

Internet Engineering Task Force (IETF)
Request for Comments: 8395
Updates: 4761
Category: Standards Track
ISSN: 2070-1721

K. Patel
Arrcus
S. Boutros
VMware
J. Liste
Cisco
B. Wen
Comcast
J. Rabadan
Nokia
June 2018

Extensions to BGP-Signaled Pseudowires to
Support Flow-Aware Transport Labels

Abstract

This document defines protocol extensions required to synchronize flow label states among Provider Edges (PEs) when using the BGP-based signaling procedures. These protocol extensions are equally applicable to point-to-point Layer 2 Virtual Private Networks (L2VPNs). This document updates RFC 4761 by defining new flags in the Control Flags field of the Layer2 Info Extended Community.

Status of This Memo

This is an Internet Standards Track document.

This document is a product of the Internet Engineering Task Force (IETF). It represents the consensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Further information on Internet Standards is available in Section 2 of RFC 7841.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <https://www.rfc-editor.org/info/rfc8395>.

Copyright Notice

Copyright (c) 2018 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction2
1.1. Requirements Language3
2. Modifications to the Layer2 Info Extended Community4
3. Signaling the Presence of the Flow Label5
4. IANA Considerations6
5. Security Considerations6
6. References7
6.1. Normative References7
6.2. Informative References7
Acknowledgements8
Contributors8
Authors’ Addresses9

1. Introduction

The mechanism described in [RFC6391] uses an additional label (Flow Label) in the MPLS label stack to allow Label Switching Routers (LSRs) to balance flows within Pseudowires (PWs) at a finer granularity than the individual PWs across the Equal Cost Multiple Paths (ECMPs) that exists within the Packet Switched Network (PSN).

Furthermore, [RFC6391] defines the LDP protocol extensions required to synchronize the flow label states between the ingress and egress PEs when using the signaling procedures defined in the [RFC8077].

A PW [RFC3985] is transported over one single network path, even if ECMPs exist between the ingress and egress PW provider edge (PE) equipment. This is required to preserve the characteristics of the emulated service.

This document introduces an optional mode of operation allowing a PW to be transported over ECMPs, for example when the use of ECMPs is known to be beneficial to the operation of the PW. This specification uses the principles defined in [RFC6391] and augments the BGP-signaling procedures of [RFC4761] and [RFC6624]. The use of a single path to preserve the packet delivery order remains the default mode of operation of a PW and is described in [RFC4385] and [RFC4928].

High-bandwidth Ethernet-based services are a prime example that use of the optional mode benefits from the ability to load-balance flows in a PW over multiple PSN paths. In general, load-balancing is applicable when the PW attachment circuit bandwidth and PSN core link bandwidth are of the same order of magnitude.

To achieve the load-balancing goal, [RFC6391] introduces the notion of an additional Label Stack Entry (LSE) (flow label) located at the bottom of the stack (right after PW LSE). LSRs commonly generate a hash of the label stack in order to discriminate and distribute flows over available ECMPs. The presence of the flow label (closely associated to a flow determined by the ingress PE) will normally provide the greatest entropy.

Furthermore, following the procedures for inter-AS scenarios described in Section 3.4 of [RFC4761], the flow label should never be handled by the ASBRs; only the terminating PEs on each AS will be responsible for popping or pushing this label. This is equally applicable to Method B as described in Section 3.4.2 of [RFC4761], where ASBRs are responsible for swapping the PW label as traffic traverses from ASBR to PE and ASBR to ASBR. Therefore, the flow label will remain untouched across AS boundaries.

1.1. Requirements Language

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

2. Modifications to the Layer2 Info Extended Community

The Layer2 Info Extended Community is used to signal control information about the PWs to be set up. The Extended Community format is described in [RFC4761]. The format of this Extended Community is described as:

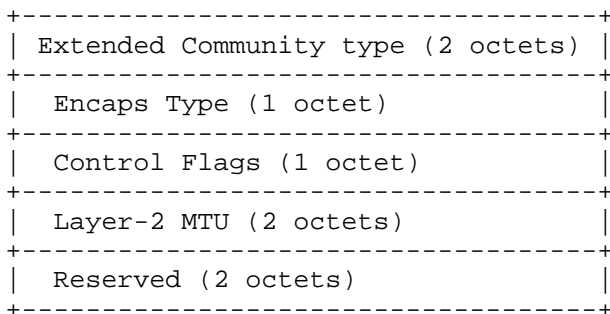


Figure 1: Layer2 Info Extended Community

Control Flags:

This field contains bit flags relating to the control information about PWs. This field is augmented with a definition of two new flags fields.

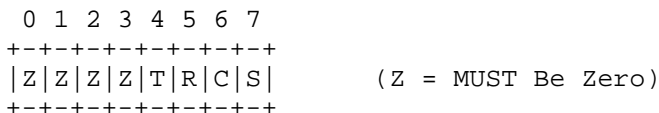


Figure 2: Control Flags Bit Vector

With reference to the Control Flags Bit Vector, the following bits in the Control Flags are defined. The remaining bits, designated "Z", MUST be set to zero when sending and MUST be ignored when receiving this Extended Community.

- T When the bit value is 1, the PE announces the ability to send a PW packet that includes a flow label. When the bit value is 0, the PE is indicating that it will not send a PW packet containing a flow label.
- R When the bit value is 1, the PE is able to receive a PW packet with a flow label present. When the bit value is 0, the PE is unable to receive a PW packet with the flow label present.

C Defined in [RFC4761].

S Defined in [RFC4761].

3. Signaling the Presence of the Flow Label

As part of the PW signaling procedures described in [RFC4761], a Layer2 Info Extended Community is advertised in the Virtual Private LAN Service (VPLS) BGP Network Layer Reachability Information (NLRI).

A PE that wishes to send a flow label in a PW packet MUST include in its VPLS BGP NLRI a Layer2 Info Extended Community using Control Flags field with T = 1.

A PE that is willing to receive a flow label in a PW packet MUST include in its VPLS BGP NLRI a Layer2 Info Extended Community using Control Flags field with R = 1.

A PE that receives a VPLS BGP NLRI containing a Layer2 Info Extended Community with R = 0 MUST NOT include a flow label in the PW packet.

Therefore, a PE sending a Control Flags field with T = 1 and receiving a Control Flags field with R = 1 MUST include a flow label in the PW packet. With any other combination, a PE MUST NOT include a flow label in the PW packet.

A PE MAY support the configuration of the flow label (T and R bits) on a per-service basis (e.g., a VPLS VPN Forwarding Instance (VFI)). Furthermore, it is also possible that on a given service, PEs may not share the same flow label settings. The presence of a flow label is therefore determined on a per-peer basis and according to the local and remote T and R bit values. For example, a PE part of a VPLS and with a local T = 1 must only transmit traffic with a flow label to those peers that signaled R = 1. If the same PE has local R = 1, it must only expect to receive traffic with a flow label from peers with T = 1. Any other traffic must not have a flow label. A PE expecting to receive traffic from a remote peer with a flow label MAY drop traffic that has no flow label. A PE expecting to receive traffic from a remote peer with no flow label MAY drop traffic that has a flow label.

Modification of flow label settings may impact traffic over a PW, as these could trigger changes in the PEs data-plane programming (i.e., imposition/disposition of the flow label). This is an implementation-specific behavior and is outside the scope of this document.

The signaling procedures in [RFC4761] state that the unspecified bits in the Control Flags field (bits 0-5) MUST be set to zero when sending and MUST be ignored when receiving. The signaling procedure described here is therefore backwards compatible with existing implementations. A PE not supporting the extensions described in this document will always advertise a value of zero in the R bit; therefore, a flow label will never be included in a packet sent to it by one of its peers. Similarly, it will always advertise a value of zero in the T bit; therefore, a peer will know that a flow label will never be included in a packet sent by it.

Note that what is signaled is the desire to include the flow LSE in the label stack. The value of the flow label is a local matter for the ingress PE, and the label value itself is not signaled.

4. IANA Considerations

Although [RFC4761] defined a Control Flags Bit Vector as part of the Layer2 Info Extended Community, it did not ask for the creation of a registry.

Per this document, IANA has created the "Layer2 Info Extended Community Control Flags Bit Vector" registry
<<https://www.iana.org/assignments/bgp-extended-communities>>.

Based on [RFC4761] and this document, the initial contents of this registry are as follows:

Value	Name	Reference
T	Request to send a flow label	This document
R	Ability to receive a flow label	This document
C	Presence of a Control Word	RFC 4761
S	Sequenced delivery of frames	RFC 4761

As per [RFC4761] and this document, the remaining bits are unassigned, and MUST be set to zero when sending and MUST be ignored when receiving the Layer2 Info Extended Community.

5. Security Considerations

This extension to BGP does not change the underlying security issues inherent in [RFC4271] and [RFC4761].

6. References

6.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC4761] Kompella, K., Ed., and Y. Rekhter, Ed., "Virtual Private LAN Service (VPLS) Using BGP for Auto-Discovery and Signaling", RFC 4761, DOI 10.17487/RFC4761, January 2007, <<https://www.rfc-editor.org/info/rfc4761>>.
- [RFC6391] Bryant, S., Ed., Filss, C., Drafz, U., Kompella, V., Regan, J., and S. Amante, "Flow-Aware Transport of Pseudowires over an MPLS Packet Switched Network", RFC 6391, DOI 10.17487/RFC6391, November 2011, <<https://www.rfc-editor.org/info/rfc6391>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

6.2. Informative References

- [RFC3985] Bryant, S., Ed., and P. Pate, Ed., "Pseudo Wire Emulation Edge-to-Edge (PWE3) Architecture", RFC 3985, DOI 10.17487/RFC3985, March 2005, <<https://www.rfc-editor.org/info/rfc3985>>.
- [RFC4385] Bryant, S., Swallow, G., Martini, L., and D. McPherson, "Pseudowire Emulation Edge-to-Edge (PWE3) Control Word for Use over an MPLS PSN", RFC 4385, DOI 10.17487/RFC4385, February 2006, <<https://www.rfc-editor.org/info/rfc4385>>.
- [RFC8077] Martini, L., Ed., and G. Heron, Ed., "Pseudowire Setup and Maintenance Using the Label Distribution Protocol (LDP)", STD 84, RFC 8077, DOI 10.17487/RFC8077, February 2017, <<https://www.rfc-editor.org/info/rfc8077>>.

- [RFC4928] Swallow, G., Bryant, S., and L. Andersson, "Avoiding Equal Cost Multipath Treatment in MPLS Networks", BCP 128, RFC 4928, DOI 10.17487/RFC4928, June 2007, <<https://www.rfc-editor.org/info/rfc4928>>.
- [RFC6624] Kompella, K., Kothari, B., and R. Cherukuri, "Layer 2 Virtual Private Networks Using BGP for Auto-Discovery and Signaling", RFC 6624, DOI 10.17487/RFC6624, May 2012, <<https://www.rfc-editor.org/info/rfc6624>>.

Acknowledgements

The authors would like to thank Bertrand Duvivier and John Drake for their review and comments.

Contributors

In addition to the authors listed above, the following individuals also contributed to this document:

Eric Lent

John Brzozowski

Steven Cotter

Authors' Addresses

Keyur Patel
Arrcus

Email: keyur@arrcus.com

Sami Boutros
VMware

Email: boutros.sami@gmail.com

Jose Liste
Cisco

Email: jliste@cisco.com

Bin Wen
Comcast

Email: bin_wen@cable.comcast.com

Jorge Rabadan
Nokia

Email: jorge.rabadan@nokia.com