G) = { PW13, AC4 }
```

Even in this example, at the end of the (S,G) / (*,G) assert election, there should be no duplicate traffic forwarded downstream and traffic should flow only to the desired CEs.

Other more complex scenarios exist. This draft should address in PIM-SM and the rules specified in this draft should ensure that assert is triggered among the CEs in all scenarios.

Authors' Addresses

Olivier Dornon
Alcatel-Lucent
50 Copernicuslaan
Antwerp, B2018

Email: olivier.dornon@alcatel-lucent.com

Jayant Kotalwar
Alcatel-Lucent
701 East Middlefield Rd.
Mountain View, CA 94043

Email: jayant.kotalwar@alcatel-lucent.com

Jeffrey Zhang
Juniper Networks, Inc.
10 Technology Park Drive
Westford, MA 01886

Email: zzhang@juniper.net

Venu Hemige
Alcatel-Lucent
701 East Middlefield Rd.
Mountain View, CA 94043

Email: Venu.hemige@alcatel-lucent.com

