

Network Working Group  
Internet Draft  
Intended status: Proposed Standard  
Expires: April 2012

S. Previdi  
Cisco Systems

S. Giacalone  
Thomson Reuters

D. Ward  
Juniper Networks

J. Drake  
Juniper Networks

A. Atlas  
Juniper Networks

C. Filsfils  
Cisco Systems

October 10, 2011

**IS-IS Traffic Engineering (TE) Metric Extensions**  
**draft-previdi-isis-te-metric-extensions-00.txt**

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 TE [[RFC5305](#)] such that network performance information can be distributed and collected in a scalable fashion. The information distributed using ISIS TE Express Path 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 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), 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>

This Internet-Draft will expire on April 10, 2011.

## Copyright Notice

Copyright (c) 2011 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

<u>1.</u>	<u>Introduction.....</u>	<u>3</u>
<u>2.</u>	<u>Conventions used in this document.....</u>	<u>4</u>
<u>3.</u>	<u>Express Path Extensions to IS-IS TE.....</u>	<u>4</u>
<u>4.</u>	<u>Sub TLV Details.....</u>	<u>6</u>
<u>4.1.</u>	<u>Unidirectional Link Delay Sub-TLV.....</u>	<u>6</u>
<u>4.2.</u>	<u>Unidirectional Delay Variation Sub-TLV.....</u>	<u>6</u>
<u>4.3.</u>	<u>Unidirectional Link Loss Sub-TLV.....</u>	<u>7</u>
<u>4.4.</u>	<u>Unidirectional Residual Bandwidth Sub-TLV.....</u>	<u>8</u>
<u>4.5.</u>	<u>Unidirectional Available Bandwidth Sub-TLV.....</u>	<u>9</u>
<u>5.</u>	<u>Announcement Thresholds and Filters.....</u>	<u>10</u>
<u>6.</u>	<u>Announcement Suppression.....</u>	<u>11</u>
<u>7.</u>	<u>Network Stability and Announcement Periodicity.....</u>	<u>11</u>
<u>8.</u>	<u>Compatibility.....</u>	<u>11</u>
<u>9.</u>	<u>Security Considerations.....</u>	<u>11</u>
<u>10.</u>	<u>IANA Considerations.....</u>	<u>11</u>
<u>11.</u>	<u>References.....</u>	<u>12</u>
<u>11.1.</u>	<u>Normative References.....</u>	<u>12</u>
<u>11.2.</u>	<u>Informative References.....</u>	<u>12</u>
<u>12.</u>	<u>Acknowledgments.....</u>	<u>12</u>
<u>13.</u>	<u>Author's Addresses.....</u>	<u>12</u>

## **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 to IS-IS Extended Reachability TLV [RFC5305](hereafter called "IS-IS TE Express Path"), 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 IS-IS TE Express Path 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), by enhancing



CSPF, or for other uses such as supplementing the data used by an Alto server [[Alto](#)]. With respect to CSPF, the data distributed by IS-IS TE Express Path can be used to setup, 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 [[Frost](#)], or acting on it once it is distributed are outside the scope of this document.

## **2. 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. Lower case uses of these words are not to be interpreted as carrying [RFC-2119](#) significance.

## **3. Express Path Extensions to IS-IS TE**

This document proposes new IS-IS TE sub-TLVs that can be announced in ISIS Extended Reachability TLV (TLV-22) 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 (TLV-22) defined in [[RFC5305](#)] has nested sub-TLVs which permit the ISIS Reachability TLV to be readily extended. This document proposes several additional sub-TLVs:

Type	Length	Value
TBD1	6	Unidirectional Link Delay
TBD2	6	Unidirectional Delay Variation
TBD3	6	Unidirectional Packet Loss
TBD4	6	Unidirectional Residual Bandwidth Sub TLV
TBD5	6	Unidirectional Available Bandwidth Sub TLV

As can be seen in the list above, the sub-TLVs described in this document carry different types of network performance information. Many (but not all) of the 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 LSP ingress nodes to:

- A) Determine whether the link referenced in the sub-TLV affects any of the LSPs for which it is ingress. If there are, then:
- B) Determine whether those LSPs still meet end-to-end performance objectives. If not, then:
- C) The node could then conceivably move affected traffic to a pre-established protection LSP or establish a new LSP 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 Anomalous bit cleared. In this case, a receiving node can conceivably do whatever re-optimization (or failback) 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.

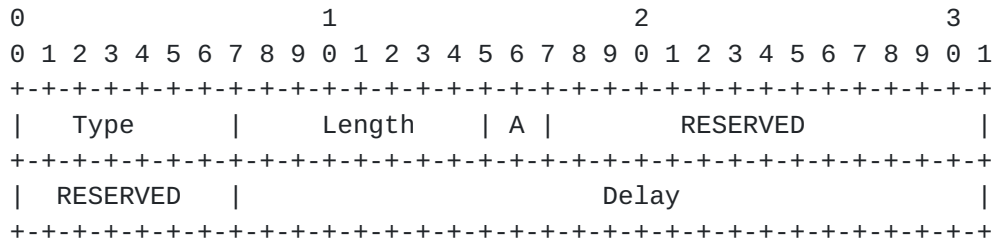
Consistent with existing IS-IS TE specifications [[RFC5305](#)], the bandwidth advertisements defined in this draft MUST be encoded as IEEE floating point values. The delay and delay variation advertisements defined in this draft MUST be encoded as integer values. Delay values MUST be quantified in units 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.



**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:



This sub-TLV has a type of TBD1.  
The length is 6.  
Where:

"A" 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.

The "Reserved" field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

"Delay Value" is a 24-bit field carries the average link delay over a configurable interval in micro-seconds, 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.

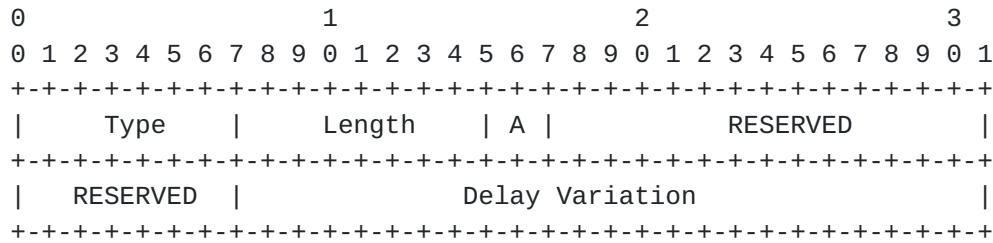
**4.2. Unidirectional Delay Variation Sub-TLV**

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:



This sub-TLV has a type of TBD2. The length is 6.

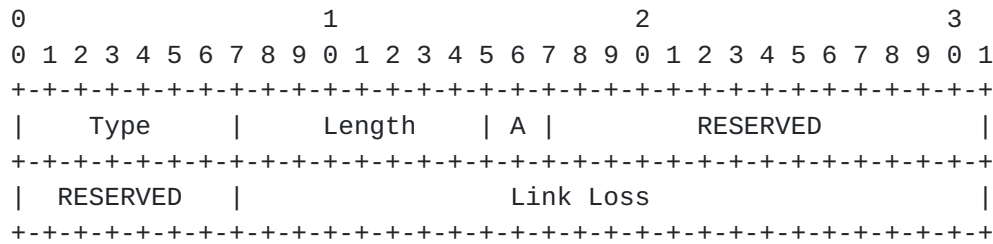
Where:

"A" 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. The "Reserved" field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

"Delay Variation" is a 24-bit field carries the average link delay variation over a configurable interval in micro-seconds, 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.

**4.3. 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:



This sub-TLV has a type of TBD3.  
The length is 6.

Where:

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" field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

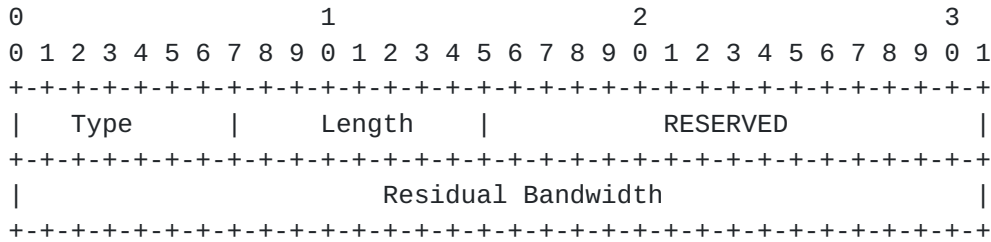
"Link Loss" is a 24-bit field carries link packet loss as a percentage of the total traffic sent over a configurable interval. The basic unit is 0.000003%, where (2^24 - 2) is 50.331642%. This value is the highest packet loss percentage that can be expressed (the assumption being that precision is more 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. When set to a value of all 1s (2^24 - 1), the link packet loss has not been measured.

**4.4. Unidirectional Residual Bandwidth Sub-TLV**

This TLV advertises the residual bandwidth (defined in section Error! Reference source not found.between two directly connected IS-IS neighbors. The residual bandwidth advertised by this sub-TLV MUST be the residual bandwidth from the system originating the sub-TLV to its neighbor.

The format of this sub-TLV is shown in the following diagram:





This sub-TLV has a type of TBD4.  
The length is 6.

"Reserved" field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

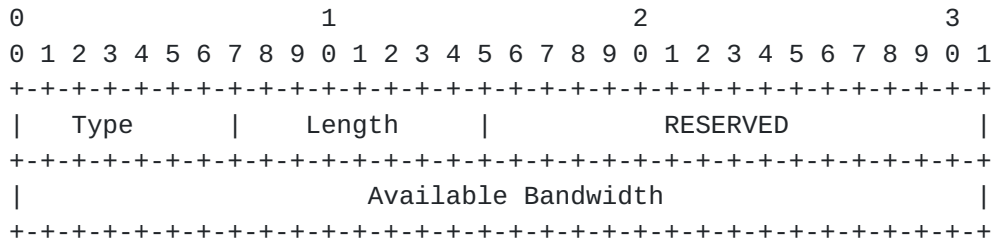
"Residual Bandwidth" is a field carries the residual bandwidth on a link, forwarding adjacency [RFC4206], or bundled link in IEEE floating point format with units of bytes per second. For a link or forwarding adjacency, residual bandwidth is defined to be Maximum Link Bandwidth [RFC5305] minus the bandwidth currently allocated to RSVP-TE LSPs. For a bundled link, residual bandwidth is defined to be the sum of the component link residual bandwidths.

Note that although it may seem possible to calculate Residual Bandwidth using the existing sub-TLVs in RFC 5305, this is not a consistently reliable approach and hence the Residual Bandwidth sub-TLV has been added here. For example, because the Maximum Reservable Bandwidth [RFC5305] can be larger than the capacity of the link, using it as part of an algorithm to determine the value of the Maximum Link Bandwidth [RFC5305] minus the bandwidth currently allocated to RSVP-TE Label Switched Paths cannot be considered reliably accurate.

**4.5. Unidirectional Available Bandwidth Sub-TLV**

This TLV advertises the available bandwidth (defined in section Error! Reference source not found.between two directly connected IS-IS neighbors. The available bandwidth advertised by this sub-TLV MUST be the available bandwidth from the system originating the Sub-TLV to

its neighbor. The format of this sub-TLV is shown in the following diagram:



This sub-TLV has a type of TBD5. The length is 6.

Where:

"Reserved" field is reserved for future use. It MUST be set to 0 when sent and MUST be ignored when received.

"Available Bandwidth" is a field that 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 4.4. minus the measured bandwidth used for the actual forwarding of non-RSVP-TE Label Switched Paths packets. For a bundled link, available bandwidth is defined to be the sum of the component link available bandwidths.

**5. Announcement Thresholds and Filters**

The values advertised in all sub-TLVs MUST be controlled using an exponential filter (i.e. a rolling average) with a configurable measurement interval and filter coefficient.

Implementations are expected to provide separately configurable advertisement thresholds. All thresholds MUST be configurable on a per sub-TLV basis.

The announcement of all sub-TLVs that do not include the A bit SHOULD be controlled by variation thresholds that govern when they are sent.

Sub-TLV that include the A bit are governed by several thresholds. Firstly, a threshold SHOULD be implemented to govern the announcement

of sub-TLVs that advertise a change in performance, but not an SLA violation (i.e. when the A bit is not set). Secondly, implementations MUST provide configurable thresholds that govern the announcement of sub-TLVs with the A bit set (for the indication of a performance violation). Thirdly, implementations SHOULD provide reuse thresholds. These thresholds govern sub-TLV re-announcement with the A bit cleared to permit fail back.

## **6. Announcement Suppression**

When link performance average values change, but fall under the threshold that would cause the announcement of a sub-TLV with the A bit set, implementations MAY suppress or throttle sub-TLV announcements. All suppression features and thresholds SHOULD be configurable.

## **7. Network Stability and Announcement Periodicity**

To mitigate concerns about stability, all values (except residual bandwidth) MUST be calculated as rolling averages where the averaging period MUST be a configurable period of time, rather than instantaneous measurements.

Announcements MUST also be able to be throttled using configurable inter-update throttle timers. The minimum announcement periodicity is 1 announcement per second.

## **8. Compatibility**

As per ([RFC5305](#)), unrecognized TLVs should be silently ignored

## **9. Security Considerations**

This document does not introduce security issues beyond those discussed in [[RFC3630](#)] and [[RFC5329](#)].

## **10. IANA Considerations**

IANA maintains the registry for the sub-TLVs. IS-IS TE Express Path will require one new type code per sub-TLV defined in this document.

## **11. References**

### **11.1. Normative References**

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC5305] Li, T., Smit, H., "IS-IS Extensions for Traffic Engineering", [RFC 3630](#), September 2003.

### **11.2. Informative References**

- [RFC3031] Rosen, E., Viswanathan, A., Callon, R., "Multiprotocol Label Switching Architecture", January 2001
- [RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", [RFC 3209](#), December 2001.
- [Frost] D. Frost, S. Bryant "A Packet Loss and Delay Measurement Profile for MPLS-based Transport Networks"
- [Alto] R. Alimi R. Penno Y. Yang, "ALTO Protocol"

## **12. Acknowledgments**

The authors would like to recognize Ayman Soliman and Les Ginsberg for their contributions.

This document was prepared using 2-Word-v2.0.template.dot.



**13. Author's Addresses**

Stefano Previdi  
Cisco Systems  
Via Del Serafico 200  
00142 Rome  
Italy

Email: sprevidi@cisco.com

Spencer Giacalone  
Thomson Reuters  
195 Broadway  
New York NY 10007, USA

Email: Spencer.giacalone@thomsonreuters.com

Dave Ward  
Juniper Networks  
1194 N. Mathilda Ave.  
Sunnyvale, CA 94089, USA

Email: dward@juniper.net

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

Email: jdrake@juniper.net

Alia Atlas  
Juniper Networks  
1194 N. Mathilda Ave.  
Sunnyvale, CA 94089, USA

Email: akatlas@juniper.net

Clarence Filsfils

Cisco Systems  
Brussels, Belgium

Email: cfilsfil@cisco.com