

Internet Engineering Task Force (IETF)
Request for Comments: 7810
Category: Standards Track
ISSN: 2070-1721

S. Previdi, Ed.
Cisco Systems, Inc.
S. Giacalone
Microsoft
D. Ward
Cisco Systems, Inc.
J. Drake
Juniper Networks
Q. Wu
Huawei
May 2016

IS-IS Traffic Engineering (TE) Metric Extensions

Abstract

In certain networks, such as, but not limited to, financial information networks (e.g., stock market data providers), network-performance criteria (e.g., latency) are becoming as critical to data-path selection as other metrics.

This document describes extensions to IS-IS Traffic Engineering Extensions (RFC 5305) such that network-performance information can be distributed and collected in a scalable fashion. The information distributed using IS-IS TE Metric Extensions can then be used to make path-selection decisions based on network performance.

Note that this document only covers the mechanisms with which network-performance information is distributed. The mechanisms for measuring network performance or acting on that information, once distributed, are outside the scope of this document.

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 5741.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <http://www.rfc-editor.org/info/rfc7810>.

Copyright Notice

Copyright (c) 2016 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	3
1.1. Conventions Used in This Document	4
2. TE Metric Extensions to IS-IS	4
3. Interface and Neighbor Addresses	5
4. Sub-TLV Details	6
4.1. Unidirectional Link Delay Sub-TLV	6
4.2. Min/Max Unidirectional Link Delay Sub-TLV	7
4.3. Unidirectional Delay Variation Sub-TLV	8
4.4. Unidirectional Link Loss Sub-TLV	9
4.5. Unidirectional Residual Bandwidth Sub-TLV	10
4.6. Unidirectional Available Bandwidth Sub-TLV	11
4.7. Unidirectional Utilized Bandwidth Sub-TLV	12
5. Announcement Thresholds and Filters	12
6. Announcement Suppression	13
7. Network Stability and Announcement Periodicity	14
8. Enabling and Disabling Sub-TLVs	14
9. Static Metric Override	14
10. Compatibility	14
11. Security Considerations	15
12. IANA Considerations	15
13. References	16
13.1. Normative References	16
13.2. Informative References	16
Acknowledgements	17
Contributors	17
Authors' Addresses	18

1. Introduction

In certain networks, such as, but not limited to, financial information networks (e.g., stock market data providers), network-performance information (e.g., latency) is becoming as critical to data-path selection as other metrics.

In these networks, extremely large amounts of money rest on the ability to access market data in "real time" and to predictably make trades faster than the competition. Because of this, using metrics such as hop count or cost as routing metrics is becoming only tangentially important. Rather, it would be beneficial to be able to make path-selection decisions based on performance data (such as latency) in a cost-effective and scalable way.

This document describes extensions (hereafter called "IS-IS TE Metric Extensions") to the IS-IS Extended Reachability TLV defined in [RFC5305], that can be used to distribute network-performance information (such as link delay, delay variation, packet loss, residual bandwidth, and available bandwidth).

The data distributed by the IS-IS TE Metric Extensions proposed in this document is meant to be used as part of the operation of the routing protocol (e.g., by replacing cost with latency or considering bandwidth as well as cost), to enhance Constrained-SPF (CSPF), or for other uses such as supplementing the data used by an ALTO server [RFC7285]. With respect to CSPF, the data distributed by IS-IS TE Metric Extensions can be used to set up, fail over, and fail back data paths using protocols such as RSVP-TE [RFC3209].

Note that the mechanisms described in this document only disseminate performance information. The methods for initially gathering that performance information, such as described in [RFC6375], or acting on it once it is distributed are outside the scope of this document. Example mechanisms to measure latency, delay variation, and loss in an MPLS network are given in [RFC6374]. While this document does not specify how the performance information should be obtained, the measurement of delay SHOULD NOT vary significantly based upon the offered traffic load. Thus, queuing delays SHOULD NOT be included in the delay measurement. For links such as Forwarding Adjacencies, care must be taken that measurement of the associated delay avoids significant queuing delay; that could be accomplished in a variety of ways, including either by measuring with a traffic class that experiences minimal queuing or by summing the measured link delays of the components of the link's path.

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 RFC 2119 [RFC2119].

In this document, these words will appear with that interpretation only when in ALL CAPS. Lowercase uses of these words are not to be interpreted as carrying the significance described in RFC 2119.

2. TE Metric Extensions to IS-IS

This document registers new IS-IS TE sub-TLVs that can be announced in the "Sub-TLVs for TLVs 22, 23, 141, 222, and 223" registry in order to distribute network-performance information. The extensions in this document build on the ones provided in IS-IS TE [RFC5305] and GMPLS [RFC4203].

IS-IS Extended Reachability TLV 22 (defined in [RFC5305]), Inter-AS Reachability Information TLV 141 (defined in [RFC5316]), and MT-ISIS TLV 222 (defined in [RFC5120]) have nested sub-TLVs that permit the TLVs to be readily extended. This document registers several sub-TLVs:

Type	Description
33	Unidirectional Link Delay
34	Min/Max Unidirectional Link Delay
35	Unidirectional Delay Variation
36	Unidirectional Link Loss
37	Unidirectional Residual Bandwidth
38	Unidirectional Available Bandwidth
39	Unidirectional Utilized Bandwidth

As can be seen in the list above, the sub-TLVs described in this document carry different types of network-performance information. The new sub-TLVs include a bit called the Anomalous (or "A") bit. When the A bit is clear (or when the sub-TLV does not include an A bit), the sub-TLV describes steady-state link performance. This information could conceivably be used to construct a steady-state performance topology for initial tunnel-path computation, or to verify alternative failover paths.

When network performance violates configurable link-local thresholds, a sub-TLV with the A bit set is advertised. These sub-TLVs could be used by the receiving node to determine whether to fail traffic to a backup path or whether to calculate an entirely new path. From an MPLS perspective, the intent of the A bit is to permit label switched path ingress nodes to determine whether the link referenced in the sub-TLV affects any of the label switched paths for which it is ingress. If they are affected, then they can determine whether those label switched paths still meet end-to-end performance objectives. If not, then the node could conceivably move affected traffic to a pre-established protection label switched path or establish a new label switched path and place the traffic in it.

If link performance then improves beyond a configurable minimum value (reuse threshold), that sub-TLV can be re-advertised with the A bit cleared. In this case, a receiving node can conceivably do whatever re-optimization (or fallback) it wishes to do (including nothing).

Note that when a sub-TLV does not include the A bit, that sub-TLV cannot be used for failover purposes. The A bit was intentionally omitted from some sub-TLVs to help mitigate oscillations. See Section 5 for more information.

Consistent with existing IS-IS TE specification [RFC5305], the bandwidth advertisements defined in this document MUST be encoded as IEEE floating-point values. The delay and delay-variation advertisements defined in this document MUST be encoded as integer values. Delay values MUST be quantified in units of microseconds, packet loss MUST be quantified as a percentage of packets sent, and bandwidth MUST be sent as bytes per second. All values (except residual bandwidth) MUST be calculated as rolling averages where the averaging period MUST be a configurable period of time. See Section 5 for more information.

3. Interface and Neighbor Addresses

The use of IS-IS TE Metric Extensions sub-TLVs is not confined to the TE context. In other words, IS-IS TE Metric Extensions sub-TLVs defined in this document can also be used for computing paths in the absence of a TE subsystem.

However, as for the TE case, Interface Address and Neighbor Address sub-TLVs (IPv4 or IPv6) MUST be present. The encoding is defined in [RFC5305] for IPv4 and in [RFC6119] for IPv6.

4. Sub-TLV Details

4.1. Unidirectional Link Delay Sub-TLV

This sub-TLV advertises the average link delay between two directly connected IS-IS neighbors. The delay advertised by this sub-TLV **MUST** be the delay from the local neighbor to the remote one (i.e., the forward-path latency). The format of this sub-TLV is shown in the following diagram:

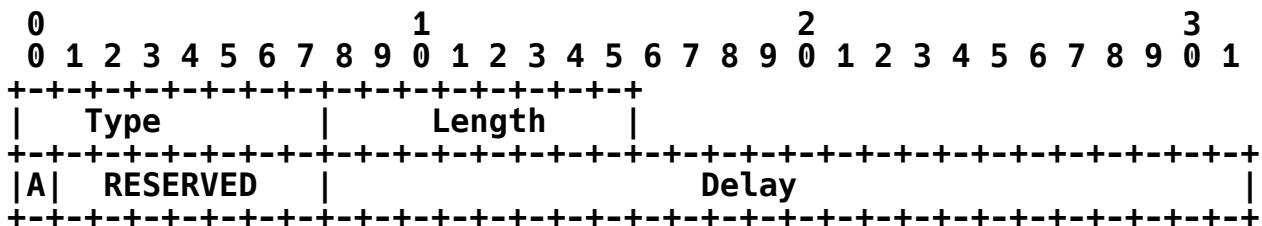


Figure 1

where:

Type: 33

Length: 4

A bit: The A bit represents the Anomalous (A) bit. The A bit is set when the measured value of this parameter exceeds its configured maximum threshold. The A bit is cleared when the measured value falls below its configured reuse threshold. If the A bit is clear, the sub-TLV represents steady-state link performance.

RESERVED: This field is reserved for future use. It **MUST** be set to 0 when sent and **MUST** be ignored when received.

Delay: This 24-bit field carries the average link delay over a configurable interval in microseconds, encoded as an integer value. When set to the maximum value 16,777,215 (16.777215 sec), then the delay is at least that value and may be larger.

4.2. Min/Max Unidirectional Link Delay Sub-TLV

This sub-TLV advertises the minimum and maximum delay values between two directly connected IS-IS neighbors. The delay advertised by this sub-TLV **MUST** be the delay from the local neighbor to the remote one (i.e., the forward-path latency). The format of this sub-TLV is shown in the following diagram:

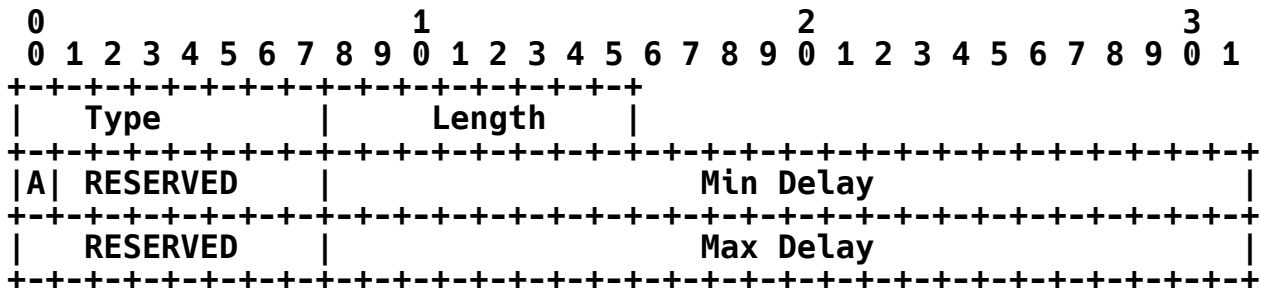


Figure 2

where:

Type: 34

Length: 8

A bit: This field represents the Anomalous (A) bit. The A bit is set when one or more measured values exceed a configured maximum threshold. The A bit is cleared when the measured value falls below its configured reuse threshold. If the A bit is clear, the sub-TLV represents steady-state link performance.

RESERVED: This field is reserved for future use. It **MUST** be set to 0 when sent and **MUST** be ignored when received.

Min Delay: This 24-bit field carries the minimum measured link delay value (in microseconds) over a configurable interval, encoded as an integer value.

Max Delay: This 24-bit field carries the maximum measured link delay value (in microseconds) over a configurable interval, encoded as an integer value.

Implementations **MAY** also permit the configuration of an offset value (in microseconds) to be added to the measured delay value, to facilitate the communication of operator-specific delay constraints.

It is possible for the Min and Max delay to be the same value.

When the delay value (Min or Max) is set to the maximum value 16,777,215 (16.777215 sec), then the delay is at least that value and may be larger.

4.3. Unidirectional Delay Variation Sub-TLV

This sub-TLV advertises the average link delay variation between two directly connected IS-IS neighbors. The delay variation advertised by this sub-TLV **MUST** be the delay from the local neighbor to the remote one (i.e., the forward-path latency). The format of this sub-TLV is shown in the following diagram:

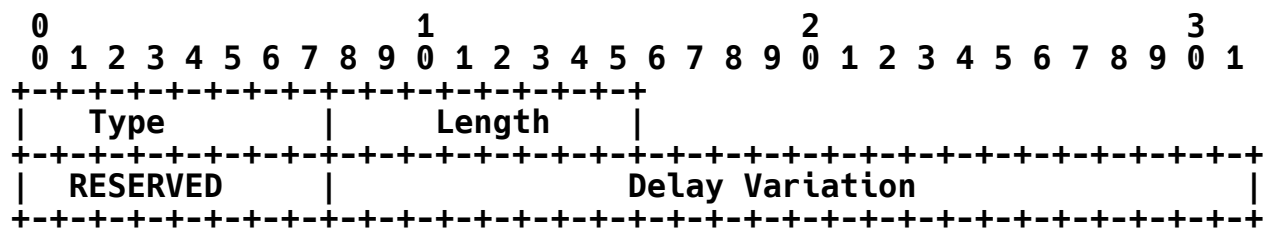


Figure 3

where

Type: 35

Length: 4

RESERVED: This field is reserved for future use. It **MUST** be set to 0 when sent and **MUST** be ignored when received.

Delay Variation: This 24-bit field carries the average link delay variation over a configurable interval in microseconds, encoded as an integer value. When set to 0, it has not been measured. When set to the maximum value 16,777,215 (16.777215 sec), then the delay is at least that value and may be larger.

4.4. Unidirectional Link Loss Sub-TLV

This sub-TLV advertises the loss (as a packet percentage) between two directly connected IS-IS neighbors. The link loss advertised by this sub-TLV **MUST** be the packet loss from the local neighbor to the remote one (i.e., the forward-path loss). The format of this sub-TLV is shown in the following diagram:

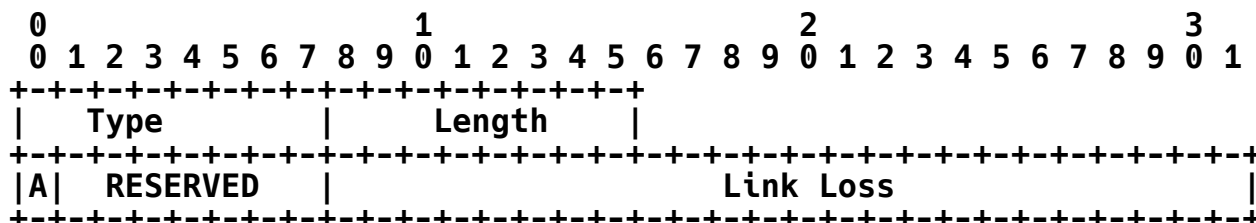


Figure 4

where:

Type: 36

Length: 4

A bit: The A bit represents the Anomalous (A) bit. The A bit is set when the measured value of this parameter exceeds its configured maximum threshold. The A bit is cleared when the measured value falls below its configured reuse threshold. If the A bit is clear, the sub-TLV represents steady-state link performance.

RESERVED: This field is reserved for future use. It **MUST** be set to 0 when sent and **MUST** be ignored when received.

Link Loss: This 24-bit field carries link packet loss as a percentage of the total traffic sent over a configurable interval. The basic unit is 0.000003%, where $(2^{24} - 2)$ is 50.331642%. This value is the highest packet-loss percentage that can be expressed (the assumption being that precision is more important on high-speed links than the ability to advertise loss rates greater than this, and that high-speed links with over 50% loss are unusable). Therefore, measured values that are larger than the field maximum **SHOULD** be encoded as the maximum value.

4.6. Unidirectional Available Bandwidth Sub-TLV

This sub-TLV advertises the available bandwidth between two directly connected IS-IS neighbors. The available bandwidth advertised by this sub-TLV **MUST** be the available bandwidth from the system originating this sub-TLV. The format of this sub-TLV is shown in the following diagram:

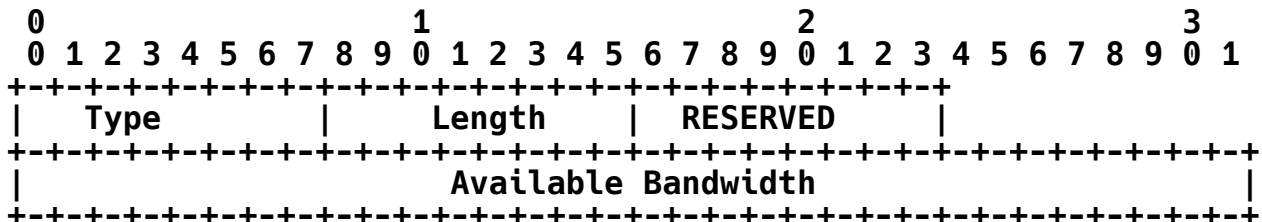


Figure 5

where:

Type: 38

Length: 4

RESERVED: This field is reserved for future use. It **MUST** be set to 0 when sent and **MUST** be ignored when received.

Available Bandwidth: This field carries the available bandwidth on a link, forwarding adjacency, or bundled link in IEEE floating-point format with units of bytes per second. For a link or forwarding adjacency, available bandwidth is defined to be residual bandwidth (see Section 4.5) minus the measured bandwidth used for the actual forwarding of non-RSVP-TE label switched path packets. For a bundled link, available bandwidth is defined to be the sum of the component link available bandwidths minus the measured bandwidth used for the actual forwarding of non-RSVP-TE label switched path packets. For a bundled link, available bandwidth is defined to be the sum of the component link available bandwidths.

The measurement interval, any filter coefficients, and any advertisement intervals **MUST** be configurable per sub-TLV.

In addition to the measurement intervals governing re-advertisement, implementations **SHOULD** provide configurable accelerated advertisement thresholds per sub-TLV, such that:

1. If the measured parameter falls outside a configured upper bound for all but the minimum delay metric (or lower bound for minimum delay metric only) and the advertised sub-TLV is not already outside that bound or,
2. If the difference between the last advertised value and current measured value exceeds a configured threshold then,
3. The advertisement is made immediately.
4. For sub-TLVs that include an A bit, an additional threshold **SHOULD** be included corresponding to the threshold for which the performance is considered anomalous (and sub-TLVs with the A bit are sent). The A bit is cleared when the sub-TLV's performance has been below (or re-crosses) this threshold for an advertisement interval(s) to permit fail back.

To prevent oscillations, only the high threshold or the low threshold (but not both) may be used to trigger any given sub-TLV that supports both.

Additionally, once outside the bounds of the threshold, any re-advertisement of a measurement within the bounds would remain governed solely by the measurement interval for that sub-TLV.

6. Announcement Suppression

When link-performance values change by small amounts that fall under thresholds that would cause the announcement of a sub-TLV, implementations **SHOULD** suppress sub-TLV re-advertisement and/or lengthen the period within which they are refreshed.

Only the accelerated advertisement threshold mechanism described in Section 5 may shorten the re-advertisement interval. All suppression and re-advertisement interval backoff timer features **SHOULD** be configurable.

7. Network Stability and Announcement Periodicity

Sections 5 and 6 provide configurable mechanisms to bound the number of re-advertisements. Instability might occur in very large networks if measurement intervals are set low enough to overwhelm the processing of flooded information at some of the routers in the topology. Therefore, care should be taken in setting these values.

Additionally, the default measurement interval for all sub-TLVs SHOULD be 30 seconds.

Announcements MUST also be able to be throttled using configurable inter-update throttle timers. The minimum announcement periodicity is 1 announcement per second. The default value SHOULD be set to 120 seconds.

Implementations SHOULD NOT permit the inter-update timer to be lower than the measurement interval.

Furthermore, it is RECOMMENDED that any underlying performance-measurement mechanisms not include any significant buffer delay, any significant buffer-induced delay variation, or any significant loss due to buffer overflow or due to active queue management.

8. Enabling and Disabling Sub-TLVs

Implementations MUST make it possible to individually enable or disable each sub-TLV based on configuration.

9. Static Metric Override

Implementations SHOULD permit static configuration and/or manual override of dynamic measurements for each sub-TLV in order to simplify migration and to mitigate scenarios where dynamic measurements are not possible.

10. Compatibility

As per [RFC5305], unrecognized sub-TLVs should be silently ignored.

11. Security Considerations

The sub-TLVs introduced in this document allow an operator to advertise state information of links (bandwidth, delay) that could be sensitive and that an operator may not want to disclose.

Section 7 describes a mechanism to ensure network stability when the new sub-TLVs defined in this document are advertised. Implementation SHOULD follow the described guidelines to mitigate the instability risk.

[RFC5304] describes an authentication method for IS-IS Link State PDUs that allows cryptographic authentication of IS-IS Link State PDUs.

It is anticipated that in most deployments, the IS-IS protocol is used within an infrastructure entirely under control of the same operator. However, it is worth considering that the effect of sending IS-IS Traffic Engineering sub-TLVs over insecure links could result in a man-in-the-middle attacker delaying real-time data to a given site or destination, which could negatively affect the value of the data for that site or destination. The use of Link State PDU cryptographic authentication allows mitigation the risk of man-in-the-middle attack.

12. IANA Considerations

IANA maintains the registry for the sub-TLVs. IANA has registered the following sub-TLVs in the "Sub-TLVs for TLVs 22, 23, 141, 222, and 223" registry:

Type	Description
33	Unidirectional Link Delay
34	Min/Max Unidirectional Link Delay
35	Unidirectional Delay Variation
36	Unidirectional Link Loss
37	Unidirectional Residual Bandwidth
38	Unidirectional Available Bandwidth
39	Unidirectional Utilized Bandwidth

13. References

13.1. 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>>.
- [RFC4206] Kompella, K. and Y. Rekhter, "Label Switched Paths (LSP) Hierarchy with Generalized Multi-Protocol Label Switching (GMPLS) Traffic Engineering (TE)", RFC 4206, DOI 10.17487/RFC4206, October 2005, <<http://www.rfc-editor.org/info/rfc4206>>.
- [RFC5120] Przygienda, T., Shen, N., and N. Sheth, "M-ISIS: Multi Topology (MT) Routing in Intermediate System to Intermediate Systems (IS-ISs)", RFC 5120, DOI 10.17487/RFC5120, February 2008, <<http://www.rfc-editor.org/info/rfc5120>>.
- [RFC5304] Li, T. and R. Atkinson, "IS-IS Cryptographic Authentication", RFC 5304, DOI 10.17487/RFC5304, October 2008, <<http://www.rfc-editor.org/info/rfc5304>>.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", RFC 5305, DOI 10.17487/RFC5305, October 2008, <<http://www.rfc-editor.org/info/rfc5305>>.
- [RFC5316] Chen, M., Zhang, R., and X. Duan, "ISIS Extensions in Support of Inter-Autonomous System (AS) MPLS and GMPLS Traffic Engineering", RFC 5316, DOI 10.17487/RFC5316, December 2008, <<http://www.rfc-editor.org/info/rfc5316>>.
- [RFC6119] Harrison, J., Berger, J., and M. Bartlett, "IPv6 Traffic Engineering in IS-IS", RFC 6119, DOI 10.17487/RFC6119, February 2011, <<http://www.rfc-editor.org/info/rfc6119>>.

13.2. Informative References

- [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>>.

- [RFC4203] Kompella, K., Ed. and Y. Rekhter, Ed., "OSPF Extensions in Support of Generalized Multi-Protocol Label Switching (GMPLS)", RFC 4203, DOI 10.17487/RFC4203, October 2005, <<http://www.rfc-editor.org/info/rfc4203>>.
- [RFC6374] Frost, D. and S. Bryant, "Packet Loss and Delay Measurement for MPLS Networks", RFC 6374, DOI 10.17487/RFC6374, September 2011, <<http://www.rfc-editor.org/info/rfc6374>>.
- [RFC6375] Frost, D., Ed. and S. Bryant, Ed., "A Packet Loss and Delay Measurement Profile for MPLS-Based Transport Networks", RFC 6375, DOI 10.17487/RFC6375, September 2011, <<http://www.rfc-editor.org/info/rfc6375>>.
- [RFC7285] Alimi, R., Ed., Penno, R., Ed., Yang, Y., Ed., Kiesel, S., Previdi, S., Roome, W., Shalunov, S., and R. Woundy, "Application-Layer Traffic Optimization (ALTO) Protocol", RFC 7285, DOI 10.17487/RFC7285, September 2014, <<http://www.rfc-editor.org/info/rfc7285>>.

Acknowledgements

The authors would like to recognize Ayman Soliman, Nabil Bitar, David McDysan, Les Ginsberg, Edward Crabbe, Don Fedyk, Hannes Gredler, Uma Chunduri, Alvaro Retana, Brian Weis, and Barry Leiba for their contribution and review of this document.

The authors also recognize Curtis Villamizar for significant comments and direct content collaboration.

Contributors

The following people contributed substantially to the content of this document and should be considered co-authors:

Alia Atlas
Juniper Networks
United States

Email: akatlas@juniper.net

Clarence Filsfils
Cisco Systems Inc.
Belgium

Email: cfilsfil@cisco.com

Authors' Addresses

Stefano Previdi (editor)
Cisco Systems, Inc.
Via Del Serafico 200
Rome 00191
Italy

Email: sprevidi@cisco.com

Spencer Giacalone
Microsoft

Email: spencer.giacalone@gmail.com

Dave Ward
Cisco Systems, Inc.
3700 Cisco Way
San Jose, CA 95134
United States

Email: wardd@cisco.com

John Drake
Juniper Networks
1194 N. Mathilda Ave.
Sunnyvale, CA 94089
United States

Email: jdrake@juniper.net

Qin Wu
Huawei
101 Software Avenue, Yuhua District
Nanjing, Jiangsu 210012
China

Email: sunseawq@huawei.com