BESS Working Group Internet-Draft Updates: <u>6513</u>,6514 (if approved) Intended status: Standards Track Expires: February 16, 2017 E. Rosen, Ed. Juniper Networks, Inc. K. Subramanian Sproute Networks Z. Zhang Juniper Networks, Inc. August 15, 2016

Ingress Replication Tunnels in Multicast VPN draft-ietf-bess-ir-05

Abstract

RFCs 6513, 6514, and other RFCs describe procedures by which a Service Provider may offer Multicast VPN service to its customers. These procedures create point-to-multipoint (P2MP) or multipoint-tomultipoint trees across the Service Provider's backbone. One type of P2MP tree that may be used is known as an "Ingress Replication (IR) tunnel". In an IR tunnel, a parent node need not be "directly connected" to its child nodes. When a parent node has to send a multicast data packet to its child nodes, it does not use layer 2 multicast, IP multicast, or MPLS multicast to do so. Rather, it makes n individual copies, and then unicasts each copy, through an IP or MPLS unicast tunnel, to exactly one child node. While the prior MVPN specifications allow the use of IR tunnels, those specifications are not always very clear or explicit about how the MVPN protocol elements and procedures are applied to IR tunnels. This document updates RFCs 6513 and 6514 by adding additional details that are specific to the use of IR tunnels.

Status of This Memo

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

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 February 16, 2017.

Rosen, et al.

Expires February 16, 2017

Copyright Notice

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>http://trustee.ietf.org/license-info</u>) 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

$\underline{1}$. Introduction
2. What is an IR P-tunnel?
$\underline{3}$. How are IR P-tunnels Identified?
$\underline{4}$. How to Join an IR P-tunnel
<u>4.1</u> . Advertised IR P-tunnels
<u>4.1.1</u> . If the 'Leaf Info Required Bit' is Set <u>10</u>
4.1.2. If the 'Leaf Info Required Bit' is Not Set <u>10</u>
<u>4.2</u> . Unadvertised IR P-tunnels <u>11</u>
5. The PTA's 'Tunnel Identifier' Field
6. A Note on IR P-tunnels and 'Discarding Packets from the Wrong
PE'
7. The PTA's 'MPLS Label' Field
7.1. Leaf A-D Route Originated by an Egress PE <u>14</u>
7.2. Leaf A-D Route Originated by an Intermediate Node <u>16</u>
<u>7.3</u> . Intra-AS I-PMSI A-D Route <u>17</u>
8. How A Child Node Prunes Itself from an IR P-tunnel <u>17</u>
9. Parent Node Actions Upon Receiving Leaf A-D Route <u>18</u>
<u>10</u> . Use of Timers when Switching UMH \ldots \ldots \ldots \ldots 19
<u>11</u> . IANA Considerations
<u>12</u> . Acknowledgments
<u>13</u> . Security Considerations
<u>14</u> . References
<u>14.1</u> . Normative References
<u>14.2</u> . Informative References
Authors' Addresses

<u>1</u>. Introduction

RFCs 6513, 6514, and others describe procedures by which a Service Provider (SP) may offer Multicast VPN (MVPN) service to its customers. These procedures create point-to-multipoint (P2MP) or multipoint-to-multipoint (MP2MP) tunnels, called "P-tunnels" (Provider-tunnels), across the SP's backbone network. Customer multicast traffic is carried through the P-tunnels.

A number of different P-tunnel technologies are supported. One of the supported P-tunnel technologies is known as "ingress replication" or "unicast replication". We will use the acronym "IR" to refer to this P-tunnel technology.

An IR P-tunnel is a P2MP tree, but a given node on the tree is not necessarily "directly attached" to its parent node or to its child nodes. To send a multicast data packet from a parent node to one of its child nodes, the parent node encapsulates the packet and then unicasts it through a tunnel to the child node. The tunnel may be a P2P (point-to-point) or MP2P (multipoint-to-point) MPLS LSP (label switched path) or a unicast IP tunnel. If a node on an IR tree has n child nodes, and has a multicast data packet that must be sent along the tree, the parent node makes n individual copies of the data packet, and then sends each copy, through a unicast tunnel, to exactly one child node. No lower layer multicast technology is used when sending traffic from a parent node to a child node; multiple copies of the packet may therefore be sent out a single interface.

With the single exception of IR, the P-tunnel technologies supported by the MVPN specifications are pre-existing IP multicast or MPLS multicast technologies. Each such technology has its own set of specifications, its own setup and maintenance protocols, its own syntax for identifying specific multicast trees, and its own procedures for enabling a router to be added to or removed from a particular multicast tree. For IR P-tunnels, on the other hand, there is no prior specification for setting up and maintaining the P2MP trees; the procedures and protocol elements used for setting up and maintaining the P2MP trees are specified in the MVPN specifications themselves, and all the signaling/setup is done by using the BGP A-D (Auto-Discovery) routes that are defined in [RFC6514]. (The unicast tunnels used to transmit multicast data from one node to another in an IR P-tunnel may of course have their own setup and maintenance protocols, e.g., [RFC5036], [RFC3209].)

Since the transmission of a multicast data packet along an IR P-tunnel is done by transmitting the packet through a unicast tunnel, previous RFCs sometimes speak of an IR P-tunnel as "consisting of" a set of unicast tunnels. However, that way of speaking is not quite

accurate. For one thing, it obscures the fact that an IR P-tunnel is really a P2MP tree, whose nodes must maintain multicast state in both the control and data planes. For another, it obscures the fact the unicast tunnels used by a particular IR P-tunnel need not be specific to that P-tunnel; a single unicast tunnel can carry the multicast traffic of many different IR P-tunnels (and can also carry unicast traffic as well).

In this document, we provide a clearer and more explicit conceptual model for IR P-tunnels, clarifying the relationship between an IR P-tunnel and the unicast tunnels that are used for data transmission along the IR P-tunnel.

Section 5 of [RFC6514] defines a BGP Path Attribute known as the "PMSI (Provider Multicast Service Interface) Tunnel attribute" (PTA). This attribute contains a field known as the "Tunnel Identifier" field. For most P-tunnel technologies, the PTA's "Tunnel Identifier" field is used to identify a P-tunnel (i.e., to identify a P2MP or MP2MP tree). However, when IR P-tunnels are used, the PTA "Tunnel Identifier" field does not actually identify an IR P-tunnel. In some cases it identifies one of the P-tunnel's constituent unicast tunnels, and in other cases it is not used to identify a tunnel at all. In this document, we provide an explicit specification for how IR P-tunnels are actually identified.

Some of the MVPN specifications use phrases like "join the identified P-tunnel", even though there has up to now not been an explicit specification of how to identify an IR P-tunnel, of how a router joins such a P-tunnel, or of how a router prunes itself from such a P-tunnel. In this document, we make these procedures more explicit.

[RFC6514] does provide a method for binding an MPLS label to a P-tunnel, but does not discuss the label allocation policies that are needed for correct operation when the P-tunnel is an IR P-tunnel. Those policies are discussed in this document.

This document does not provide any new protocol elements, or any fundamentally new procedures; its purpose is to make explicit just how a router is to use the protocol elements and procedures of [<u>RFC6513</u>] and [<u>RFC6514</u>] to identify an IR P-tunnel, to join an IR P-tunnel, and to prune itself from an IR P-tunnel.

This document also discusses the MPLS label allocation policies that need to be supported when binding MPLS labels to IR P-tunnels, and the timer policies that need to be supported when switching a customer multicast flow from one IR P-tunnel to another. These are procedures that are not clearly specified in [<u>RFC6513</u>] or [<u>RFC6514</u>]. As the material in this document must be understood in order to

properly implement IR P-tunnels, this document is considered to update [<u>RFC6513</u>] and [<u>RFC6514</u>].

This document also discusses the application of "seamless multicast" [RFC7524] and "extranet" [RFC7900] procedures to IR P-tunnels.

This draft does not discuss the use of IR P-tunnels to support a VPN customer's use of Bidirectional Protocol Independent Multicast (BIDIR-PIM). [RFC7740] explains how to adapt the procedures of [RFC6513], [RFC6514], and [RFC7582] so that a customer's use of BIDIR-PIM can be supported by IR P-tunnels.

In the event of any conflict between this document and either [<u>RFC6513</u>] or [<u>RFC6514</u>], this document takes precedence.

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL", when and only when appearing in all capital letters, are to be interpreted as described in [RFC2119].

2. What is an IR P-tunnel?

An IR P-tunnel is a P2MP tree. Its nodes are BGP speakers that support the MVPN procedures of [RFC6514] and related RFCs. In general, the nodes of an IR P-tunnel are either Provider Edge (PE) routers, Autonomous System Border Routers (ASBRs), or (if [RFC7524] is supported) Area Border Routers (ABRs). (MVPN procedures are sometimes used to support non-MVPN, or "global table" multicast; one way of doing this is defined in [RFC7524]. Another way is defined in [RFC7716]. In such cases, IR P-tunnels can be used outside the context of MVPN.)

MVPN P-tunnels may be either "segmented" or "non-segmented" (as these terms are defined in [<u>RFC6513</u>] and [<u>RFC6514</u>]).

A "non-segmented" IR P-tunnel is a two-level P2MP tree, consisting only of a root node and a set of nodes that are children of the root node. When used in an MVPN context, the root is an ingress PE, and the child nodes of the root are the egress PEs.

In a segmented P-tunnel, IR may be used for some or all of the segments. If a particular segment of a segmented P-tunnel uses IR, then the root of that segment may have child nodes that are ABRs or ASBRs, rather than egress PEs.

As with any type of P2MP tree, each node of an IR P-tunnel holds "multicast state" for the P-tunnel. That is, each node knows the identity of its parent node on the tree, and each node knows the

identities of its child nodes on the tree. In the MVPN specs, the "parent" node is also known as the "Upstream Multicast Hop" or "UMH". Note that the "UMH" may be a PE, an ASBR, or (if procedures from [<u>RFC7524</u>] are being used) an ABR. (In [<u>RFC7524</u>], the term "upstream node" is used instead of "UMH".)

What distinguishes an IR P-tunnel from any other kind of P2MP tree is the method by which a data packet is transmitted from a parent node to a child node. To transmit a multicast data packet from a parent node to a child node along a particular IR P-tunnel, the parent node does the following:

- o It labels the packet with a label (call it a "P-tunnel label") that the child node has assigned to that P-tunnel,
- o It then places the packet in a unicast encapsulation and unicasts the packet to the child node. That is, the parent node sends the packet through a "unicast tunnel" to a particular child node. This unicast tunnel need not be specially created to be part of the IR P-tunnel; it can be any P2P or MP2P unicast tunnel that will get the packets from the parent node to the child node. A single such unicast tunnel may be carrying multicast data packets of several different P2MP trees, and may also be carrying unicast data packets.

The parent node repeats this process for each child node, creating one copy for each child node, and sending each copy through a unicast tunnel to corresponding child node. It does not use layer 2 multicast, IP multicast, or MPLS multicast to transmit packets to its child nodes. As a result, multiple copies of each packet may be sent out a single interface; this may happen, e.g., if that interface is the next hop interface, according to unicast routing, from the parent node to several of the child nodes.

Since data traveling along an IR P-tunnel is always unicast from parent node to child node, it can be convenient to think of an IR P-tunnel as a P2MP tree whose arcs are unicast tunnels. However, it is important to understand that the unicast tunnels need not be specific to any particular IR P-tunnel. If R1 is the parent node of R2 on two different IR P-tunnels, a single unicast tunnel from R1 to R2 may be used to carry data along both IR P-tunnels. All that is required is that when the data packets arrive at R2, R2 will see the "P-tunnel label" at the top of the packets' label stack; R2's further processing of the packets will depend upon that label. Note that the same unicast tunnel between R1 and R2 may also be carrying unicast data packets.

Typically the unicast tunnels are the Label Switched Paths (LSPs) that already exist to carry unicast traffic; either MP2P LSPs created by LDP (Label Distribution Protocol, [RFC5036]) or P2P LSPs created by RSVP-TE (Resource Reservation Protocol - Traffic Engineering, [RFC3209]). However, any other kind of unicast tunnel may be used. A unicast tunnel may have an arbitrary number of intermediate routers; those routers do not maintain any multicast state for the IR P-tunnel, and in general are not even aware of its existence.

As with all other P-tunnel types, an IR P-tunnel may be used to instantiate either an Inclusive PMSI or a Selective PMSI. See <u>Section 3.2 of [RFC6513]</u> for an explanation of those concepts.

3. How are IR P-tunnels Identified?

There are four MVPN BGP route types in which P-tunnels can be identified: Intra-AS I-PMSI A-D (Intra Autonomous System Inclusive PMSI A-D) routes, Inter-AS I-PMSI A-D routes, S-PMSI (Selective PMSI) A-D routes, and Leaf A-D routes. (These route types are all defined in [<u>RFC6514</u>]).

Whenever it is necessary to identify a P-tunnel in a route of one of these types, a "PMSI Tunnel Attribute" (PTA) is added to the route. As defined in [RFC6514] section 5, the PTA contains four fields: "Tunnel Type", "MPLS Label", "Tunnel Identifier", and "Flags". [RFC6514] defines only one bit in the "Flags" field, the "Leaf Information Required" bit.

If a route identifies an IR P-tunnel, the "Tunnel Type" field of its PTA is set to the value 6, meaning "Ingress Replication".

Most types of P-tunnel are associated with specific protocols that are used to set up and maintain tunnels of that type. For example, if the "Tunnel Type" field is set to 2, meaning "mLDP P2MP LSP", the associated setup protocol is mLDP [RFC6388]. The associated setup protocol always has a method of identifying the tunnels that it sets up. For example, mLDP uses a "FEC element" (Forwarding Equivalence Class Element) to identify a tree. If the "Tunnel type" field is set to 3, meaning "PIM SSM Tree" (Protocol Independent Multicast Source-Specific Tree), the associated setup protocol is PIM, and "(S,G)" is used to identify the tree. In these cases, the "Tunnel Identifier" field of the PTA carries a tree identifier as defined by the setup protocol used for the particular tunnel type.

IR P-tunnels, on the other hand, are entirely setup and maintained by the use of BGP A-D routes, and are not associated with any other setup protocol. (The unicast tunnels used to transmit multicast data along an IR P-tunnel may have their own setup and maintenance

protocols, of course.) The means of identifying a P-tunnel is very different for IR P-tunnels than for other types of P-tunnel:

When an IR P-tunnel is identified in an S-PMSI A-D route, an Intra-AS I-PMSI A-D route, or an Inter-AS I-PMSI A-D route (we will refer to these three route types as "advertising A-D routes"), its identifier is hereby defined to be the NLRI (Network Layer Reachability Information) of that route. See sections <u>4.1</u>, 4.2, and 4.3 of [<u>RFC6514</u>] for the specification of these NLRIS. Note that the IR P-tunnel identifier includes the "route type" and "length" octets of the NLRI.

To reiterate:

The identifier of the IR P-tunnel does not appear in the PTA at all; the "Tunnel Identifier" field of the PTA does not contain the identifier of the IR P-tunnel.

Rather, the identifier of the IR P-tunnel appears in the "Network Layer Reachability Information" (NLRI) field of the A-D routes that are used to advertise and to setup the IR P-tunnel.

Note that an advertising A-D route is considered to identify an IR P-tunnel only if it carries a PTA whose "Tunnel Type" field is set to "IR".

When an IR P-tunnel is identified in an S-PMSI A-D route or in an Inter-AS I-PMSI A-D route, the "Leaf Info Required" bit of the Flags field of the PTA MUST be set.

In an advertising A-D route:

o If the "Leaf Info Required" bit of the Flags field of the PTA is set, then the "Tunnel Identifier" field of the PTA has no significance whatsoever, and MUST be ignored upon reception.

Note that, per <u>RFC6514</u>, the length of the "Tunnel Identifier" field of the PTA is variable, and is inferred from the length of the PTA. Even when this field is of no significance, its length MUST be the length of an IP address in the address space of the SP's backbone, as specified in <u>section 4.2 of [RFC6515]</u>. In this case, it is RECOMMENDED that it be set to a routable address of the router that constructed the PTA. (While it might make more sense to allow or even require the field to be omitted entirely, that might raise issues of backwards compatibility with implementations that were designed prior to the publication of this document.)

o If the "Leaf Info Required" bit is not set, the "Tunnel Identifier" field of the PTA does have significance, but it does not identify the IR P-tunnel. The use of the PTA's "Tunnel Identifier" field in this case is discussed in <u>Section 5</u> of this document.

Note that according to the above definition, there is no way for two different advertising A-D routes (i.e., two advertising A-D routes with different NLRIS) to advertise the same IR P-tunnel. In the terminology of [RFC6513], an IR P-tunnel can instantiate only a single PMSI. If an ingress PE, for example, wants to bind two customer multicast flows to a single IR P-tunnel, it must advertise that IR P-tunnel either in an I-PMSI A-D route or in an S-PMSI A-D route whose NLRI contains wildcards ([RFC6625]).

When an IR P-tunnel is identified in a Leaf A-D route, its identifier is the "route key" field of the route's NLRI. See <u>section 4.4 of</u> [RFC6514].

A Leaf A-D route is considered to identify an IR P-tunnel only if it carries a PTA whose "Tunnel Type" field is set to "IR". In this type of route, the "Tunnel Identifier" field of the PTA does have significance, but it does not identify the IR P-tunnel. The use of the PTA's "Tunnel Identifier" field in this case is discussed in <u>Section 5</u>.

<u>4</u>. How to Join an IR P-tunnel

The procedures for joining an IR P-tunnel depend upon whether the P-tunnel has been previously advertised, and if so, upon how the P-tunnel was advertised. Note that joining an unadvertised IR P-tunnel is only possible when using the "Global Table Multicast" procedures of [RFC7524].

4.1. Advertised IR P-tunnels

The procedures in this section apply when the IR P-tunnel to be joined has been advertised in an S-PMSI A-D route, an Inter-AS I-PMSI A-D route, or an Intra-AS I-PMSI A-D route.

The procedures for joining an advertised IR P-tunnel depend upon whether the A-D route that advertises the IR P-tunnel has the "Leaf Info Required" bit set in its PTA.

4.1.1. If the 'Leaf Info Required Bit' is Set

The procedures in this section apply when the P-tunnel to be joined has been advertised in a route whose PTA has the "Leaf Info Required Bit" set.

The router joining a particular IR P-tunnel must determine its UMH for that P-tunnel. If the route that advertised the IR P-tunnel contains a P2MP Segmented Next Hop Extended Community, the UMH is determined from the value of this community (see [RFC7524]). Otherwise the UMH is determined from the route's next hop (see [RFC6514]).

Once the UMH is determined, the router joining the IR P-tunnel originates a Leaf A-D route. The NLRI of the Leaf A-D route is formed following the procedures of [RFC6514]. As a result, the NLRI of the Leaf A-D route will contain the IR P-tunnel identifier defined in <u>Section 3</u> above as its "route key". The UMH MUST be identified by attaching an "IP Address Specific Route Target" (or an "IPv6 Address Specific Route Target") to the Leaf A-D route. The IP address of the UMH appears in the "global administrator" field of the Route Target (RT). Details can be found in [RFC6514] and [RFC7524].

The Leaf A-D route MUST also contain a PTA whose fields are set as follows:

- o The "Tunnel Type" field is set to "IR".
- o The "Tunnel Identifier" field is set as described in <u>Section 5</u> of this document. (Note that this field does not contain the IR P-tunnel Identifier that is defined in <u>Section 3</u>.)
- o The "MPLS Label" field is set to a non-zero value. This is the "P-tunnel label". The value must be chosen so as to satisfy various constraints, as discussed in <u>Section 7</u> this document.

4.1.2. If the 'Leaf Info Required Bit' is Not Set

The procedures in this section apply when the IR P-tunnel to be joined has been advertised in a route whose PTA does not have the "Leaf Info Required Bit" set. This can only be the case if the IR P-tunnel was advertised in an Intra-AS I-PMSI A-D route.

If an IR P-tunnel is advertised in the Intra-AS I-PMSI A-D routes originated by the PE routers of a given MVPN, the Intra-AS I-PMSI can be thought of as being instantiated by a set of IR P-tunnels. Each PE is the root of one such IR P-tunnel, and the other PEs are children of the root. A PE simultaneously joins all these P-tunnels

by originating (if it hasn't already done so) an Intra-AS I-PMSI A-D route with a PTA whose fields are set as follows:

- o The "Tunnel Type" field is set to "IR".
- The "Tunnel Identifier" field is set as described in <u>Section 5</u> of this document. (Note that this field does not contain the IR P-tunnel Identifier that defined in <u>Section 3</u>.)
- o The "MPLS Label" field MUST be set to a non-zero value. This label value will be used by the child node to associate a received packet with the I-PMSI of a particular MVPN. The MPLS label allocation policy must be such as to ensure that the binding from label to I-PMSI is one-to-one.

The NLRI and the RTs of the originated I-PMSI A-D route are set as specified in [RFC6514].

4.2. Unadvertised IR P-tunnels

In [RFC7524], a procedure is defined for "Global Table Multicast", in which a P-tunnel can be joined even if the P-tunnel has not been previously advertised. See the sections of that document entitled "Leaf A-D Route for Global Table Multicast" and "Constructing the Rest of the Leaf A-D Route". The route key of the Leaf A-D route has the form of the "S-PMSI Route-Type Specific NLRI" in this case, and that should be considered to be the IR P-tunnel identifier. Note that the procedure for finding the UMH is different in this case; the UMH is the next hop of the best UMH-eligible route towards the "ingress PE". See the section of that document entitled "Determining the Upstream ABR/PE/ASBR (Upstream Node)".

5. The PTA's 'Tunnel Identifier' Field

As discussed in <u>Section 1</u>, when the "Tunnel Type" field of a PTA is set to "IR", the "Tunnel Identifier" field of that PTA does not contain the IR P-tunnel identifier. This section (<u>Section 5</u>) specifies the procedures for setting the "Tunnel Identifier" field of the PTA when the "Tunnel Type" field of the PTA is set to "IR".

If the "Tunnel Type" field of a PTA is set to "IR", its "Tunnel Identifier" field is significant only when one of the following two conditions holds:

- o The PTA is carried by a Leaf A-D route, or
- o The "Leaf Information Required" bit of the "Flags" field of the PTA is not set.

If one of these conditions holds, then the "Tunnel Identifier" field must contain a routable IP address of the originator of the route. (See [RFC6514] sections 9.2.3.2.1 and 9.2.3.4.1 for the detailed specification of the contents of this field.) This address is used by the UMH to determine the unicast tunnel that it will use in order to send data, along the IR P-tunnel identified by the route key, to the originator of the Leaf A-D route.

The means by which the unicast tunnel is determined from this IP address is outside the scope of this document. The means by which the unicast tunnel is set up and maintained is also outside the scope of this document.

<u>Section 4 of [RFC6515]</u> MUST be applied when a PTA is carried in a Leaf A-D route, and describes how to determine whether the "Tunnel Identifier" field carries an IPv4 or an IPv6 address.

If neither of the above conditions hold, then the "Tunnel Identifier" field is of no significance, and MUST be ignored upon reception.

6. A Note on IR P-tunnels and 'Discarding Packets from the Wrong PE'

<u>Section 9.1.1 of [RFC6513]</u> specifies a procedure known as "Discarding Packets from the Wrong PE". When an egress PE receives a multicast data packet, this procedure requires it to determine the packet's ingress PE.

In this document, we assume that when a packet has reached an egress PE via an IR P-tunnel, the egress PE will infer the identity of the packet's ingress PE by examining the packet's P-tunnel label.

<u>Section 7</u> specifies certain constraints on the way in which the P-tunnel label is allocated for a given P-tunnel. In general, if these constraints are followed, an egress PE will be able to infer the identity of a packet's ingress PE from the P-tunnel label, and hence will be able to apply the procedures of <u>Section 9.1.1 of</u> [<u>RFC6513</u>]. This method of identifying a packet's ingress PE works exactly the same when the unicast tunnels are IP tunnels as it does when the unicast tunnels are MPLS LSPs.

However, if the egress PE joined a particular IR P-tunnel using the procedures of <u>Section 4.1.2</u>, then when the egress PE receives a packet through that P-tunnel, it will not be able to infer the identity of the packet's ingress PE from the P-tunnel label, and thus will not be able to apply the procedures of <u>Section 9.1.1 of [RFC6513]</u>.

One might think that if a particular IR P-tunnel uses IP unicast tunnels rather than MPLS LSPs, an egress PE could identify the ingress PE by inspecting the IP source address field of the encapsulating IP header. However, there are several reasons why this procedure is not desirable:

- o When segmented P-tunnels are being used, the IP source address field of the encapsulating IP header might not contain the address of the ingress PE.
- o Even if the IP source address field of the encapsulating IP header does identify the ingress PE, there is no guarantee that the IP source address in that header is the same as the IP address used by the ingress PE for the MVPN signaling procedures.
- o To apply the procedures of <u>Section 9.1.1 of [RFC6513]</u> when extranet functionality [<u>RFC7900</u>] is supported, it is necessary to infer a packet's ingress VRF (Virtual Routing and Forwarding table), not merely its ingress PE. This can be inferred from the P-tunnel label (assuming that the label is allocated following the procedures of <u>Section 7</u>), but can not be inferred from the IP source address of the encapsulating IP header.

We therefore assume in this document that if the procedures of <u>Section 9.1.1 of [RFC6513]</u> are to be applied to packets traveling through IR P-tunnels, those procedures will be based on the P-tunnel label, even if the IR P-tunnel is using IP unicast tunnels.

This means that if an egress PE joined a particular IR P-tunnel using the procedures of <u>Section 4.1.2</u>, duplicate prevention on that IR P-tunnel requires the use of either Single Forwarder Selection ([<u>RFC6513</u>] section 9.1.2) or native PIM procedures ([<u>RFC6513</u>] section 9.1.3).

7. The PTA's 'MPLS Label' Field

When the "Tunnel Type" field of a PTA is set to "IR", the "MPLS Label" field is not always significant. It is significant only under the following conditions:

- 1. Either the PTA is being carried in a Leaf A-D route, or
- 2. the "Leaf Information Required" flag of the PTA is NOT set.

Note that the "Leaf Information Required" flag of the PTA is always set when a PTA specifying an IR P-tunnel is carried in an S-PMSI A-D route or in an Inter-AS I-PMSI A-D route; thus the "MPLS Label" field of the PTA is never significant when the PTA is carried by one of

these route types. The "MPLS Label" field is significant only when the PTA appears either in a Leaf A-D route or in an Intra-AS I-PMSI A-D route that does not have the "Leaf Information Required" bit set. In these cases, the MPLS label is the label that the originator of the route is assigning to the IR P-tunnel(s) identified by the route's NLRI. (That is, the MPLS label assigned in the PTA is what we have called the "P-tunnel label".)

In those cases where the "MPLS Label" field is not significant, it SHOULD be set to zero upon transmission and MUST be ignored upon reception.

7.1. Leaf A-D Route Originated by an Egress PE

As previously stated, when a Leaf A-D route is used to join an IR P-tunnel, the "route key" of the Leaf A-D route is the P-tunnel identifier.

We now define the notion of the "root of an IR P-tunnel".

- o If the identifier of an IR P-tunnel is of the form of an S-PMSI NLRI, the "root" of the IR P-tunnel is the router identified in the "Originating Router's IP Address" field of that NLRI.
- o If the identifier of an IR P-tunnel is of the form specified in Section "Leaf A-D Route for Global Table Multicast" of [<u>RFC7524</u>], the "root" of the IR P-tunnel is the router identified in the "Ingress PE's IP Address" field of that NLRI.
- o If the identifier of an IR P-tunnel is of the form of an Intra-AS I-PMSI NLRI, the "root" of the IR P-tunnel is the router identified in the "Originating Router's IP Address" field of that NLRI.
- o If the identifier of an IR P-tunnel is of the form of an Inter-AS I-PMSI NLRI, the "root" of the IR P-tunnel is same as the identifier of the IR P-tunnel, i.e., the combination of an RD and an AS.

Note that if an IR P-tunnel is segmented, the root of the IR P-tunnel, by this definition, is actually the root of the entire P-tunnel, not the root of the local segment. In this case, there may be segments upstream that are not themselves IR P-tunnels. However, the egress PE is aware only of the final segment of the P-tunnel, and hence considers the P-tunnel to be an IR P-tunnel.

In order to apply the procedures of <u>RFC 6513 Section 9.1.1</u> ("Discarding Packets from Wrong PE"), the following condition MUST be met by the MPLS label allocation policy:

Suppose an egress PE originates two Leaf A-D routes, each with a different route key in its NLRI, and each with a PTA specifying a "Tunnel Type" of "IR". Thus each of the Leaf A-D routes identifies a different IR P-tunnel. Suppose further that each of those IR P-tunnels has a different root. Then the egress PE MUST NOT specify the same MPLS label in both PMSI Tunnel attributes.

That is, to apply the "Discarding Packets from the Wrong PE" duplicate prevention procedures (<u>[RFC6513] section 9.1.1</u>), the same MPLS label MUST NOT be assigned to two IR P-tunnels that have different roots.

If segmented P-tunnels are in use, the above rule is necessary but not sufficient to prevent a PE from forwarding duplicate data to the CEs. For various reasons, a given egress PE or egress ABR or egress ASBR may decide to change its parent node, on a given segmented P-tunnel, from one router to another. It does this by changing the RT of the Leaf A-D route that it originated in order to join that P-tunnel. Once the RT is changed, there may be a period of time during which the old parent node and the new parent node are both sending data of the same multicast flow. To ensure that the egress node not forward duplicate data, whenever the egress node changes the RT that it attaches to a Leaf A-D route, it MUST also change the "MPLS Label" specified in the Leaf A-D route's PTA. This allows the egress router to distinguish between packets arriving on a given P-tunnel from the old parent and packets arriving on that same P-tunnel from the new parent. At any given time, a router MUST consider itself to have only a single parent node on a given P-tunnel, and MUST discard traffic that arrives on that P-tunnel from a different parent node.

If extranet functionality [RFC7900] is not implemented in a particular egress PE, or if an egress PE is provisioned with the knowledge that extranet functionality is not needed, the PE may adopt the policy of assigning a label that is unique for the ordered triple <root, parent node, egress VRF>. This will enable the egress PE to apply the duplicate prevention procedures discussed above, and to determine the VRF to which an arriving packet must be directed.

However, this policy is not sufficient to support the "Discard Packets from the Wrong P-tunnel" procedures that are specified in [<u>RFC7900</u>]. To support those procedures, the labels specified in the PTA of Leaf A-D routes originated by a given egress PE MUST be unique for the ordered triple <root, root RD, parent node>, where the "root

RD" is taken from the RD field of the IR P-tunnel identifier. (All forms of IR P-tunnel identifier contain an embedded "RD" field.) This policy is also sufficient for supporting non-extranet cases, but in some cases may result in the use of more labels than the policy of the previous paragraph.

7.2. Leaf A-D Route Originated by an Intermediate Node

When a P-tunnel is segmented, there will be "intermediate nodes", i.e., nodes that have a parent and also have children on the P-tunnel. Each intermediate node is a leaf node of an "upstream segment" and a root node of one or more "downstream segments". The intermediate node needs to set up its forwarding state so that data it receives on the upstream segment gets transmitted on the proper downstream segments.

If the upstream segment is instantiated by IR, the intermediate node will need to originate a Leaf A-D route to join that segment, and will need to allocate a downstream-assigned MPLS label to advertise in the MPLS label field of the Leaf A-D route's PTA. <u>Section 7.1</u> specifies constraints on the label allocation policy for egress PEs; this section specifies constraints on the label allocation policy for intermediate nodes.

Suppose intermediate node N originates two Leaf A-D routes, one whose route key is K1, and one whose route key is K2, where K1 != K2. The respective PTAs of these Leaf A-D routes MUST specify distinct nonzero MPLS labels, UNLESS the following conditions all hold:

- N's parent node for P-tunnel K1 is the same as N's parent node for P-tunnel K2.
- N's forwarding state is such that any packet it receives from P-tunnel K1 is forwarded to the exact same set of downstream neighbors as any packet it receives from P-tunnel K2.
- 3. For each downstream neighbor D to which N sends the packets it receives from P-tunnels K1 and K2, N's forwarding state is such that it applies the exact same encapsulation to packets it forwards from either tunnel to D. (E.g., if N uses MPLS to forward the packets to D, it pushes the exact same set of labels on packets from P-tunnel K1 as it pushes on packets from P-tunnel K2.)

Of course, N MAY always specify distinct non-zero labels in each of the Leaf A-D routes that it originates.

Note that the rules of this section apply whenever the upstream P-tunnel segment is an IR P-tunnel. These rules hold whether or not some or all of the downstream segments are other types of P-tunnels.

If the P-tunnels from N to a particular downstream neighbor D are IR P-tunnels, then condition 3 above will hold with respect to D only if the following conditions all hold as well:

- o N has received and installed a Leaf A-D route from D, whose route key is K1, and which carries an IP-address-specific RT identifying N,
- N has received and installed a Leaf A-D route from D, whose route key is K2, and which carries an IP-address-specific RT identifying N,
- o Those two Leaf A-D routes specify the same MPLS label in their respective PTAs.

7.3. Intra-AS I-PMSI A-D Route

When a router joins a set of IR P-tunnels using the procedures of <u>Section 4.1.2</u> of this document, the procedures of <u>section 9.1.1 of</u> [RFC6513] cannot be applied, no matter what the label allocation policy is. In this case, the ingress PE is the same as the UMH, but it is not possible to assign a label uniquely to a particular ingress PE or UMH. However, the label in the MPLS label field of the PTA MUST NOT appear in the MPLS label field of the PTA carried by any other route originated by the same router.

8. How A Child Node Prunes Itself from an IR P-tunnel

If a particular IR P-tunnel was joined via the procedures of <u>Section 4.1.2</u> of this document, a router can prune itself from the P-tunnel by withdrawing the Intra-AS I-PMSI A-D route it used to join the P-tunnel. This is not usually done unless the router is removing itself entirely from a particular MVPN.

The procedures in the remainder of this section apply when a router joined a particular IR P-tunnel by originating a Leaf A-D route (as described in <u>Section 4.1.1</u> or <u>Section 4.2</u> of this document).

If a router no longer has a need to receive any multicast data from a given IR P-tunnel, it may prune itself from the P-tunnel by withdrawing the Leaf A-D route it used to join the tunnel. This is done, e.g., if the router no longer needs any of the flows traveling over the P-tunnel, or if all the flows the router does need are being received over other P-tunnels.

A router that is attached to a particular IR P-tunnel via a particular parent node may determine that it needs to stay joined to that IR P-tunnel, but via a different parent node. This can happen, for example, if there is a change in the Next Hop or the P2MP Segmented Next Hop Extended Community of the S-PMSI A-D route in which that P-tunnel was advertised. In this case, the router changes the Route Target of the Leaf A-D route it used to join the IR P-tunnel, so that the Route Target now identifies the new parent node.

A parent node must notice when a child node has been pruned from a particular tree, as this will affect the parent node's multicast data state. Note that the pruning of a child node may appear to the parent node as the explicit withdrawal of a Leaf A-D route, or it may appear as a change in the Route Target of a Leaf A-D route. If the Route Target of a particular Leaf A-D route previously identified a particular parent node, but changes so that it no longer does so, the effect on the multicast state of the parent node is the same as if the Leaf A-D route had been explicitly withdrawn.

9. Parent Node Actions Upon Receiving Leaf A-D Route

These actions are detailed in [RFC6514] and [RFC7524]. Two points of clarification are made:

 o If a router R1 receives and installs a Leaf A-D route originated by router R2, R1's multicast state is affected only if the Leaf A-D route carries an "IP Address Specific RT" (or "IPv6 Address Specific RT") whose "global administrator" field identifies R1.

(This is as specified in [RFC6514] and [RFC7524].) If a Leaf A-D route's RT does not identify R1, but then changes so that it does identify R1, R1 must take the same actions it would take if the Leaf A-D route were newly received.

- o It is possible that router R1 will receive and install a Leaf A-D route originated by router R2, where:
 - * the route's RT identifies R1,
 - the route's NLRI contains a route key whose first octet indicates that it is identifying a P-tunnel advertised in an S-PMSI A-D route,
 - * R1 has neither originated nor installed any such S-PMSI A-D route.

If at some later time, R1 installs the corresponding S-PMSI A-D route, and the Leaf A-D route is still installed, and the Leaf A-D route's RT still identifies R1, then R1 MUST follow the same procedures it would have followed if the S-PMSI A-D route had been installed before the Leaf A-D route was installed. Implementers must not assume that events occur in the "usual" or "expected" order.

10. Use of Timers when Switching UMH

Consider a child node that has joined a particular IR P-tunnel via a particular UMH. To do so, it will have originated a Leaf A-D route with an RT that identifies the UMH. Suppose the child node now determines (for whatever reason) that it needs to change its UMH for that P-tunnel. It does this by:

- o modifying the RT of the Leaf A-D route, so that the RT now identifies the new parent rather than the old one, and by
- o modifying the PTA of the Leaf A-D route, changing the MPLS Label field as discussed in <u>Section 7</u>.

Note that, in accordance with the procedures of [<u>RFC6514</u>] and of <u>Section 4</u> of this document, the NLRI of the Leaf A-D route is not modified; only the RT and the PTA are changed.

It is desirable for such a "switch of UMH" to be done using a "make before break" technique, so that the old UMH does not stop transmitting packets of the given P-tunnel to the child until the new UMH has a chance to start transmitting packets of the given P-tunnel to the child. However, the control plane operation (i.e., modifying the RT and PTA of the Leaf A-D route) does not permit the child node to first join the IR P-tunnel via the new UMH, and then later prune itself from the old UMH. Rather, a single control plane operation has both effects.

Therefore, the old UMH MUST continue transmitting to the child node for a period of time after it sees the child's Leaf A-D route being withdrawn (or its RT changing to identify a different UMH). This timer (the "parent-continues" timer) SHOULD have a default value of 60 seconds, and SHOULD be configurable.

By the procedures of <u>Section 7</u>, the child node will have advertised a different label for the IR P-tunnel to the new UMH than it had advertised to the old UMH. This allows it to distinguish the packets of that IR P-tunnel transmitted by the new UMH from packets of that IR P-tunnel transmitted by the old UMH. At any given time, the child node will accept packets of that IR P-tunnel from only one parent node, and will discard packets of that IR P-tunnel that are received

from the other. To achieve "make before break" functionality, the child node needs to continue to accept packets from the old UMH for a period of time. After this period, it will discard any packets from the given IR P-tunnel that it receives from the old UMH, and will only accept such packets from the new UMH.

Once the child node modifies the RT of its Leaf A-D route, it MUST run a timer (the "switch-parents-delay" timer). This timer SHOULD default to 30 seconds, and SHOULD be configurable. The child node MUST continue to accept packets of the given IR P-tunnel from the old UMH until the timer expires. However, once the child node receives a packet of the given IR P-tunnel from the new UMH, it MAY consider the switch-parents-delay timer to have expired.

The "parent-continues" timer MUST be longer than the "switch-parentsdelay" timer. Note that both timers are specific to a given IR P-tunnel.

<u>11</u>. IANA Considerations

This document contains no actions for IANA.

12. Acknowledgments

The authors wish to thank Yakov Rekhter for his contributions to this work. We also wish to thank Huajin Jeng and Samir Saad for their contributions, and to thank Thomas Morin for pointing out (both before and after the document was written) some of the issues that needed further elaboration. We also thank Lucy Yong for her review and comments.

<u>Section 7.1</u> discusses the importance of having an MPLS label allocation policy that, when ingress replication is used, allows an egress PE to infer the identity of a received packet's ingress PE. This issue was first raised in earlier work by Xu Xiaohu.

<u>13</u>. Security Considerations

No security considerations are raised by this document beyond those already discussed in [RFC6513] and [RFC6514].

<u>14</u>. References

<u>**14.1</u>**. Normative References</u>

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, DOI 10.17487/RFC2119, March 1997, <<u>http://www.rfc-editor.org/info/rfc2119</u>>.
- [RFC6513] Rosen, E., Ed. and R. Aggarwal, Ed., "Multicast in MPLS/ BGP IP VPNs", <u>RFC 6513</u>, DOI 10.17487/RFC6513, February 2012, <<u>http://www.rfc-editor.org/info/rfc6513</u>>.
- [RFC6514] Aggarwal, R., Rosen, E., Morin, T., and Y. Rekhter, "BGP Encodings and Procedures for Multicast in MPLS/BGP IP VPNs", <u>RFC 6514</u>, DOI 10.17487/RFC6514, February 2012, <<u>http://www.rfc-editor.org/info/rfc6514</u>>.
- [RFC6515] Aggarwal, R. and E. Rosen, "IPv4 and IPv6 Infrastructure Addresses in BGP Updates for Multicast VPN", <u>RFC 6515</u>, DOI 10.17487/RFC6515, February 2012, <<u>http://www.rfc-editor.org/info/rfc6515</u>>.

<u>14.2</u>. Informative References

- [RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", <u>RFC 3209</u>, DOI 10.17487/RFC3209, December 2001, <<u>http://www.rfc-editor.org/info/rfc3209</u>>.
- [RFC5036] Andersson, L., Ed., Minei, I., Ed., and B. Thomas, Ed., "LDP Specification", <u>RFC 5036</u>, DOI 10.17487/RFC5036, October 2007, <<u>http://www.rfc-editor.org/info/rfc5036</u>>.
- [RFC6388] Wijnands, IJ., Ed., Minei, I., Ed., Kompella, K., and B. Thomas, "Label Distribution Protocol Extensions for Pointto-Multipoint and Multipoint-to-Multipoint Label Switched Paths", <u>RFC 6388</u>, DOI 10.17487/RFC6388, November 2011, <<u>http://www.rfc-editor.org/info/rfc6388></u>.
- [RFC6625] Rosen, E., Ed., Rekhter, Y., Ed., Hendrickx, W., and R. Qiu, "Wildcards in Multicast VPN Auto-Discovery Routes", <u>RFC 6625</u>, DOI 10.17487/RFC6625, May 2012, <<u>http://www.rfc-editor.org/info/rfc6625</u>>.
- [RFC7524] Rekhter, Y., Rosen, E., Aggarwal, R., Morin, T., Grosclaude, I., Leymann, N., and S. Saad, "Inter-Area Point-to-Multipoint (P2MP) Segmented Label Switched Paths (LSPs)", <u>RFC 7524</u>, DOI 10.17487/RFC7524, May 2015, <<u>http://www.rfc-editor.org/info/rfc7524</u>>.

- [RFC7582] Rosen, E., Wijnands, IJ., Cai, Y., and A. Boers, "Multicast Virtual Private Network (MVPN): Using Bidirectional P-Tunnels", <u>RFC 7582</u>, DOI 10.17487/RFC7582, July 2015, <<u>http://www.rfc-editor.org/info/rfc7582</u>>.
- [RFC7716] Zhang, J., Giuliano, L., Rosen, E., Ed., Subramanian, K., and D. Pacella, "Global Table Multicast with BGP Multicast VPN (BGP-MVPN) Procedures", <u>RFC 7716</u>, DOI 10.17487/RFC7716, December 2015, <<u>http://www.rfc-editor.org/info/rfc7716</u>>.
- [RFC7740] Zhang, Z., Rekhter, Y., and A. Dolganow, "Simulating Partial Mesh of Multipoint-to-Multipoint (MP2MP) Provider Tunnels with Ingress Replication", <u>RFC 7740</u>, DOI 10.17487/RFC7740, January 2016, <<u>http://www.rfc-editor.org/info/rfc7740</u>>.
- [RFC7900] Rekhter, Y., Ed., Rosen, E., Ed., Aggarwal, R., Cai, Y., and T. Morin, "Extranet Multicast in BGP/IP MPLS VPNs", <u>RFC 7900</u>, DOI 10.17487/RFC7900, June 2016, <<u>http://www.rfc-editor.org/info/rfc7900</u>>.

Authors' Addresses

Eric C. Rosen (editor) Juniper Networks, Inc. 10 Technology Park Drive Westford, Massachusetts 01886 United States

Email: erosen@juniper.net

Karthik Subramanian Sproute Networks

Email: karthik@sproute.com

Zhaohui Zhang Juniper Networks, Inc. 10 Technology Park Drive Westford, Massachusetts 01886 United States

Email: zzhang@juniper.net