

TEAS Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 3, 2018

V. Beeram, Ed.
Juniper Networks
I. Minei
R. Shakir
Google, Inc
D. Pacella
Verizon
T. Saad
Cisco Systems
July 2, 2017

**Implementation Recommendations to Improve the Scalability of RSVP-TE
Deployments
draft-ietf-teas-rsvp-te-scaling-rec-05**

Abstract

The scale at which RSVP-TE Label Switched Paths (LSPs) get deployed is growing continually and the onus is on RSVP-TE implementations across the board to keep up with this increasing demand.

This document introduces a couple of techniques - "Refresh-Interval Independent RSVP (RI-RSVP)" and "Per-Peer Flow-Control" - to help RSVP-TE deployments push the envelope on scaling.

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 3, 2018.

Copyright Notice

Copyright (c) 2017 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	2
1.1.	Conventions used in this document	3
2.	Recommendations	3
2.1.	"RFC2961 Specific" Recommendations	3
2.1.1.	Basic Prerequisites	3
2.1.2.	Making Acknowledgements Mandatory	4
2.1.3.	Clarifications On Reaching Rapid Retry Limit (Rl)	4
2.2.	Refresh-Interval Independent RSVP	5
2.2.1.	Capability Advertisement	5
2.2.2.	Compatibility	6
2.3.	Per-Peer RSVP Flow-Control	6
2.3.1.	Capability Advertisement	7
2.3.2.	Compatibility	7
3.	Acknowledgements	7
4.	Contributors	7
5.	IANA Considerations	8
5.1.	Capability Object Values	8
6.	Security Considerations	8
7.	References	8
7.1.	Normative References	8
7.2.	Informative References	9
Appendix A.	Recommended Defaults	9
	Authors' Addresses	10

[1.](#) Introduction

The scale at which RSVP-TE [[RFC3209](#)] Label Switched Paths (LSPs) get deployed is growing continually and there is considerable onus on RSVP-TE implementations across the board to keep up with this increasing demand in scale.

The set of RSVP Refresh Overhead Reduction procedures [[RFC2961](#)] serves as a powerful toolkit for RSVP-TE implementations to help cover a majority of the concerns about soft-state scaling. However, even with these tools in the toolkit, analysis of existing

implementations [[RFC5439](#)] indicates that the processing required under certain scale may still cause significant disruption to an LSR.

This document builds on the scaling work and analysis that has been done so far and makes a set of concrete implementation recommendations to help RSVP-TE deployments push the envelope further on scaling - push higher the threshold above which an LSR struggles to achieve sufficient processing to maintain LSP state.

This document advocates the use of a couple of techniques - "Refresh-Interval Independent RSVP (RI-RSVP)" and "Per-Peer Flow-Control" - for significantly cutting down the amount of processing cycles required to maintain LSP state. "RI-RSVP" helps completely eliminate RSVP's reliance on refreshes and refresh-timeouts while "Per-Peer Flow-Control" enables a busy RSVP speaker to apply back pressure to its peer(s). In order to reap maximum scaling benefits, it is strongly RECOMMENDED that implementations support both the techniques.

1.1. Conventions used in this document

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](#)]

2. Recommendations

2.1. "[RFC2961](#) Specific" Recommendations

The implementation recommendations discussed in this section are based on the proposals made in [[RFC2961](#)] and act as prerequisites for implementing the techniques discussed in Sections [2.2](#) and [2.3](#).

2.1.1. Basic Prerequisites

An implementation that supports the techniques discussed in Sections 2.2 and 2.3 must meet certain basic prerequisites.

- o It MUST indicate support for RSVP Refresh Overhead Reduction extensions (as specified in [Section 2 of \[RFC2961\]\(#\)](#)).
- o It MUST support receipt of any RSVP Refresh Overhead Reduction message as defined in [[RFC2961](#)].
- o It SHOULD initiate all RSVP Refresh Overhead Reduction mechanisms as defined in [[RFC2961](#)] (including the SRefresh message) with the default behavior being to initiate the mechanisms but offering a configuration override.

- o It MUST support reliable delivery of Path/Resv and the corresponding Tear/Err messages (as specified in [Section 4 of \[RFC2961\]](#)).
- o It MUST support retransmission of all unacknowledged RSVP-TE messages using exponential-backoff (as specified in [Section 6 of \[RFC2961\]](#)).

[2.1.2.](#) Making Acknowledgements Mandatory

The reliable message delivery mechanism specified in [\[RFC2961\]](#) states that "Nodes receiving a non-out of order message containing a MESSAGE_ID object with the ACK_Desired flag set, SHOULD respond with a MESSAGE_ID_ACK object."

In an implementation that supports the techniques discussed in Sections [2.2](#) and [2.3](#), nodes receiving a non-out of order message containing a MESSAGE ID object with the ACK-Desired flag set, MUST respond with a MESSAGE_ID_ACK object. This MESSAGE_ID_ACK object can be packed along with other MESSAGE_ID_ACK or MESSAGE_ID_NACK objects and sent in an Ack message (or piggy-backed in any other RSVP message). This improvement to the predictability of the system in terms of reliable message delivery is key for being able to take any action based on a non-receipt of an ACK.

[2.1.3.](#) Clarifications On Reaching Rapid Retry Limit (Rl)

According to [section 6 of \[RFC2961\]](#) "The staged retransmission will continue until either an appropriate MESSAGE_ID_ACK object is received, or the rapid retry limit, Rl, has been reached." The following clarifies what actions, if any, a router should take once Rl has been reached.

If Rl has been reached for the retransmission of a message that is neither a Path nor a Resv message, then the router need not take any further action. If Rl has been reached for the retransmission of a Path or a Resv message, then the router starts periodic retransmission of these on a slower timer. The retransmitted messages MUST carry MESSAGE_ID object with ACK_Desired flag set. This periodic retransmission SHOULD continue until an appropriate MESSAGE_ID_ACK object is received indicating acknowledgement of the (retransmitted) Path/Resv message. The configurable periodic retransmission interval for this slower timer SHOULD be less than the regular refresh interval. A default periodic retransmission interval (for this slower timer) of 30 seconds is RECOMMENDED by this document.

2.2. Refresh-Interval Independent RSVP

The RSVP protocol relies on periodic refreshes for state synchronization between RSVP neighbors and for recovery from lost RSVP messages. It relies on refresh timeout for stale state cleanup. The primary motivation behind introducing the notion of "Refresh Interval Independent RSVP" (RI-RSVP) is to completely eliminate RSVP's reliance on refreshes and refresh timeouts. This is done by simply increasing the refresh interval to a fairly large value. [\[RFC2961\]](#) and [\[RFC5439\]](#) do talk about increasing the value of the refresh-interval to provide linear improvement on transmission overhead, but also point out the degree of functionality that is lost by doing so. This section revisits this notion, but also proposes sufficient recommendations to make sure that there is no loss of functionality incurred by increasing the value of the refresh interval.

An implementation that supports RI-RSVP:

- o MUST support all the recommendations made in [Section 2.1](#).
- o MUST make the default value of the configurable refresh interval be a large value (10s of minutes). A default value of 20 minutes is RECOMMENDED by this document.
- o MUST implement coupling the state of individual LSPs with the state of the corresponding RSVP-TE signaling adjacency. When an RSVP-TE speaker detects RSVP-TE signaling adjacency failure, the speaker MUST act as if the all the Path and Resv state learnt via the failed signaling adjacency has timed out.
- o MUST make use of Node-ID based Hello Session ([\[RFC3209\]](#), [\[RFC4558\]](#)) for detection of RSVP-TE signaling adjacency failures; A default value of 9 seconds is RECOMMENDED by this document for the configurable node hello interval (as opposed to the 5ms default value proposed in [Section 5.3 of \[RFC3209\]](#)).
- o MUST indicate support for RI-RSVP via the CAPABILITY object in Hello messages.

2.2.1. Capability Advertisement

An implementation supporting the RI-RSVP recommendations MUST set a new flag "RI-RSVP Capable" in the CAPABILITY object signaled in Hello messages.

Bit Number TBA1 (TBA2) - RI-RSVP Capable (I-bit):

Indicates that the sender supports RI-RSVP.

Any node that sets the new I-bit in its CAPABILITY object MUST also set Refresh-Reduction-Capable bit in common header of all RSVP-TE messages. If a peer sets the I-bit in the CAPABILITY object but does not set the Refresh-Reduction-Capable bit, then the RI-RSVP functionality MUST NOT be activated for that peer.

2.2.2. Compatibility

The RI-RSVP functionality MUST NOT be activated with a peer that does not indicate support for this functionality.

2.3. Per-Peer RSVP Flow-Control

The set of recommendations discussed in this section provide an RSVP speaker with the ability to apply back pressure to its peer(s) to reduce/eliminate RSVP-TE control plane congestion.

An implementation that supports "Per-Peer RSVP Flow-Control":

- o MUST support all the recommendations made in Sections [2.1](#) and [2.2](#).
- o MUST treat lack of ACKs from a peer as an indication of peer's RSVP- TE control plane congestion. If congestion is detected, the local system MUST throttle RSVP-TE messages to the affected peer. This MUST be done on a per-peer basis. (Per-peer throttling MAY be implemented by a traffic shaping mechanism that proportionally reduces the RSVP signaling packet rate as the number of outstanding Acks increases. And when the number of outstanding Acks decreases, the send rate would be adjusted up again.)
- o SHOULD use a Retry Limit (Rl) value of 7 ([Section 6.2 of \[RFC2961\]](#), suggests using 3).
- o SHOULD prioritize Hello messages and messages carrying Acknowledgements over other RSVP messages.
- o SHOULD prioritize Tear/Error over trigger Path/Resv (messages that bring up new LSP state) sent to a peer when the local system detects RSVP-TE control plane congestion in the peer.
- o MUST indicate support for all recommendations in this section via the CAPABILITY object in Hello messages.

2.3.1. Capability Advertisement

An implementation supporting the "Per-Peer Flow-Control" recommendations MUST set a new flag "Per-Peer Flow-Control Capable" in the CAPABILITY object signaled in Hello messages.

Bit Number TBA3 (TBA4) - Per-Peer Flow-Control Capable (F-bit):

Indicates that the sender supports Per-Peer RSVP Flow-Control.

Any node that sets the new F-bit in its CAPABILITY object MUST also set Refresh-Reduction-Capable bit in common header of all RSVP-TE messages. If a peer sets the F-bit in the CAPABILITY object but does not set the Refresh-Reduction-Capable bit, then the Per-Peer Flow-Control functionality MUST NOT be activated for that peer.

2.3.2. Compatibility

The Per-Peer Flow-Control functionality MUST NOT be activated with a peer that does not indicate support for this functionality. If a peer hasn't indicated that it is capable of participating in "Per-Peer Flow-Control", then it is risky to assume that the peer would always acknowledge a non-out of order message containing a MESSAGE ID object with the ACK-Desired flag set.

3. Acknowledgements

The authors would like to thank Yakov Rekhter for initiating this work and providing valuable inputs. They would like to thank Raveendra Torvi and Chandra Ramachandran for participating in the many discussions that led to the recommendations made in this document. They would also like to thank Adrian Farrel and Lou Berger for providing detailed review comments.

4. Contributors

Markus Jork
Juniper Networks
Email: mjork@juniper.net

Ebben Aries
Juniper Networks
Email: exa@juniper.net

5. IANA Considerations

5.1. Capability Object Values

IANA maintains all the registries associated with "Resource Reservation Protocol (RSVP) Parameters" (see <http://www.iana.org/assignments/rsvp-parameters/rsvp-parameters.xhtml>). "Capability Object Values" Registry (introduced by [RFC5063]) is one of them.

IANA is requested to assign two new Capability Object Value bit flags as follows:

Bit Number	Hex Value	Name	Reference

TBA1	TBA2	RI-RSVP Capable (I)	Section 2.2.1
TBA3	TBA4	Per-Peer Flow-Control Capable (F)	Section 2.3.1

6. Security Considerations

This document does not introduce new security issues. The security considerations pertaining to the original RSVP protocol [RFC2205] and RSVP-TE [RFC3209] and those that are described in [RFC5920] remain relevant.

7. References

7.1. 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, <<http://www.rfc-editor.org/info/rfc2119>>.
- [RFC2205] Braden, R., Ed., Zhang, L., Berson, S., Herzog, S., and S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", [RFC 2205](#), DOI 10.17487/RFC2205, September 1997, <<http://www.rfc-editor.org/info/rfc2205>>.
- [RFC2961] Berger, L., Gan, D., Swallow, G., Pan, P., Tommasi, F., and S. Molendini, "RSVP Refresh Overhead Reduction Extensions", [RFC 2961](#), DOI 10.17487/RFC2961, April 2001, <<http://www.rfc-editor.org/info/rfc2961>>.

- [RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", [RFC 3209](#), DOI 10.17487/RFC3209, December 2001, <<http://www.rfc-editor.org/info/rfc3209>>.
- [RFC4558] Ali, Z., Rahman, R., Prairie, D., and D. Papadimitriou, "Node-ID Based Resource Reservation Protocol (RSVP) Hello: A Clarification Statement", [RFC 4558](#), DOI 10.17487/RFC4558, June 2006, <<http://www.rfc-editor.org/info/rfc4558>>.
- [RFC5063] Satyanarayana, A., Ed. and R. Rahman, Ed., "Extensions to GMPLS Resource Reservation Protocol (RSVP) Graceful Restart", [RFC 5063](#), DOI 10.17487/RFC5063, October 2007, <<http://www.rfc-editor.org/info/rfc5063>>.

7.2. Informative References

- [RFC5439] Yasukawa, S., Farrel, A., and O. Komolafe, "An Analysis of Scaling Issues in MPLS-TE Core Networks", [RFC 5439](#), DOI 10.17487/RFC5439, February 2009, <<http://www.rfc-editor.org/info/rfc5439>>.
- [RFC5920] Fang, L., Ed., "Security Framework for MPLS and GMPLS Networks", [RFC 5920](#), DOI 10.17487/RFC5920, July 2010, <<http://www.rfc-editor.org/info/rfc5920>>.

Appendix A. Recommended Defaults

(a) Refresh-Interval (R)- 20 minutes ([Section 2.2](#)) Given that an implementation supporting RI-RSVP doesn't rely on refreshes for state sync between peers, the RSVP refresh interval is sort of analogous to IGP refresh interval, the default of which is typically in the order of 10s of minutes. Choosing a default of 20 minutes allows the refresh timer to be randomly set to a value in the range [10 minutes (0.5R), 30 minutes (1.5R)].

(b) Node Hello-Interval - 9 Seconds ([Section 2.2](#)) [[RFC3209](#)] defines the hello timeout as 3.5 times the hello interval. Choosing 9 seconds for the node hello-interval gives a hello timeout of $3.5 \times 9 = 31.5$ seconds. This puts the hello timeout value to be in the same ballpark as the IGP hello timeout value.

(c) Retry-Limit (Rl) - 7 ([Section 2.3](#)) Choosing 7 as the retry-limit results in an overall rapid retransmit phase of 31.5 seconds. This nicely matches up with the 31.5 seconds hello timeout.

(d) Periodic Retransmission Interval - 30 seconds ([Section 2.1.3](#))
If the Retry-Limit (RL) is 7, then it takes about 30 (31.5 to be precise) seconds for the 7 rapid retransmit steps to max out. (The last delay from message 6 to message 7 is 16 seconds). The 30 seconds interval also matches the traditional default refresh time.

Authors' Addresses

Vishnu Pavan Beeram (editor)
Juniper Networks

Email: vbeeram@juniper.net

Ina Minei
Google, Inc

Email: inaminei@google.com

Rob Shakir
Google, Inc

Email: rjs@rob.sh

Dante Pacella
Verizon

Email: dante.j.pacella@verizon.com

Tarek Saad
Cisco Systems

Email: tsaad@cisco.com

