

Graceful Restart Mechanism for BGP with MPLS

[draft-ietf-mpls-bgp-mpls-restart-02.txt](#)

Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of [Section 10 of RFC2026](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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.''

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

Abstract

A mechanism for BGP that would help minimize the negative effects on routing caused by BGP restart is described in "Graceful Restart Mechanism for BGP" (see [1]). This document extends this mechanism to also minimize the negative effects on MPLS forwarding caused by the Label Switching Router's (LSR's) control plane restart, and specifically by the restart of its BGP component when BGP is used to carry MPLS labels and the LSR is capable of preserving the MPLS forwarding state across the restart.

The mechanism described in this document is agnostic with respect to the types of the addresses carried in the BGP Network Layer Reachability Information (NLRI) field. As such it works in conjunction with any of the address families that could be carried in BGP (e.g., IPv4, IPv6, etc...)

The mechanism described in this document is applicable to all LSRs, both those with the ability to preserve their forwarding state during BGP restart and those without (although the latter need to implement only a subset of the mechanism described in this document).

Supporting (a subset of) the mechanism described here by the LSRs that can not preserve their MPLS forwarding state across the restart would not reduce the negative impact on MPLS traffic caused by their control plane restart, but it would minimize the impact if their neighbor(s) are capable of preserving the forwarding state across the restart of their control plane and implement the mechanism described here.

The mechanism makes minimalistic assumptions on what has to be preserved across restart - the mechanism assumes that only the actual MPLS forwarding state has to be preserved; the mechanism does not require any of the BGP-related state to be preserved across the restart.

Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [[RFC2119](#)].

Summary for Sub-IP Area

(This section to be removed before publication.)

[0.1](#) Summary

This document describes a mechanism that helps to minimize the negative effects on MPLS forwarding caused by LSR's control plane restart, and specifically by the restart of its BGP component in the case where BGP is used to carry MPLS labels and LSR is capable of preserving its MPLS forwarding state across the restart.

0.2 Related documents

See the Reference Section

0.3 Where does it fit in the Picture of the Sub-IP Work

This work fits squarely in MPLS box.

0.4 Why is it Targeted at this WG

The specifications on carrying MPLS Labels in BGP is a product of the MPLS WG. This document specifies procedures to minimize the negative effects on MPLS forwarding caused by the restart of the control plane BGP module in the case where BGP is used to carry MPLS labels. Since the procedures described in this document are directly related to MPLS forwarding and carrying MPLS labels in BGP, it would be logical to target this document at the MPLS WG.

0.5 Justification

The WG should consider this document, as it allows to minimize the negative effects on MPLS forwarding caused by the restart of the control plane BGP module in the case where BGP is used to carry MPLS labels.

1. Motivation

For the sake of brevity in the context of this document by "MPLS forwarding state" we mean either <incoming label -> (outgoing label, next hop)>, or <address prefix -> (outgoing label, next hop)> mapping. In the context of this document the forwarding state that is referred to in [1] means MPLS forwarding state.

In the case where a Label Switching Router (LSR) could preserve its MPLS forwarding state across restart of its control plane, and specifically its BGP component, and BGP is used to carry MPLS labels (as specified in [2]), it may be desirable not to perturb the LSPs going through that LSR (and specifically, the LSPs established by BGP). In this document, we describe a mechanism that allows to accomplish this goal. The mechanism described in this document works in conjunction with the mechanism specified in [1]. The mechanism described in this document places no restrictions on the types of addresses (address families) that it can support.

The mechanism described in this document is applicable to all LSRs, both those with the ability to preserve forwarding state during BGP restart and those without (although the latter need to implement only a subset of the mechanism described in this document). Supporting (a subset of) the mechanism described here by the LSRs that can not preserve their MPLS forwarding state across the restart would not reduce the negative impact on MPLS traffic caused by their control plane restart, but it would minimize the impact if their neighbor(s) are capable of preserving the forwarding state across the restart of their control plane and implement the mechanism described here.

2. Assumptions

First of all we assume that an LSR implements the Graceful Restart Mechanism for BGP, as specified in [1]. Second, we assume that the LSR is capable of preserving its MPLS forwarding state across the restart of its control plane (including the restart of BGP).

The mechanism makes minimalistic assumptions on what has to be preserved across restart - the mechanism assumes that only the actual MPLS forwarding state has to be preserved; the mechanism does not require any of the BGP-related state to be preserved across the restart.

In the scenario where label binding on an LSR is created/maintained not just by the BGP component of the control plane, but by other protocol components as well (e.g., LDP, RSVP-TE), and the LSR supports restart of the individual components of the control plane that create/maintain label binding (e.g., restart of BGP, but no restart of LDP) the LSR needs to preserve across the restart the information about which protocol has assigned which labels.

3. Capability Advertisement

An LSR that supports the mechanism described in this document advertises this to its peer by using the Graceful Restart Capability, as specified in [1]. The Subsequent Address Family Identifier (SAFI) in the advertised capability MUST indicate that the Network Layer Reachability Information (NLRI) field carries not just addressing information but labels as well (see [2]).

4. Procedures for the restarting LSR

After the LSR restarts, it follows the procedures as specified in [1]. In addition, if the LSR is able to preserve its MPLS forwarding state across the restart, the LSR advertises this to its neighbors by appropriately setting the Flag field in the Graceful Restart Capability for all applicable AFI/SAFI pairs.

Once the restarting LSR completes its route selection (as specified in Section "Procedures for the Restarting Speaker" of [1]), then in addition to the procedures specified in [1], the restarting LSR performs one of the following:

4.1. Case 1

The following applies when (a) the best route selected by the restarting LSR was received with a label, (b) that label is not an Implicit NULL, and (c) the LSR advertises this route with itself as the next hop.

In this case the restarting LSR searches its MPLS forwarding state (the one preserved across the restart) for an entry with <outgoing label, Next-Hop> equal to the one in the received route. If such an entry is found, the LSR no longer marks the entry as stale. In addition if the entry is of type <incoming label, (outgoing label, next hop)> rather than <prefix, (outgoing label, next hop)>, the LSR uses the incoming label from the entry when advertising the route to its neighbors. If the found entry has no incoming label, or if no such entry is found, the LSR just picks up some unused label when advertising the route to its neighbors (assuming that there are neighbors to which the LSR has to advertise the route with a label).

4.2. Case 2

The following applies when (a) the best route selected by the restarting LSR was received either without a label, or with an Implicit NULL label, or the route is originated by the restarting LSR, (b) the LSR advertises this route with itself as the next hop, and (c) the LSR has to generate a (non Implicit NULL) label for the route.

In this case the LSR searches its MPLS forwarding state for an entry that indicates that the LSR has to perform label pop, and the next hop equal to the next hop of the route in consideration. If such an entry is found, then the LSR uses the incoming label from the entry when advertising the route to its neighbors. If no such entry is

found, the LSR just picks up some unused label when advertising the route to its neighbors.

The description in the above paragraph assumes that the restarting LSR generates the same label for all the routes with the same next hop. If this is not the case, and the restarting LSR generates a unique label per each such route, then the LSR needs to preserve across the restart not just <incoming label, (outgoing label, next hop)> mapping, but also the prefix associated with this mapping. In such case the LSR would search its MPLS forwarding state for an entry that (a) indicates Label pop (means no outgoing label), (b) the next hop equal to the next hop of the route and (c) has the same prefix as the route. If such an entry is found, then the LSR uses the incoming label from the entry when advertising the route to its neighbors. If no such entry is found, the LSR just picks up some unused label when advertising the route to its neighbors.

4.3. Case 3

The following applies when the restarting LSR does not set BGP Next Hop to self.

In this case the restarting LSR, when advertising its best route for a particular NLRI just uses the label that was received with that route. And if the route was received with no label, the LSR advertises the route with no label as well.

5. Alternative procedures for the restarting LSR

In this section we describe an alternative to the procedures described in Section "Procedures for the restarting LSR".

The procedures described in this section assume that the restarting LSR has (at least) as many unallocated as allocated labels. The latter forms the MPLS forwarding state that the LSR managed to preserve across the restart. The former is used for allocating labels after the restart.

After the LSR restarts, it follows the procedures as specified in [1]. In addition, if the LSR is able to preserve its MPLS forwarding state across the restart, the LSR advertises this to its neighbors by appropriately setting the Flag field in the Graceful Restart Capability.

To create local label bindings the LSR uses unallocated labels (this is pretty much the normal procedure). That means that as long as the

LSR retains the MPLS forwarding state that the LSR preserved across the restart, the labels from that state are not used for creating local label bindings.

The restarting LSR SHOULD retain the MPLS forwarding state that the LSR preserved across the restart at least until the LSR sends End-of-RIB marker to all of its neighbors (by that time the LSR already completed its route selection process, and also advertised its Adj-RIB-Out to its neighbors). The restarting LSR MAY retain the forwarding state even a bit longer, as to allow the neighbors to receive and process the routes that have been advertised by the restarting LSR. After that, the restarting LSR MAY delete the MPLS forwarding state that it preserved across the restart.

Note that while an LSR is in the process of restarting, the LSR may have not one, but two local label bindings for a given BGP route - one that was retained from prior to restart, and another that was created after the restart. Once the LSR completes its restart, the former will be deleted. Both of these bindings though would have the same outgoing label (and the same next hop).

6. Procedures for a neighbor of a restarting LSR

The neighbor of a restarting LSR (the receiving router in terminology used in [1]) follows the procedures specified in [1]. In addition, the neighbor treats the MPLS labels received from the restarting LSR the same way as it treats the routes received from the restarting LSR (both prior and after the restart).

Replacing the stale routes by the routing updates received from the restarting LSR involves replacing/updating the appropriate MPLS labels.

In addition, if the Flags in the Graceful Restart Capability received from the restarting LSR indicate that the LSR wasn't able to retain its MPLS state across the restart, the neighbor SHOULD immediately remove all the NLRI and the associated MPLS labels that it previously acquired via BGP from the restarting LSR.

An LSR, once it creates a binding between a label and a Forwarding Equivalence Class (FEC), SHOULD keep the value of the label in this binding for as long as the LSR has a route to the FEC in the binding. If the route to the FEC disappears, and then re-appears again later, then this may result in using a different label value, as when the route re-appears, the LSR would create a new <label, FEC> binding.

To minimize the potential mis-routing caused by the label change,

when creating a new <label, FEC> binding the LSR SHOULD pick up the least recently used label. Once an LSR releases a label, the LSR SHALL NOT re-use this label for advertising a <label, FEC> binding to a neighbor that supports graceful restart for at least the Restart Time, as advertised by the neighbor to the LSR.

7. Security Consideration

The security considerations pertaining to the original BGP protocol remain relevant.

In addition, the mechanism described here renders LSRs that implement it to additional denial-of-service attacks as follows:

An intruder may impersonate a BGP peer in order to force a failure and reconnection of the TCP connection, but where the intruder sets the Forwarding State (F) bit (as defined in [1]) to 0 on reconnection. This forces all labels received from the peer to be released.

An intruder could intercept the traffic between BGP peers and override the setting of the Forwarding State (F) bit to be set to 0. This forces all labels received from the peer to be released.

All of these attacks may be countered by use of an authentication scheme between BGP peers, such as the scheme outlined in [RFC2385].

As with BGP carrying labels, a security issue may exist if a BGP implementation continues to use labels after expiration of the BGP session that first caused them to be used. This may arise if the upstream LSR detects the session failure after the downstream LSR has released and re-used the label. The problem is most obvious with the platform-wide label space and could result in mis-routing of data to other than intended destinations and it is conceivable that these behaviors may be deliberately exploited to either obtain services without authorization or to deny services to others.

In this document, the validity of the BGP session may be extended by the Restart Time, and the session may be re-established in this period. After the expiry of the Restart Time the session must be considered to have failed and the same security issue applies as described above.

However, the downstream LSR may declare the session as failed before the expiration of its Restart Time. This increases the period during which the downstream LSR might reallocate the label while the upstream LSR continues to transmit data using the old usage of the

label. To reduce this issue, this document requires that labels are not re-used until for at least the Restart Time.

8. Intellectual Property Considerations

This section is taken from [Section 10.4 of \[RFC2026\]](#).

The IETF takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Information on the IETF's procedures with respect to rights in standards-track and standards-related documentation can be found in [BCP-11](#). Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementors or users of this specification can be obtained from the IETF Secretariat.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this standard. Please address the information to the IETF Executive Director.

The IETF has been notified of intellectual property rights claimed in regard to some or all of the specification contained in this document. For more information consult the online list of claimed rights.

9. Copyright Notice

Copyright (C) The Internet Society (date). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of

developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

10. Acknowledgments

We would like to thank Chaitanya Kodeboyina and Loa Andersson for their review and comments. The approach described in Section "Alternative procedures for the restarting LSR" is based on the idea suggested by Manoj Leelanivas.

11. Normative References

- [1] "Graceful Restart Mechanism for BGP", [draft-ietf-idr-restart-01.txt](#)
- [2] Rekhter, Y., Rosen, E., "Carrying Label Information in BGP-4", [RFC3107](#)
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#)
- [RFC2385] Heffernan, A., "Protection of BGP Sessions via the TCP MD5 Signature Option", [RFC2385](#)
- [RFC2026] Bradner, S., "The Internet Standards Process -- Revision 3", [RFC2026](#)

12. Author Information

Yakov Rekhter
Juniper Networks
1194 N.Mathilda Ave
Sunnyvale, CA 94089
e-mail: yakov@juniper.net

Rahul Aggarwal
Redback Networks
350 Holger Way
San Jose, CA 95134
e-mail: rahul@redback.com

