

Networking Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 3, 2019

N. Shen
Cisco Systems
S. Amante
Apple, Inc.
M. Abrahamsson
T-Systems Nordic
July 2, 2018

IS-IS Routing with Reverse Metric
draft-ietf-isis-reverse-metric-11

Abstract

This document describes a mechanism to allow IS-IS routing to quickly and accurately shift traffic away from either a point-to-point or multi-access LAN interface during network maintenance or other operational events. This is accomplished by signaling adjacent IS-IS neighbors with a higher reverse metric, i.e., the metric towards the signaling IS-IS router.

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, 2019.

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 (<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.	Node and Link Isolation	2
1.2.	Distributed Forwarding Planes	3
1.3.	Spine-Leaf Applications	3
1.4.	LDP IGP Synchronization	3
1.5.	IS-IS Reverse Metric	3
1.6.	Specification of Requirements	4
2.	IS-IS Reverse Metric TLV	4
3.	Elements of Procedure	6
3.1.	Processing Changes to Default Metric	6
3.2.	Multi-Topology IS-IS Support on Point-to-point links . .	7
3.3.	Multi-Access LAN Procedures	7
3.4.	Point-To-Point Link Procedures	8
3.5.	LDP/IGP Synchronization on LANs	8
3.6.	Operational Guidelines	9
4.	Security Considerations	9
5.	IANA Considerations	10
6.	Acknowledgments	10
7.	References	10
7.1.	Normative References	10
7.2.	Informative References	11
Appendix A.	Node Isolation Challenges	11
Appendix B.	Link Isolation Challenges	12
Appendix C.	Contributors' Addresses	13
	Authors' Addresses	13

[1. Introduction](#)

The IS-IS [[ISO10589](#)] routing protocol has been widely used in Internet Service Provider IP/MPLS networks. Operational experience with the protocol, combined with ever increasing requirements for lossless operations have demonstrated some operational issues. This document describes the issues and a mechanism for mitigating them.

[1.1. Node and Link Isolation](#)

IS-IS routing mechanism has the overload-bit, which can be used by operators to perform disruptive maintenance on the router. But in many operational maintenance cases, it is not necessary to divert all the traffic away from this node. It is necessary to avoid only a single link during the maintenance. More detailed descriptions of

the challenges can be found in [Appendix A](#) and [Appendix B](#) of this document.

1.2. Distributed Forwarding Planes

In a distributed forwarding platform, different forwarding line-cards may have interfaces and IS-IS connections to neighbor routers. If one of the line-card's software resets, it may take some time for the forwarding entries to be fully populated on the line-card, in particular if the router is a PE (Provider Edge) router in ISP's MPLS VPN. An IS-IS adjacency may be established with a neighbor router long before the entire BGP VPN prefixes are downloaded to the forwarding table. It is important to signal to the adjacent IS-IS routers to raise metric values and not to use the corresponding IS-IS adjacency inbound to this router if possible. Temporarily signaling the 'Reverse Metric' over this link to discourage the traffic via the corresponding line-card will help to reduce the traffic loss in the network. In the meantime, the remote PE routers will select a different set of PE routers for the BGP best path calculation or use a different link towards the same PE router on which a line-card is resetting.

1.3. Spine-Leaf Applications

In the IS-IS Spine-Leaf extension [[I-D.shen-isis-spine-leaf-ext](#)], the leaf nodes will perform equal-cost or unequal-cost load sharing towards all the spine nodes. In certain operational cases, for instance, when one of the backbone links on a spine node is congested, a spine node can push a higher metric towards the connected leaf nodes to reduce the transit traffic through the corresponding spine node or link.

1.4. LDP IGP Synchronization

In the [[RFC5443](#)], a mechanism is described to achieve LDP IGP synchronization by using the maximum link metric value on the interface. But in the case of a new IS-IS node joining the broadcast network (LAN), it is not optimal to change all the nodes on the LAN to the maximum link metric value, as described in [[RFC6138](#)]. In this case, the Reverse Metric can be used to discourage both outbound and inbound traffic without affecting the traffic of other IS-IS nodes on the LAN.

1.5. IS-IS Reverse Metric

This document uses the routing protocol itself as the transport mechanism to allow one IS-IS router to advertise a "reverse metric" in an IS-IS Hello (IIH) PDU to an adjacent node on a point-to-point

or multi-access LAN link. This would allow the provisioning to be performed only on a single node, setting a "reverse metric" on a link and have traffic bidirectionally shift away from that link gracefully to alternate, viable paths.

This Reverse Metric mechanism is used for both point-to-point and multi-access LAN links. Unlike the point-to-point links, the IS-IS protocol currently does not have a way to influence the traffic towards a particular node on LAN links. This mechanism provides IS-IS routing the capability of altering traffic in both directions on either a point-to-point link or a multi-access link of an IS-IS node.

The metric value in the "reverse metric" TLV and the TE metric in the sub-TLV being advertised is an offset or relative metric to be added to the existing local link and TE metric values of the receiver.

1.6. Specification of Requirements

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. IS-IS Reverse Metric TLV

The Reverse Metric TLV is composed of a 1 octet field of Flags, a 3 octet field containing an IS-IS Metric Value, and a 1 octet Traffic Engineering (TE) sub-TLV length field representing the length of a variable number of Extended Intermediate System (IS) Reachability sub-TLVs. If the "sub-TLV len" is non-zero, then the Value field MUST also contain one or more Extended IS Reachability sub-TLVs.

The Reverse Metric TLV is optional. The Reverse Metric TLV may be present in any IS-IS Hello PDU. A sender MUST only transmit a single Reverse Metric TLV in a IS-IS Hello PDU. If a received IS-IS Hello PDU contains more than one Reverse Metric TLV, an implementation SHOULD ignore all the Reverse Metric TLVs and treat it as an error condition.

U bit (0x02): The "Unreachable" bit specifies that the metric calculated by addition of the reverse metric value to the "default metric" is limited to $(2^{24}-1)$. This "U" bit applies to both the default metric in the Extended IS Reachability TLV and the TE default-metric sub-TLV of the link. This is only relevant to the IS-IS "wide" metric mode.

The Reverse Metric TLV can include sub-TLVs when an IS-IS router wishes to signal to its neighbor to raise its Traffic Engineering (TE) Metric over the link. In this document, only the "Traffic Engineering Default Metric" sub-TLV [RFC5305], sub-TLV Type 18, is defined and MAY be included in the Reverse Metric TLV, because that is a similar 'reverse metric' operation to be used in TE computations. Upon receiving this TE METRIC sub-TLV in a Reverse Metric TLV, a node SHOULD add the received TE metric offset value to its existing, configured TE default metric within its Extended IS Reachability TLV. Use of other sub-TLVs is outside the scope of this document. The "sub-TLV Len" value MUST be set to zero when an IS-IS router does not have TE sub-TLVs that it wishes to send to its IS-IS neighbor.

3. Elements of Procedure

3.1. Processing Changes to Default Metric

The Metric field, in the Reverse Metric TLV, is a "reverse offset metric" that will either be in the range of 0 - 63 when a "narrow" IS-IS metric is used (IS Neighbors TLV, Pseudonode LSP) [RFC1195] or in the range of 0 - $(2^{24} - 2)$ when a "wide" Traffic Engineering metric value is used, (Extended IS Reachability TLV) [RFC5305] [RFC5817]. It is important to use the same IS-IS metric mode on both ends of the link. On the receiving side of the 'reverse-metric' TLV, the accumulated value of configured metric and the reverse-metric needs to be limited to 63 in "narrow" metric mode and to $(2^{24} - 2)$ in "wide" metric mode. This applies to both the default metric of Extended IS Reachability TLV and the TE default-metric sub-TLV in LSP or Pseudonode LSP for the "wide" metric mode case. If the "U" bit is present in the flags, the accumulated metric value is to be limited to $(2^{24} - 1)$ for both the normal link metric and TE metric in IS-IS "wide" metric mode.

If an IS-IS router is configured to originate a TE Default Metric sub-TLV for a link, but receives a Reverse Metric TLV from its neighbor that does not contain a TE Default Metric sub-TLV, then the IS-IS router MUST NOT change the value of its TE Default Metric sub-TLV for that link.

3.2. Multi-Topology IS-IS Support on Point-to-point links

The Reverse Metric TLV is applicable to Multi-Topology IS-IS (M-ISIS) [[RFC5120](#)]. On point-to-point links, if an IS-IS router is configured for M-ISIS, it MUST send only a single Reverse Metric TLV in IIH PDUs toward its neighbor(s) on the designated link. When an M-ISIS router receives a Reverse Metric TLV, it MUST add the received Metric value to its default metric in all Extended IS Reachability TLVs for all topologies. If an M-ISIS router receives a Reverse Metric TLV with a TE Default Metric sub-TLV, then the M-ISIS router MUST add the received TE Default Metric value to each of its TE Default Metric sub-TLVs in all of its MT Intermediate Systems TLVs. If an M-ISIS router is configured to advertise TE Default Metric sub-TLVs for one or more topologies, but does not receive a TE Default Metric sub-TLV in a Reverse Metric TLV, then the M-ISIS router MUST NOT change the value in each of the TE Default Metric sub-TLVs for all topologies.

3.3. Multi-Access LAN Procedures

On a Multi-Access LAN, only the DIS SHOULD act upon information contained in a received Reverse Metric TLV. All non-DIS nodes MUST silently ignore a received Reverse Metric TLV. The decision process of the routers on the LAN MUST follow the procedure in [section 7.2.8.2](#) of [[ISO10589](#)], and use the "Two-way connectivity check" during the topology and route calculation.

The Reverse Metric TE sub-TLV also applies to the DIS. If a DIS is configured to apply TE over a link and it receives TE metric sub-TLV in a Reverse Metric TLV, it should update the TE Default Metric sub-TLV value of the corresponding Extended IS Reachability TLV or insert a new one if not present.

In the case of multi-access LANs, the "W" Flags bit is used to signal from a non-DIS to the DIS whether to change the metric and, optionally Traffic Engineering parameters for all nodes in the Pseudonode LSP or solely the node on the LAN originating the Reverse Metric TLV.

A non-DIS node, e.g., Router B, attached to a multi-access LAN will send the DIS a Reverse Metric TLV with the W bit clear when Router B wishes the DIS to add the Metric value to the default metric contained in the Pseudonode LSP specific to just Router B. Other non-DIS nodes, e.g., Routers C and D, may simultaneously send a Reverse Metric TLV with the W bit clear to request the DIS to add their own Metric value to their default metric contained in the Pseudonode LSP. When the DIS receives a properly formatted Reverse Metric TLV with the W bit clear, the DIS MUST only add the default

metric contained in its Pseudonode LSP for the specific neighbor that sent the corresponding Reverse Metric TLV.

As long as at least one IS-IS node on the LAN sending the signal to DIS with the W bit set, the DIS would add the metric value in the Reverse Metric TLV to all neighbor adjacencies in the Pseudonode LSP, regardless if some of the nodes on the LAN advertise the Reverse Metric TLV without the W bit set. The DIS MUST use the reverse metric of the highest source MAC address Non-DIS advertising the Reverse Metric TLV with the W bit set. The DIS MUST use the metric value towards the nodes which explicitly advertise the Reverse Metric TLV.

Local provisioning on the DIS to adjust the default metric(s) contained in the Pseudonode LSP MUST take precedence over received Reverse Metric TLVs. For instance, local policy on the DIS may be provisioned to ignore the W bit signaling on a LAN.

Multi-Topology IS-IS [[RFC5120](#)] specifies there is no change to construction of the Pseudonode LSP, regardless of the Multi-Topology capabilities of a multi-access LAN. If any MT capable node on the LAN advertises the Reverse Metric TLV to the DIS, the DIS should update, as appropriate, the default metric contained in the Pseudonode LSP. If the DIS updates the default metric in and floods a new Pseudonode LSP, those default metric values will be applied to all topologies during Multi-Topology SPF calculations.

[3.4.](#) Point-To-Point Link Procedures

On a point-to-point link, there is already a "configured" IS-IS interface metric to be applied over the link towards the IS-IS neighbor.

When IS-IS receives the IIH PDU with the "Reverse Metric" on a point-to-point link and if the local policy allows the supporting of "Reverse Metric", it MUST add the metric value in "reverse metric" TLV according to the rules described in [Section 3.1](#) and [Section 3.2](#).

[3.5.](#) LDP/IGP Synchronization on LANs

As described in [[RFC6138](#)] when a new IS-IS node joins a broadcast network, it is unnecessary and sometimes even harmful for all IS-IS nodes on the LAN to advertise maximum link metric. [[RFC6138](#)] proposes a solution to have the new node not advertise its adjacency towards the pseudo-node when it is not in a "cut-edge" position.

With the introduction of Reverse Metric in this document, a simpler alternative solution to the above mentioned problem can be used. The

Reverse Metric allows the new node on the LAN to advertise its inbound metric value to be the maximum and this puts the link of this new node in the last resort position without impacting the other IS-IS nodes on the same LAN.

Specifically, when IS-IS adjacencies are being established by the new node on the LAN, besides setting the maximum link metric value ($2^{24} - 2$) on the interface of the LAN for LDP IGP synchronization as described in [\[RFC5443\]](#), it SHOULD advertise the maximum metric offset value in the Reverse Metric TLV in its IIH PDU sent on the LAN. It SHOULD continue this advertisement until it completes all the LDP label binding exchanges with all the neighbors over this LAN, either by receiving the LDP End-of-LIB [\[RFC5919\]](#) for all the sessions or by exceeding the provisioned timeout value for the node LDP/IGP synchronization.

3.6. Operational Guidelines

A router MUST advertise a Reverse Metric TLV toward a neighbor only for the period during which it wants a neighbor to temporarily update its IS-IS metric or TE parameters towards it.

The use of Reverse Metric does not alter IS-IS metric parameters stored in a router's persistent provisioning database.

Routers that receive a Reverse Metric TLV MAY send a syslog message or SNMP trap, in order to assist in rapidly identifying the node in the network that is advertising an IS-IS metric or Traffic Engineering parameters different from that which is configured locally on the device.

When the link TE metric is raised to ($2^{24} - 1$) [\[RFC5817\]](#), either due to the reverse-metric mechanism or by explicit user configuration, this SHOULD immediately trigger the CSPF re-calculation to move the TE traffic away from that link. It is RECOMMENDED also that the CSPF does the immediate CSPF re-calculation when the TE metric is raised to ($2^{24} - 2$) to be the last resort link.

It is RECOMMENDED that implementations provide a capability to disable any changes to a node's individual interface default metric or Traffic Engineering parameters based upon receiving a properly formatted Reverse Metric TLVs.

4. Security Considerations

The enhancement in this document makes it possible for one IS-IS router to manipulate the IS-IS default metric and, optionally, Traffic Engineering parameters of adjacent IS-IS neighbors. Although

IS-IS routers within a single Autonomous System nearly always are under the control of a single administrative authority, it is highly RECOMMENDED that operators configure authentication of IS-IS PDUs to mitigate use of the Reverse Metric TLV as a potential attack vector, particularly on multi-access LANs.

5. IANA Considerations

This document requests that IANA allocate from the IS-IS TLV Codepoints Registry a new TLV, referred to as the "Reverse Metric" TLV, possibly from the "Unassigned" range of 244-250, with the following attributes: IIH = y, LSP = n, SNP = n, Purge = n.

6. Acknowledgments

The authors would like to thank Mike Shand, Dave Katz, Guan Deng, Ilya Varlashkin, Jay Chen, Les Ginsberg, Peter Ashwood-Smith, Uma Chunduri, Alexander Okonnikov, Jonathan Harrison, Dave Ward, Himanshu Shah, Wes George, Danny McPherson, Ed Crabbe, Russ White, Robert Razsuk, Tom Petch and Acee Lindem for their comments and contributions.

This document was produced using Marshall Rose's xml2rfc tool.

7. References

7.1. Normative References

- [ISO10589] ISO, "Intermediate system to Intermediate system routing information exchange protocol for use in conjunction with the Protocol for providing the Connectionless-mode Network Service (ISO 8473)", ISO/IEC 10589:2002.
- [RFC1195] Callon, R., "Use of OSI IS-IS for routing in TCP/IP and dual environments", [RFC 1195](#), DOI 10.17487/RFC1195, December 1990, <<https://www.rfc-editor.org/info/rfc1195>>.
- [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>>.
- [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, <<https://www.rfc-editor.org/info/rfc5120>>.

- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", [RFC 5305](#), DOI 10.17487/RFC5305, October 2008, <<https://www.rfc-editor.org/info/rfc5305>>.
- [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>>.

7.2. Informative References

- [I-D.shen-isis-spine-leaf-ext]
Shen, N., Ginsberg, L., and S. Thyamagundalu, "IS-IS Routing for Spine-Leaf Topology", [draft-shen-isis-spine-leaf-ext-03](#) (work in progress), March 2017.
- [RFC5443] Jork, M., Atlas, A., and L. Fang, "LDP IGP Synchronization", [RFC 5443](#), DOI 10.17487/RFC5443, March 2009, <<https://www.rfc-editor.org/info/rfc5443>>.
- [RFC5817] Ali, Z., Vasseur, JP., Zamfir, A., and J. Newton, "Graceful Shutdown in MPLS and Generalized MPLS Traffic Engineering Networks", [RFC 5817](#), DOI 10.17487/RFC5817, April 2010, <<https://www.rfc-editor.org/info/rfc5817>>.
- [RFC5919] Asati, R., Mohapatra, P., Chen, E., and B. Thomas, "Signaling LDP Label Advertisement Completion", [RFC 5919](#), DOI 10.17487/RFC5919, August 2010, <<https://www.rfc-editor.org/info/rfc5919>>.
- [RFC6138] Kini, S., Ed. and W. Lu, Ed., "LDP IGP Synchronization for Broadcast Networks", [RFC 6138](#), DOI 10.17487/RFC6138, February 2011, <<https://www.rfc-editor.org/info/rfc6138>>.

Appendix A. Node Isolation Challenges

On rare occasions, it is necessary for an operator to perform disruptive network maintenance on an entire IS-IS router node, i.e., major software upgrades, power/cooling augments, etc. In these cases, an operator will set the IS-IS Overload Bit (OL-bit) within the Link State Protocol Data Units (LSPs) of the IS-IS router about to undergo maintenance. The IS-IS router immediately floods its updated LSPs to all IS-IS routers in the IS-IS domain. Upon receipt of the updated LSPs, all IS-IS routers recalculate their Shortest Path First (SPF) tree excluding IS-IS routers whose LSPs have the OL-bit set. This effectively removes the IS-IS router about to undergo maintenance from the topology, thus preventing it from receiving any transit traffic during the maintenance period.

After the maintenance activity has completed, the operator resets the IS-IS Overload Bit within the LSPs of the original IS-IS router causing it to flood updated IS-IS LSPs throughout the IS-IS domain. All IS-IS routers recalculate their SPF tree and now include the original IS-IS router in their topology calculations, allowing it to be used for transit traffic again.

Isolating an entire IS-IS router from the topology can be especially disruptive due to the displacement of a large volume of traffic through an entire IS-IS router to other, sub-optimal paths, (e.g., those with significantly larger delay). Thus, in the majority of network maintenance scenarios, where only a single link or LAN needs to be augmented to increase its physical capacity or is experiencing an intermittent failure, it is much more common and desirable to gracefully remove just the targeted link or LAN from service, temporarily, so that the least amount of user-data traffic is affected during the link-specific network maintenance.

Appendix B. Link Isolation Challenges

Before network maintenance events are performed on individual physical links or LANs, operators substantially increase the IS-IS metric simultaneously on both devices attached to the same link or LAN. In doing so, the devices generate new Link State Protocol Data Units (LSPs) that are flooded throughout the network and cause all routers to gradually shift traffic onto alternate paths with very little or no disruption to in-flight communications by applications or end-users. When performed successfully, this allows the operator to confidently perform disruptive augmentation, fault diagnosis or repairs on a link without disturbing ongoing communications in the network.

There are a number of challenges with the above solution. First, it is quite common to have routers with several hundred interfaces and individual interfaces that are from several hundred Gigabits/second to Terabits/second of traffic. Thus, it is imperative that operators accurately identify the same point-to-point link on two, separate devices in order to increase (and, afterward, decrease) the IS-IS metric appropriately. Second, the aforementioned solution is very time consuming and even more error-prone to perform when it's necessary to temporarily remove a multi-access LAN from the network topology. Specifically, the operator needs to configure ALL devices that have interfaces attached to the multi-access LAN with an appropriately high IS-IS metric, (and then decrease the IS-IS metric to its original value afterward). Finally, with respect to multi-access LANs, there is currently no method to bidirectionally isolate only a single node's interface on the LAN when performing more fine-grained diagnosis and repairs to the multi-access LAN.

In theory, use of a Network Management System (NMS) could improve the accuracy of identifying the appropriate subset of routers attached to either a point-to-point link or a multi-access LAN as well as signaling from the NMS to those devices, using a network management protocol to adjust the IS-IS metrics on the pertinent set of interfaces. The reality is that NMSs are, to a very large extent, not used within Service Provider's networks for a variety of reasons. In particular, NMSs do not interoperate very well across different vendors or even separate platform families within the same vendor.

The risks of misidentifying one side of a point-to-point link or one or more interfaces attached to a multi-access LAN and subsequently increasing its IS-IS metric and potentially increased latency, jitter or packet loss. This is unacceptable given the necessary performance requirements for a variety of reasons including the customer perception for near lossless operations and the associated demanding Service Level Agreement's (SLAs) for all network services.

[Appendix C](#). Contributors' Addresses

Tony Li

Email: tony.li@tony.li

Authors' Addresses

Naiming Shen
Cisco Systems
560 McCarthy Blvd.
Milpitas, CA 95035
USA

Email: naiming@cisco.com

Shane Amante
Apple, Inc.
1 Infinite Loop
Cupertino, CA 95014
USA

Email: samante@apple.com

Mikael Abrahamsson
T-Systems Nordic
Kistagangen 26
Stockholm
SE

Email: Mikael.Abrahamsson@t-systems.se