Network Working Group Internet-Draft

Intended status: Experimental

Expires: April 7, 2013

Dino Farinacci
Greg Shepherd
Stig Venaas
Cisco Systems
Yiqun Cai
Microsoft
October 4, 2012

Population Count Extensions to PIM draft-ietf-pim-pop-count-07.txt

Abstract

This specification defines a method for providing multicast distribution-tree accounting data. Simple extensions to the PIM protocol allow a rough approximation of tree-based data in a scalable fashion.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of \underline{BCP} 78 and \underline{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 April 7, 2013.

Copyright Notice

Copyright (c) 2012 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> .	Introduction	<u>3</u>
1	1. Requirements Notation	<u>4</u>
1	2. Terminology	<u>4</u>
<u>2</u> .	Pop-Count-Supported Hello Option	<u>5</u>
<u>3</u> .	New Pop-Count Join Attribute Format	<u>6</u>
3	${ t 1}$. Options	9
	3.1.1. Link Speed Encoding	<u>10</u>
3	$\underline{2}$. Example message layouts	<u>11</u>
<u>4</u> .	How to use Pop-Count Encoding	<u>13</u>
<u>5</u> .	Implementation Approaches	<u>14</u>
<u>6</u> .	Caveats	<u>15</u>
<u>7</u> .	IANA Considerations	<u>16</u>
<u>8</u> .	Security Considerations	<u>17</u>
<u>9</u> .	Acknowledgments	<u>18</u>
	References	
10	<u>.1</u> . Normative References	<u>19</u>
10	<u>.2</u> . Informative References	<u>19</u>
Auth	ors' Addresses	20

1. Introduction

This document specifies a mechanism to convey accounting information using the PIM protocol [RFC4601] [RFC5015]. Putting the mechanism in PIM allows efficient distribution and maintenance of such accounting information. Previous mechanisms require data to be correlated from multiple router sources.

This mechanism allows a single router to be queried to obtain accounting and statistic information for a multicast distribution tree as a whole or any distribution sub-tree downstream from a queried router. The amount of information is fixed and does not increase as multicast membership, tree diameter, or branching increase.

The sort of accounting data this specification provides, on a per multicast route basis, are:

- 1. The number of branches in a distribution tree.
- 2. The membership type of the distribution tree, that is Source-Specific Multicast (SSM) or Any-Source Multicast (ASM).
- 3. Routing domain and time zone boundary information.
- 4. On-tree node and tree diameter counters.
- 5. Effective MTU and bandwidth.

This document defines a new PIM Join Attribute type [RFC5384] to the Join/Prune message as well as a new Hello option. The mechanism is applicable to IPv4 and IPv6 multicast.

This is a new extension to PIM, and it is not completely understood what impact collecting information using PIM would have on the operation of PIM. This is an entirely new concept. Many PIM features (including the core protocols) were first introduced as Experimental RFCs, and it seems appropriate to advance this work as Experimental. Reports of implementation and deployment across whole distribution trees or within sub-trees (see Section 6) will enable an assessment of the desirability and stability of this specification. The PIM working group will then consider whether to move this work to the Standards Track.

This document does not specify how an administrator or user can access this information. It is expected that an implementation may have a command line interface or other ways of requesting and displaying this information. As this is currently an Experimental

document, defining a MIB module has not been considered. If the PIM working group finds that this should move on to Standards Track, a MIB module should be considered.

1.1. Requirements Notation

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].

1.2. Terminology

This section defines the terms used in this document.

Multicast Route: A (S,G) or (*,G) entry regardless if the route is in ASM, SSM, or Bidir mode of operation.

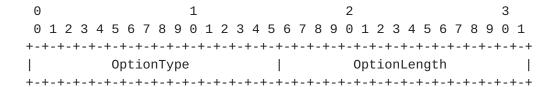
Stub Link: A link with members joined to the group via IGMP or Multicast Listener Discovery (MLD)

Transit Link: A link put in the oif-list (outgoing interface list) for a multicast route because it was joined by PIM routers.

Note that a link can be both a Stub Link and a Transit Link at the same time.

2. Pop-Count-Supported Hello Option

A PIM router indicates that it supports the mechanism specified in this document by including the Pop-Count-Supported Hello option in its PIM Hello message. Note that it also needs to include the Join-Attribute Hello option as specified in [RFC5384]. The format of the Pop-Count-Supported Hello option is defined to be:



OptionType = TBD1, OptionLength = 0. Note that there is no option value included. In order to allow future updates of this specification that may include an option value, implementations of this document MUST accept and process this option also if the length is non-zero. Implementations of this specification MUST accept and process the option ignoring any option value that may be included.

3. New Pop-Count Join Attribute Format

When a PIM router supports this mechanism and has determined from a received Hello, that the neighbor supports this mechanism, and also that all the neighbors on the interface support the use of join attributes, it will send Join/Prune messages that MAY include a Pop-Count Join Attribute. The mechanism to process a PIM Join Attribute is described in [RFC5384]. The format of the new attribute is specified in the following.

0	1	2 3	}
0 1 2 3 4 5 6 7 8 9	0 1 2 3 4 5 6 7 8 9	0 1 2 3 4 5 6 7 8 9 6) 1
+-+-+-+-	+-+-+-+-+-+-+-+-	+-+-+-+-+-	-+-+
F E Attr Type	Length	Effective MTU	
+-+-+-+-	+-+-+-+-+-	+-+-+-+-+-	-+-+
Flags	1	Options Bitmap	
+-+-+-+-	+-+-+-+-+-	+-+-+-+-+-	-+-+
1	Options		- 1
i	i		
•	•		
•	•		
+-+-+-+-+-	+-+-+-+-+-+-+-+-	+-+-+-+-+-+-	-+-+

The above format is used only for entries in the join-list section of the Join/Prune message.

F bit: 0 Non-Transitive Attribute.

E bit: As specified by [RFC5384].

Attr Type: TBD2.

Length: The minimum length is 6.

Effective MTU: This contains the minimum MTU for any link in the oif-list. The sender of Join/Prune message takes the minimum value for the MTU (in bytes) from each link in the oif-list. If this value is less than the value stored for the multicast route (the one received from downstream joiners) then the value should be reset and sent in Join/Prune message. Otherwise, the value should remain unchanged.

This provides one to obtain the MTU supported by multicast distribution tree when examined at the first-hop router(s) or for sub-tree for any router on the distribution tree.

The flags field has the following format: Flags:

- Unallocated/Reserved Flags: The flags which are currently not defined. If a new flag is defined and used by a new implementation, an old implementation should preserve the bit settings. This means that a router MUST preserve the settings of all Unallocated/Reserved Flags in PIM Join messages received from downstream routers in any PIM Join sent upstream.
- If an IGMPv3 or MLDv2 report with an INCLUDE Mode group S flag: record was received on any oif-list entry or the bit was set from any PIM Join message. This bit should only be cleared when the above becomes untrue.
- A flag: If an IGMPv3 or MLDv2 report with an EXCLUDE Mode group record, or an IGMPv1, IGMPv2, or MLDv1 report, was received on any oif-list entry or the bit was set from any PIM Join message. This bit should only be cleared when the above becomes untrue.

A combination of settings for these bits indicate:

A-flag	S-flag	Description
0	0	There are no members for the group
		('Stub Oif-List Count' is 0)
0	1	All group members are using SSM
1	Θ	All group members are using ASM
1	1	A mixture of SSM and ASM group members

If there are any tunnels on the distribution tree. If a tunnel is in the oif-list, a router should set this bit in its Join/Prune messages. Otherwise, it propagates the bit setting from downstream joiners.

- a flag: If there are any auto-tunnels on the distribution tree. If an auto-tunnel is in the oif-list, a router should set this bit in its Join/Prune messages. Otherwise, it propagates the bit setting from downstream joiners. An example of an autotunnel is an tunnel setup by the AMT [I-D.ietf-mboned-auto-multicast] protocol.
- P flag: This flag is set by a router if all downstream routers support this specification. That is, they are all PIM popcount capable. If a downstream router does not support this specification it MUST be cleared. This allows one to tell if the entire sub-tree is completely accounting capable.

Options Bitmap: This is a bitmap that shows which options are present. The format of the bitmap is as follows:

0										1						
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	
+	+ - +	- -	+		- - +	- -	- -		- - +	- - +	- -	- -		-	- - +	-
T	s	m	M	d	n	D	z	ļι	Jna	al]	Loc	C/F	RSI	٢V٥	k	
+	+ - +	H – H	 	- -	- - +	- - +	H – H	H - H	- - +	- - +	- - +	H – H	H - H	- -	- +	-

Each one of the bits T, s, m, M, d, n, D and z is associated with one option, where the option is included if and only if the respective bit is set. Included options MUST be in the same order as these bits are listed. The bits denote the following options:

bit	Option
Т	Transit Oif-List Count
S	Stub Oif-List Count
m	Minimum Speed Link
M	Maximum Speed Link
d	Domain Count
n	Node Count
D	Diameter Count
Z	TZ Count

See Section 3.1 for details on the different options. The unallocated bits are reserved. Any unknown bits MUST be set to 0 when a message is sent, and treated as 0 (ignored) when received. This means that unknown options which are denoted by unknown bits are ignored.

By using this bitmap we can specify at most 16 options. If there becomes a need for more than 16 options, one can define a new option that contains a bitmap, which can then be used to specify which further options are present. The last bit in the current bitmap could be used for that option. The exact definition of this is however left for future documents.

Options: This field contains options. Which options are present are determined by the flag bits. As new flags and options may be defined in the future, any unknown/reserved flags MUST be ignored, and any additional trailing options MUST be ignored. See Section 3.1 for details on the options defined in this document.

3.1. Options

There are several options defined in this document. For each option, there is also a related flag that shows whether the option is present. See the Options Bitmap above for a list of the options and their respective bits. Each option has a fixed size. Note that there is no alignment requirements for the options, so an implementation cannot assume they are aligned.

Transit Oif-List Count: This is filled in by a router sending a Join/Prune message indicating the number of transit links on the multicast distribution tree. The value is the number of oifs (outgoing interfaces) for the multicast route that have been joined by PIM plus the sum of the values advertised by each of the downstream PIM routers that have joined on this oif. Length 4 octets.

Stub Oif-List Count: This is filled in by a router sending a Join/ Prune message indicating the number of stub links (links where there are host members) on the multicast distribution tree. The value is the number of of oifs for the multicast route that have been joined by IGMP or MLD plus the sum of the values advertised by each of the downstream PIM routers that have joined on this oif. Length 4 octets.

Minimum Speed Link: This contains the minimum bandwidth rate for any link in the oif-list and is encoded as specified in Section 3.1.1. The sender of Join/Prune message takes the minimum value for each link in the oif-list for the multicast route. If this value is less than the value stored for the multicast route (the smallest value received from downstream joiners) then the value should be reset and sent in Join/Prune message. Otherwise, the value should remain unchanged. This together with the Maximum Speed Link option provides a way to obtain the lowest and highest speed link for the multicast distribution tree. Length 2 octets.

Maximum Speed Link: This contains the maximum bandwidth rate for any link in the oif-list and is encoded as specified in Section 3.1.1. The sender of Join/Prune message takes the maximum value for each link in the oif-list for the multicast route. If this value is greater than the value stored for the multicast route (the largest value received from downstream joiners) then the value should be reset and sent in Join/Prune message. Otherwise, the value should remain unchanged. This together with the Minimum Speed Link option provides a way to obtain the lowest and highest speed link for the multicast distribution tree. Length 2 octets.

Domain Count: This indicates the number of routing domains the distribution tree traverses. A router should increment this value if it is sending a Join/Prune message over a link which traverses a domain boundary. For this to work, an implementation needs a way of knowing that a neighbor or an interface is in a different domain. There is no standard way of doing this. Length 1 octet.

Node Count: This indicates the number of routers on the distribution tree. Each router will sum up all the Node Counts from all joiners on all oifs and increment by 1 before including this value in the Join/Prune message. Length 1 octet.

Diameter Count: This indicates the longest length of any given branch of the tree in router hops. Each router that sends a Join increments the max value received by all downstream joiners by 1. Length 1 octet.

This indicates the number of timezones the distribution TZ Count: tree traverses. A router should increment this value if it is sending a Join/Prune message over a link which traverses a time zone. This can be a configured link attribute or use other means to determine the timezone is acceptable. Length 1 octet.

3.1.1. Link Speed Encoding

The speed is encoded using 2 octets as follows:

1 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 | Exponent | Significand

Using this format, the speed of the link is Significand * 10 ^ Exponent kbps. This allows specifying link speeds with up to 3 decimal digits precision and speeds from 1 kbps to 10 $^{\wedge}$ 67 kbps. A computed speed of 0 kbps means the link speed is < 1 kbps.

Here are some examples how this is used:

Link Speed	Exponent	Significand
500 kbps	0	500
500 kbps	2	5
155 Mbps	3	155
40 Gpbs	6	40
100 Gpbs	6	100
100 Gpbs	8	1

3.2. Example message layouts

We will here give a few examples to illustrate the use of flags and options.

A minimum size message has no option flags set, and looks like this:

0	1	2	3
0 1 2 3 4 5	6 7 8 9 0 1 2 3 4 5	5 6 7 8 9 0 1 2 3	4 5 6 7 8 9 0 1
+-+-+-+-+-+	-+-+-+-+-+-	+-+-+-+-+-	+-+-+-+-+-+-+-+
F E Attr Ty	rpe Length = 6	Effectiv	/e MTU
+-+-+-+-+	-+-+-+-+-	+-+-+-+-+-	+
Unalloc/Re	served P a t A S	8 0 0 0 0 0 0 0	Unalloc/Rsrvd
+-+-+-+-+-+	-+-+-+-+-+-+-	+-+-+-+-+-+-	+-+-+-+-+-+-+

A message containing all the options defined in this document would look like this:

0	1	2	3			
0 1 2 3 4 5 6 7 8 9	9 0 1 2 3 4 5 6 7 8 9	0 1 2 3 4 5 6 7 8 9	0 1			
+-+-+-+-+-+-+-+-+-+-	+-+-+-+-	+-+-+-+-	+-+-+			
F E Attr Type	_ength = 18	Effective MTU				
+-+-+-+-+-+-+-+-+-+-+-+-+-+-	+-+-+-+-	+-+-+-+-	+-+-+			
Unalloc/Reserved	P a t A S 1 1 1 1	1 1 1 1 Unalloc/Rs	rvd			
+-+-+-+-+-+-+-+-+-+-+-+-+-	+-+-+-+-	+-+-+-+-	+-+-+			
Transit Oif-List Count						
+-						
Stub Oif-List Count						
+-+-+-+-+-+-+-+-+-+-+-+-+-	+-+-+-+-	+-+-+-+-	+-+-+			
Minimum Speed	d Link M	aximum Speed Link	- 1			
+-						
Domain Count 1	Node Count Diamet	er Count TZ Coun	t			
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-++	-+-+-+-+-+-+-+-+-	+-+-+-+-	+-+-+			

A message containing only Stub Oif-List Count and Node Count would look like this:

```
2
                          3
         1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
|F|E| Attr Type | Length = 9 | Effective MTU |
Unalloc/Reserved |P|a|t|A|S|0|1|0|0|0|1|0|0| Unalloc/Rsrvd |
Stub Oif-List Count
| Node count |
+-+-+-+-+-+-+
```

4. How to use Pop-Count Encoding

A router supporting this mechanism MUST, unless administratively disabled, include the PIM Join Attribute option in its PIM Hellos. See [RFC5384] and [HELLO] for details.

It is RECOMMENDED that implementations allow for administrative control whether to make use of this mechanism. Implementations MAY also allow further control of what information to store and send upstream.

It is very important to note that any changes to the values maintained by this mechanism MUST NOT trigger a new Join/Prune message. Due to the periodic nature of PIM, the values can be accurately obtained at 1 minute intervals (or whatever Join/Prune interval used).

When a router removes a link from an oif-list, it need to be able to reevaluate the values that it will advertise upstream. This happens when an oif-list entry is timed out or a Prune is received.

It is RECOMMENDED that the Join Attribute defined in this document be used only for entries in the join-list part of the Join/Prune message. If the attribute is used in the prune-list, an implementation MUST ignore it and process the Prune as if the attribute was not present.

It is also RECOMMENDED that join suppression be disabled on a LAN when Pop-Count is used.

It is RECOMMENDED that when triggered Join/Prune messages are sent by a downstream router, that the accounting information not be included in the message. This way when convergence is important, avoiding the processing time to build an accounting record in a downstream router and processing time to parse the message in the upstream router will help reduce convergence time. An upstream router SHOULD NOT interpret a Join/Prune message received with no accounting data to mean clearing or resetting what accounting data it has cached.

5. Implementation Approaches

This section offers some non-normative suggestions for how pop-count may be be implemented.

An implementation can decide how the accounting attributes are maintained. The values can be stored as part of the multicast route data structure by combining the local information it has with the joined information on a per oif basis. So when it is time to send a Join/Prune message, the values stored in the multicast route can be copied to the message.

Or, an implementation could store the accounting values per oif and when a Join/Prune message is sent, it can combine the oifs with its local information. Then the combined information can be copied to the message.

When a downstream joiner stops joining, accounting values cached must be evaluated. There are two approaches which can be taken. One is to keep values learned from each joiner so when the joiner goes away the count/max/min values are known and the combined value can be adjusted. The other approach is to set the value to 0 for the oif, and then start accumulating new values as subsequent Joins are received.

The same issue arises when an oif is removed from the oif-list. Keeping per-oif values allows you to adjust the per-route values when an oif goes away. Or, alternatively, a delay for reporting the new set a values from the route can occur while all oif values are zeroed (where accumulation of new values from subsequent Joins cause repopulation of values and a new max/min/count can be reevaluated for the route).

6. Caveats

This specification requires each router on a multicast distribution tree to support this specification or else the accounting attributes for the tree will not be known.

However, if there are a contiguous set of routers downstream in the distribution tree, they can maintain accounting information for the sub-tree.

If there are a set of contiguous routers supporting this specification upstream on the multicast distribution tree, accounting information will be available but it will not represent an accurate assessment of the entire tree. Also, it will not be clear for how much of the distribution tree the accounting information covers.

7. IANA Considerations

A new PIM Hello Option type, 29, has been assigned temporarily. The string TBD1 needs to be replaced with the permanently assigned value. See [HELLO] for details. Although the length is specified as 0 in this specifications, non-zero length is allowed, so IANA should list the length as being variable.

A new PIM Join Attribute type needs to be assigned. The string TBD2 needs to be replaced with the assigned value.

8. Security Considerations

The use of this specification requires some additional processing of PIM Join/Prune messages. However, the additional amount of processing is fairly limited, so this is not believed to be a significant concern.

The use of this mechanism includes information like the number of receivers. This information is assumed to not be of sensitive nature. If an operator has concerns about revealing this information to upstream routers, or other routers/hosts that may potentially inspect this information, there should be a way to disable the mechanism, or alternatively more detailed control of what information to include.

9. Acknowledgments

The authors would like to thank John Zwiebel, Amit Jain, and Clayton Wagar for their review comments on the initial versions of this document. Adrian Farrel did a detailed review of the document and proposed textual changes that have been incorporated. Further review and comments were provided by Thomas Morin and Zhaohui (Jeffrey) Zhang.

10. References

10.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4601] Fenner, B., Handley, M., Holbrook, H., and I. Kouvelas, "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)", RFC 4601, August 2006.
- [RFC5015] Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano, "Bidirectional Protocol Independent Multicast (BIDIR-PIM)", RFC 5015, October 2007.
- [RFC5384] Boers, A., Wijnands, I., and E. Rosen, "The Protocol Independent Multicast (PIM) Join Attribute Format", RFC 5384, November 2008.

10.2. Informative References

- [HELLO] IANA, "PIM-Hello Options", <http://www.iana.org/assignments/pim-parameters>.
- [I-D.ietf-mboned-auto-multicast] Bumgardner, G., "Automatic Multicast Tunneling", draft-ietf-mboned-auto-multicast-14 (work in progress), June 2012.

Authors' Addresses

Dino Farinacci Cisco Systems Tasman Drive San Jose, CA 95134 USA

Email: dino@cisco.com

Greg Shepherd Cisco Systems Tasman Drive San Jose, CA 95134 USA

Email: gjshep@gmail.com

Stig Venaas Cisco Systems Tasman Drive San Jose, CA 95134 USA

Email: stig@cisco.com

Yiqun Cai Microsoft 1065 La Avenida Mountain View, CA 94043 USA

Email: yiqunc@microsoft.com