COIN Internet-Draft Intended status: Standards Track Expires: June 11, 2020

New IPv6 Multicast Addresses for Switch ML draft-hsingh-ipv6-coin-ml-03

Abstract

Recently, in-network aggregation to scale distributed machine learning (ML) has been presented. A network switch implementation uses IPv4 broadcast messages from switch to the hosts to send updates to all workers. IPv6 does not support broadcast addresses. This document proposes, IPv6 implementations use the IPv6 link-local allnodes multicast address, until a new IPv6 link-local multicast address is assigned by IANA for switch to hosts multicast communications.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of $\underline{\text{BCP 78}}$ and $\underline{\text{BCP 79}}$.

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 <u>https://datatracker.ietf.org/drafts/current/</u>.

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 June 11, 2020.

Copyright Notice

Copyright (c) 2019 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>https://trustee.ietf.org/license-info</u>) 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> .	Requirements Language	•	•	 	•	•	•		•	<u>2</u>
<u>2</u> .	Introduction			 						<u>2</u>
<u>3</u> .	Additional Information			 						<u>2</u>
<u>4</u> .	P4 Considerations			 						<u>3</u>
<u>5</u> .	Security Considerations			 						<u>3</u>
<u>6</u> .	IANA Considerations			 						<u>3</u>
<u>7</u> .	Acknowledgements			 						<u>4</u>
<u>8</u> .	References			 						<u>4</u>
8	3.1. Normative References			 						<u>4</u>
8	3.2. Informative References	•		 						<u>4</u>
Auth	hor's Address			 						<u>4</u>

<u>1</u>. 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 <u>RFC 2119</u> [<u>RFC2119</u>].

2. Introduction

New computing in the network for ML uses IPv4 broadcast communications between switch to hosts. [Switch-ML]. With IPv6, multicast communications would be used. This document proposes a new well-known multicast address be defined for such communications. Until a new address is defined, the IPv6 link-local all-nodes multicast address may be used. By definition, a layer-2 switch operates in the link-local subnet. Thus, the IPv6 link-local multicast address defined by this document suffices for switch to host multicast communications.

3. Additional Information

It is common when new networking protocols such as RPL [RFC6550] and Babel [RFC6126] were developed, each protocol requested IANA for a new IPv6 link-local multicast address for use. Switch ML does not have a protocol defined by IETF just yet and may never define one. However, experiments in switch ML are already using IP or layer-2 broadcast communications. For IPv6, switch ML experiments should use IPv6 link-local multicast communications. A new IPv6 link-local multicast address for switch ML facilitates efficient filtering by hosts. Singh

[Page 2]

If a switch is configured in layer-3 mode and if switch ML communicates with hosts to another IPv6 subnet, an IPv6 Site-Local Scope Multicast address is recommended for communications.

4. P4 Considerations

P4 issues arise when implementing the paper. A P4 header stack (array) is used to represent the vector. P4 has no for loops, so how does one add elements of the vector? There is the one way. You can have one action that adds 1 header, another that adds 2, another that adds 3, etc., and choose between them at run time via a table lookup. Even with a vector with 20-50 elements, writing such P4 code manually is risky. One has to write a Python script to generate such code. The paper did not think of it, but an alternative exists. With every element in the vector, use a new 1-byte field. The field has values 0 or 1. 1 represents bos (bottom of stack) like what MPLS uses. The last element has bos set to 1. Then, one has a termination condition in the P4 parser to terminate recursive parsing of vector. As the parser is recursing through the vector, add the vector elements and save sum in metadata. One does not use a bit for bos because the p4c bmv2 model accepts a header on byte boundary.

The same p4c bmv2 model, also, does not accept an index into array if the index is a run-time value. The p4-16 specification allows runtime index. Algorithm 1 in the paper uses "p.idx" which is a runtime value. Again, one has to use P4 tables. Include the run-time variable index as a table search key, which selects between several different actions that are identical, except for the constant array index values they use. This gets cumbersome with 128 or 512 indices in the paper. Again, one would have to generate P4 code to avoid mistakes cut-and-pasting code. For 128 indices with some bit shifting, maybe we reduces indices to 32. If an asic allows run-time index into array, then it is better to use the asic simulator rather than bmv2 simulator.

5. Security Considerations

Use IPSec [RFC4301].

6. IANA Considerations

This document requests IANA to assign a new IPv6 link-local multicast address for use by network ML. This multicast address name is Switch_ML_Host. An interface on the host MUST join this well-known multicast address.

Additionally, IANA is requested to assign a new IPv6 Site-Local Scope Multicast address for switch ML to host communication across IPv6 Singh

[Page 3]

subnets. If configured to do so, an interface on the host MUST join this well-known multicast address.

7. Acknowledgements

Thanks (in alphabetical order by first name) to Marco Canini for his review.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, DOI 10.17487/RFC2119, March 1997, <<u>https://www.rfc-editor.org/info/rfc2119</u>>.
- [RFC4301] Kent, S. and K. Seo, "Security Architecture for the Internet Protocol", <u>RFC 4301</u>, DOI 10.17487/RFC4301, December 2005, <<u>https://www.rfc-editor.org/info/rfc4301</u>>.
- [RFC6126] Chroboczek, J., "The Babel Routing Protocol", <u>RFC 6126</u>, DOI 10.17487/RFC6126, April 2011, <<u>https://www.rfc-editor.org/info/rfc6126</u>>.
- [RFC6550] Winter, T., Ed., Thubert, P., Ed., Brandt, A., Hui, J., Kelsey, R., Levis, P., Pister, K., Struik, R., Vasseur, JP., and R. Alexander, "RPL: IPv6 Routing Protocol for Low-Power and Lossy Networks", <u>RFC 6550</u>, DOI 10.17487/RFC6550, March 2012, <<u>https://www.rfc-editor.org/info/rfc6550</u>>.

<u>8.2</u>. Informative References

[Switch-ML]

Sapio, A., Canini, M., Ho, C., Nelson, J., Kalnis, P., Kim, C., Krishnamurthy, A., Moshref, M., Ports, D. R., and P. Richtarik, "SwitchML: Scaling Distributed Machine Learning with In-Network Aggregation", February 2019, <<u>https://arxiv.org/pdf/1903.06701.pdf</u>>.

Author's Address

Singh

Expires June 11, 2020 [Page 4]

Hemant Singh MNK Consulting 7 Caldwell Drive Westford, MA 01886 USA

Phone: +1 978 692 2340 Email: hemant@mnkcg.com URI: <u>http://mnkcg.com/</u>