

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: 14 March 2026

P. Mohapatra
Google LLC
R. Das, Ed.
Juniper Networks, Inc.
S. Mohanty, Ed.
Zscaler
S. Krier
Cisco Systems
R.J. Szarecki
Google LLC
A. Gattani
Arista Networks
10 September 2025

BGP Link Bandwidth Extended Community
draft-ietf-idr-link-bandwidth-17

Abstract

This document specifies a type of BGP Extended Community that enables routers to perform weighted load-balancing in multipath scenarios.

Requirements Language

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].

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 <https://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 14 March 2026.

Copyright Notice

Copyright (c) 2025 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	2
2. Link Bandwidth Extended Community	3
3. Protocol Procedures	4
3.1. Sender (Originating Link Bandwidth Extended Community)	4
3.2. Receiver (Receiving Link Bandwidth Extended Community)	4
3.3. Re-advertisement Procedures	5
3.3.1. Re-advertisement with Next hop Self	5
3.3.2. Re-advertisement with Next Hop Unchanged	5
3.4. Link Bandwidth Extended Community Arithmetic and BGP Multipath	5
4. Error Handling	6
5. IANA Considerations	6
6. Security Considerations	7
7. Operational Considerations	7
7.1. Inconsistent Deployment	7
8. Contributors	8
9. Acknowledgments	8
10. Normative References	8
Appendix A. Document History	9
Authors' Addresses	9

1. Introduction

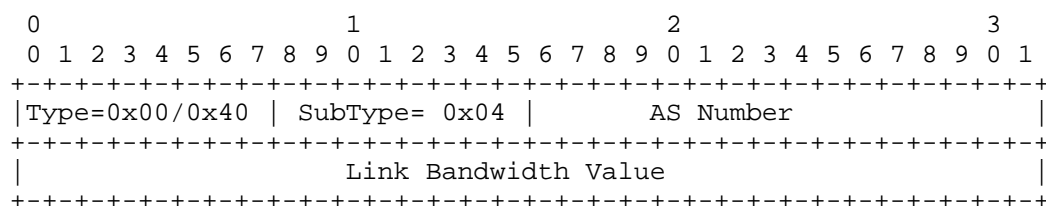
Load balancing is a critical aspect of network design, enabling efficient utilization of available bandwidth and improving overall network performance. Traditional equal-cost multi-path (ECMP) routing does not account for the varying capacities of different paths. This document suggests that the link bandwidth be carried in the network using one of two new extended communities [RFC4360] - the transitive and non-transitive Link Bandwidth Extended Community. The Link Bandwidth Extended Community provides a mechanism for routers to advertise the bandwidth of their downstream path that may either be a

directly connected link or multi-hop/multipath nexthop. This mechanism facilitates maximizing utilization of network resources.

2. Link Bandwidth Extended Community

The Link Bandwidth Extended Community is defined as a BGP extended community that carries the bandwidth information of a router, represented by BGP Next Hop, connecting to a remote network. This community can be used to inform other routers about the available bandwidth through a given route.

The Link Bandwidth Extended Community can be either transitive or non-transitive. Therefore the value of the high-order octet of the extended Type Field can be 0x00 or 0x40, respectively. The value of the low-order octet of the extended type field for this communities is 0x04. The value of the Global Administrator subfield in the Value Field SHOULD represent the Autonomous System of the router that attaches the Link Bandwidth Extended Community, but it can be set to any 2-byte value. If the Autonomous System number cannot be represented in two octets, AS_TRANS [RFC6793], SHOULD be used in the Global Administrator subfield. The encoding of 4-octet ASN is out of scope of this document. The bandwidth of the link is expressed as 4 octets in [IEEE.754-2019] floating point format, units being bytes (not bits!) per second. It is carried in the Local Administrator subfield of the Value Field.



Type: 1-octet field MUST be set to 0x00 or 0x40 to indicate transitive/non-transitive.

SubType: 1-octet field MUST be set to 0x04 to indicate 'Link-Bandwidth'.

Global Administrator sub-field:
2-octet represent the Autonomous System.

Local Administrator sub-field:
Bandwidth value (bytes per sec) encoded as 4 octets in IEEE floating point format.

Figure 1: Link Bandwidth Extended Community

3. Protocol Procedures

The procedures cover both the transitive and non-transitive variants of the Link Bandwidth Extended Community so that implementations can handle both variants in a way that supports existing deployments. Please refer to Section 5 and Appendix A for more details.

3.1. Sender (Originating Link Bandwidth Extended Community)

A BGP speaker that attaches a Link Bandwidth Extended Community SHOULD be able to advertise either a transitive or a non-transitive Link Bandwidth Extended Community. Implementations SHOULD provide configuration to set the transitivity type of the Link Bandwidth Extended Community, as well as the Global Administrator and bandwidth values in the Local Administrator field, using local policy. Different implementations MAY use different default values for the transitivity type of the Link Bandwidth Extended Community. The provided configuration SHOULD allow operators to override the default transitivity value as needed. An implementation MAY advertise a link bandwidth value as zero.

Generally, a single Link Bandwidth Extended Community of the transitivity type that is desired in a deployment is attached to a route. However during transition (refer Section 7 for details), a BGP speaker MAY attach one Link Bandwidth Extended Community per transitivity (transitive/non-transitive) both having the same 'Link Bandwidth Value' field.

A Link Bandwidth Extended Community MAY be attached or updated for a BGP route upon receipt during Adj-RIB-In processing. The Link Bandwidth Extended Community MAY be attached or updated for a BGP route's Adj-RIB-Out entry while being advertised to a neighboring BGP speaker.

Implementations MAY provide a configuration option to send non-transitive Link Bandwidth Extended Communities on external BGP sessions.

3.2. Receiver (Receiving Link Bandwidth Extended Community)

A BGP receiver MUST be able to process Link Bandwidth Extended Community of both transitive and non-transitive types. The receiver MUST NOT flap or treat the route as malformed based on the transitivity of the Link Bandwidth Extended Community and/or BGP session type (internal vs. external).

Implementations MAY provide configuration to accept non-transitive Link Bandwidth Extended Communities from external BGP sessions.

A BGP update with an attached Link Bandwidth Extended Community with a bandwidth value of zero is valid. Weighted ECMP (WECMP) described in section 6.3 [RFC7938] can be utilized when all contributing paths have a non-zero value in the Link Bandwidth Extended Community. However, in the case where the paths have a mix of zero and non-zero values, or all zero values, the behavior is determined by local policy. For example, implementations MAY exclude the paths with zero value from WECMP formation as long as at least one path with non-zero value exists or they MAY fallback to ECMP.

3.3. Re-advertisement Procedures

This section describes the procedures to be followed when a BGP speaker receives a route with an attached Link Bandwidth Extended Community and subsequently re-advertises that route.

3.3.1. Re-advertisement with Next hop Self

When a BGP speaker re-advertises a route with Link Bandwidth Extended Community and sets the next hop to itself, it SHOULD follow the same procedures as outlined in Section 3.1.

In the absence of any route policies that alter the Link Bandwidth Extended Community, any received Link Bandwidth Extended Community on the route will be re-advertised unchanged. Please also refer to Section 3.4 for use in a BGP multipath environment.

3.3.2. Re-advertisement with Next Hop Unchanged

A BGP speaker that receives a route with a Link Bandwidth Extended Community and re-advertises or reflects the same without changing its next hop, SHOULD NOT change the Link Bandwidth Extended Community in any way.

3.4. Link Bandwidth Extended Community Arithmetic and BGP Multipath

In a BGP multipath ECMP environment, the link bandwidth value that is sent or re-advertised may be calculated based on the Link Bandwidth Extended Community of the routes contributing to multipath in the Local Routing Information Base (Local-RIB). This topic is beyond the scope of this document.

4. Error Handling

If a BGP speaker receives a route with more than one Link Bandwidth Extended Communities and uses the route to compute WECP, it SHOULD use the extended community with the lowest "Link Bandwidth Value", ignoring the transitivity. Implementations MAY provide configuration to change the above preference.

Between transitive and non-transitive types of Link Bandwidth Extended Communities that have the same 'Link Bandwidth Value', the transitivity doesn't matter for purpose of computing WECP or programming to FIB (Forwarding Information Base).

Note that these procedures mean that a BGP speaker reflecting a route with next hop unchanged (e.g. RR) will re-advertise the Link Bandwidth Extended Communities received on the route as-is without any modification, while following the extended community transitivity rules.

Link Bandwidth Extended Communities with a negative value SHALL be ignored and MUST NOT be advertised.

Link Bandwidth Extended Communities with a zero value MUST NOT be considered malformed.

If any of the paths lack a valid Link Bandwidth Extended Community, ECMP (Equal-Cost Multi-Path) MUST be used instead.

5. IANA Considerations

IANA is requested to update the Transitive Two-Octet AS-Specific Extended Community Sub-Types registry (Type 0x00) and Sub-Type 0x04 to:

```
Name
----
transitive Link Bandwidth Extended Community
```

IANA is requested to update the Non-Transitive Two-Octet AS-Specific Extended Community Sub-Types registry (Type 0x40) and Sub-Type 0x04 to:

```
Name
----
non-transitive Link Bandwidth Extended Community
```

Both updates are to reference this document.

6. Security Considerations

This extension to BGP has similar security implications as BGP Extended Communities [RFC4360]

The Link Bandwidth Extended Community conveys bandwidth and capacity information that may be sensitive. Exporting this community outside of an administrative domain can expose private network resource details. When propagating the routes with Link Bandwidth Extended Community towards an untrusted network or outside of an administrative domain, it is recommended operators use policy to filter out this community.

7. Operational Considerations

7.1. Inconsistent Deployment

Prior deployments of the feature specified in this document have involved implementations that only understood one of the two extended community transitivity types. As a result, such implementations would treat the use of the other transitivity type in a "ships in the night" fashion. The procedures in this document govern how multiple transitivity types for link bandwidth should operate.

In circumstances where networks have deployed a mixture of implementations supporting this document's procedures for both transitivity types, and older implementations that only understand one transitivity type, inconsistent behavior could result. A prime example is when a route received by a BGP speaker contains both a transitive and a non-transitive Link Bandwidth Extended Community and that BGP speaker performs an operation that updates only one of the Link Bandwidth Extended Communities, the other community may have an inconsistent value. As a result, downstream BGP speakers that may receive such routes may perform inappropriate WECMP load balancing.

To mitigate such issues, when operators are aware that older implementations are present in their networks, they may wish to take actions to address such inconsistencies. One option would be to filter either at advertisement time on the older BGP speaker the unsupported transitivity type of Link Bandwidth Extended Community - if the implementation is capable of such filtering. Alternatively, a receiving BGP speaker, knowing that the sending speaker is incapable of doing such operations, could strip the Link Bandwidth Extended Community type that is unsupported by the sender.

Ideally this operational consideration is short-lived until all the routers in the network have been upgraded to implementations that consistently support the procedures in this document.

8. Contributors

Kaliraj Vairavakkalai
Juniper Networks, Inc.
1133 Innovation Way,
Sunnyvale, CA 94089
United States of America
Email: kaliraj@juniper.net

Natrajan Venkataraman
Juniper Networks, Inc.
1133 Innovation Way,
Sunnyvale, CA 94089
United States of America
Email: natv@juniper.net

Rex Fernando
Cisco Systems
170 W. Tasman Drive
San Jose, CA 95134
United States of America
Email: rex@cisco.com

9. Acknowledgments

The authors would like to thank Yakov Rekhter, Srihari Sangli and Dan Tappan for proposing unequal cost load balancing as one possible application of the extended community attribute. The authors would like to thank Jeff Haas for all the discussions and providing text for operational considerations.

The authors would like to thank Bruno Decraene, Robert Raszuk, Joel Halpern, Aleksi Suhonen, Randy Bush, Stephane Litkowski, Mankamana Mishra, Moshiko Nayman, Yingzhen Qu, Anoop Ghanwani, Dongjie (Jimmy) and John Scudder for their comments and contributions.

10. Normative References

- [IEEE.754-2019]
IEEE, "IEEE Standard for Floating-Point Arithmetic", 22
July 2019, <<https://ieeexplore.ieee.org/document/8766229>>.

- [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>>.
- [RFC4360] Sangli, S., Tappan, D., and Y. Rekhter, "BGP Extended Communities Attribute", RFC 4360, DOI 10.17487/RFC4360, February 2006, <<https://www.rfc-editor.org/info/rfc4360>>.
- [RFC6793] Vohra, Q. and E. Chen, "BGP Support for Four-Octet Autonomous System (AS) Number Space", RFC 6793, DOI 10.17487/RFC6793, December 2012, <<https://www.rfc-editor.org/info/rfc6793>>.
- [RFC7938] Lapukhov, P., Premji, A., and J. Mitchell, Ed., "Use of BGP for Routing in Large-Scale Data Centers", RFC 7938, DOI 10.17487/RFC7938, August 2016, <<https://www.rfc-editor.org/info/rfc7938>>.

Appendix A. Document History

BGP Link Bandwidth Extended Community has evolved over several versions of the IETF draft. In the earlier versions up to draft-ietf-idr-link-bandwidth-08, only the non-transitive version of Link Bandwidth Extended Community was supported. However, starting from draft-ietf-idr-link-bandwidth-09, both transitive and non-transitive versions of Link Bandwidth Extended Community are supported.

A BGP speaker (Sender or Receiver) needs to be upgraded to support the procedures defined in this document to provide full interoperability for both transitive and non-transitive versions of Link Bandwidth Extended Community. In order to simplify implementations, it is not a goal to provide interoperability by upgrading only the RR.

Authors' Addresses

Pradosh Mohapatra
Google LLC
Email: pradosh@google.com

Reshma Das (editor)
Juniper Networks, Inc.
1133 Innovation Way,
Sunnyvale, CA 94089
United States of America
Email: dreshma@juniper.net

Satya Mohanty (editor)
Zscaler
120 Holger Way,
San Jose, CA 95134
United States of America
Email: smohanty@zscaler.com

Serge Krier
Cisco Systems
Pegasus Parc, De Kleetlaan 6a
Belgium
Email: sekrier@cisco.com

Rafal Jan Szarecki
Google LLC
1160 N Mathilda Ave,
Sunnyvale, CA 94089
United States of America
Email: rszarecki@gmail.com

Akshay Gattani
Arista Networks
5453 Great America Parkway
Santa Clara, CA 95054
United States of America
Email: akshay@arista.com