

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: May 4, 2017

IJ. Wijnands
S. Venaas
Cisco Systems, Inc.
M. Brig
Aegis BMD Program Office
A. Jonasson
Swedish Defence Material Administration (FMV)
October 31, 2016

PIM flooding mechanism and source discovery
draft-ietf-pim-source-discovery-bsr-05

Abstract

PIM Sparse-Mode uses a Rendezvous Point and shared trees to forward multicast packets from new sources. Once last hop routers receive packets from a new source, they may join the Shortest Path Tree for the source for optimal forwarding. This draft defines a new protocol that provides a way to support PIM Sparse Mode (SM) without the need for PIM registers, RPs or shared trees. Multicast source information is flooded throughout the multicast domain using a new generic PIM flooding mechanism. This allows last hop routers to learn about new sources without receiving initial data packets.

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 May 4, 2017.

Copyright Notice

Copyright (c) 2016 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.	Conventions used in this document	3
1.2.	Terminology	3
2.	Testing and deployment experiences	3
3.	A generic PIM flooding mechanism	4
3.1.	PFM message format	4
3.2.	Processing PFM messages	5
3.2.1.	Initial checks	5
3.2.2.	Processing and forwarding of PFM messages	6
4.	Distributing Source to Group Mappings	6
4.1.	Group Source Holdtime TLV	6
4.2.	Originating PFM messages	7
4.3.	Processing GSH TLVs	8
4.4.	The first packets and bursty sources	8
4.5.	Resiliency to network partitioning	9
5.	Security Considerations	10
6.	IANA considerations	10
7.	Acknowledgments	10
8.	References	10
8.1.	Normative References	10
8.2.	Informative References	11
	Authors' Addresses	11

[1.](#) Introduction

PIM Sparse-Mode uses a Rendezvous Point (RP) and shared trees to forward multicast packets to Last Hop Routers (LHR). After the first packet is received by a LHR, the source of the multicast stream is learned and the Shortest Path Tree (SPT) can be joined. This draft defines a new mechanism that provides a way to support PIM Sparse Mode (SM) without the need for PIM registers, RPs or shared trees. Multicast source information is flooded throughout the multicast domain using a new generic PIM flooding mechanism. This mechanism is defined in this document, and is modeled after the Bootstrap Router mechanism [[RFC5059](#)]. By removing the need for RPs and shared trees, the PIM-SM procedures are simplified, improving router operations,

management and making the protocol more robust. Also the data packets are only sent on the SPTs, providing optimal forwarding.

1.1. 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](#)].

1.2. Terminology

RP: Rendezvous Point.

BSR: Bootstrap Router.

RPF: Reverse Path Forwarding.

SPT: Shortest Path Tree.

FHR: First Hop Router, directly connected to the source.

LHR: Last Hop Router, directly connected to the receiver.

PFM: PIM Flooding Mechanism.

PFM-SA: PFM Source Announcement.

SG Mapping: Multicast source to group mapping.

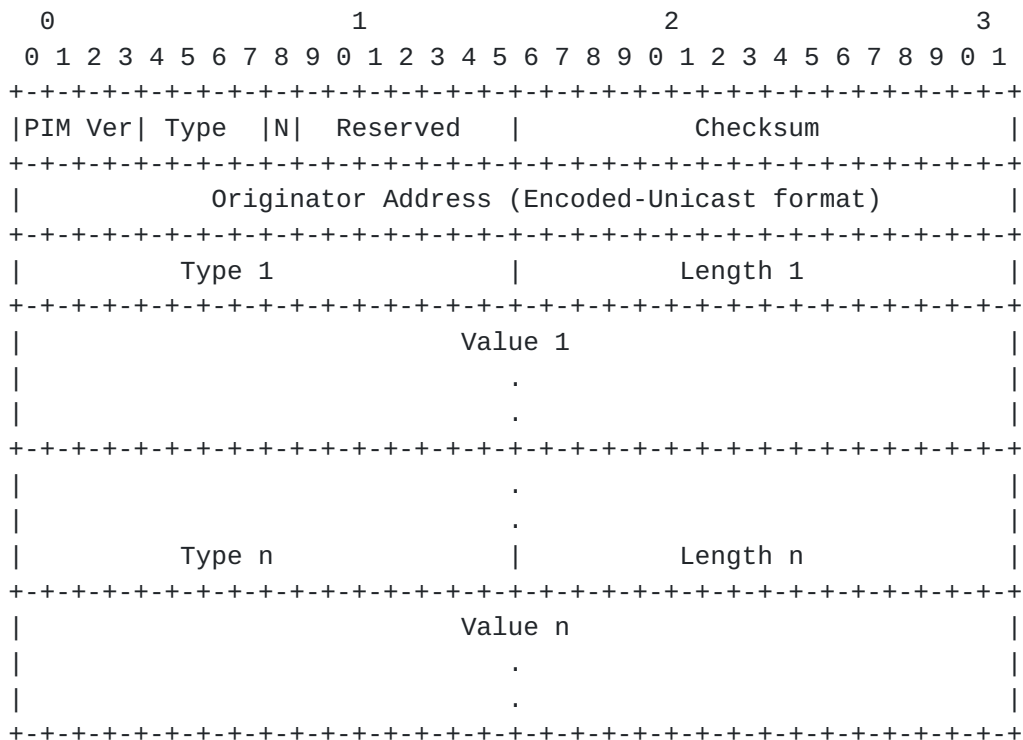
2. Testing and deployment experiences

A prototype of this specification has been implemented and there has been some limited testing in the field. The prototype was tested in a network with low bandwidth radio links. In this network with frequent topology changes and link or router failures, PIM-SM with RP election is found to be too slow. With PIM-DM, issues were observed with new multicast sources starving low bandwidth links even when there are no receivers, in some cases such that there was no bandwidth left for prune message. For the tests, all routers were configured to send PFM-SA for directly connected source and to cache received announcements. Applications such as SIP with multicast subscriber discovery, multicast voice conferencing, position tracking and NTP were successfully tested. The tests went quite well. Packets were rerouted as needed and there were no unnecessary forwarding of packets. Ease of configuration was seen as a plus.

3. A generic PIM flooding mechanism

The Bootstrap Router mechanism (BSR) [[RFC5059](#)] is a commonly used mechanism for distributing dynamic Group to RP mappings in PIM. It is responsible for flooding information about such mappings throughout a PIM domain, so that all routers in the domain can have the same information. BSR as defined, is only able to distribute Group to RP mappings. We are defining a more generic mechanism that can flood any kind of information throughout a PIM domain. It is not necessarily a domain though, it depends on the administrative boundaries being configured. The forwarding rules are identical to BSR, except that there is no BSR election and that one can control whether routers should forward unsupported data types. For some types of information it is quite useful that it can be distributed without all routers having to support the particular type, while there may also be types where it is necessary for every single router to support it. The mechanism includes an originator address which is used for RPF checking to restrict the flooding, and prevent loops, just like BSR. Just like BSR it is also sent hop by hop. Note that there is no built in election mechanism as in BSR, so there can be multiple originators. We call this mechanism the PIM Flooding Mechanism (PFM).

3.1. PFM message format



PIM Version: Reserved, Checksum Described in [\[RFC7761\]](#).

Type: PIM Message Type. Value (pending IANA) for a PFM message.

[N]o-Forward bit: When set, this bit means that the PFM message is not to be forwarded.

Originator Address: The address of the router that originated the message. This can be any address assigned to the originating router, but MUST be routable in the domain to allow successful forwarding. The format for this address is given in the Encoded-Unicast address in [[RFC7761](#)].

Type 1..n: A message contains one or more TLVs, in this case n TLVs. The Type specifies what kind of information is in the Value.

Length 1..n: The length of the the value field.

Value 1..n: The value associated with the type and of the specified length.

[3.2.](#) Processing PFM messages

A router that receives a PFM message MUST perform the initial checks specified here. If the checks fail, the message MUST be dropped. An error MAY be logged, but otherwise the message MUST be dropped silently. If the checks pass, the contents is processed according to the processing rules of the included TLVs.

[3.2.1.](#) Initial checks

In order to do further processing, a message MUST meet the following requirements. The message MUST be from a directly connected neighbor for which we have active Hello state, and it MUST have been sent to the ALL-PIM-ROUTERS group. Also, the interface MUST NOT be an administrative boundary for PFM. If No-Forward is not set, it MUST have been sent by the RPF neighbor for the originator address. If No-Forward is set, we MUST have restarted within 60 seconds. In pseudo-code the algorithm is as follows:


```
if ((DirectlyConnected(PFM.src_ip_address) == FALSE) OR
    (we have no Hello state for PFM.src_ip_address) OR
    (PFM.dst_ip_address != ALL-PIM-ROUTERS) OR
    (Incoming interface is admin boundary for PFM)) {
    drop the message silently, optionally log error.
}
if (PFM.no_forward_bit == 0) {
    if (PFM.src_ip_address !=
        RPF_neighbor(PFM.originator_ip_address)) {
        drop the message silently, optionally log error.
    }
} else if (more than 60 seconds elapsed since startup)) {
    drop the message silently, optionally log error.
}
```

3.2.2. Processing and forwarding of PFM messages

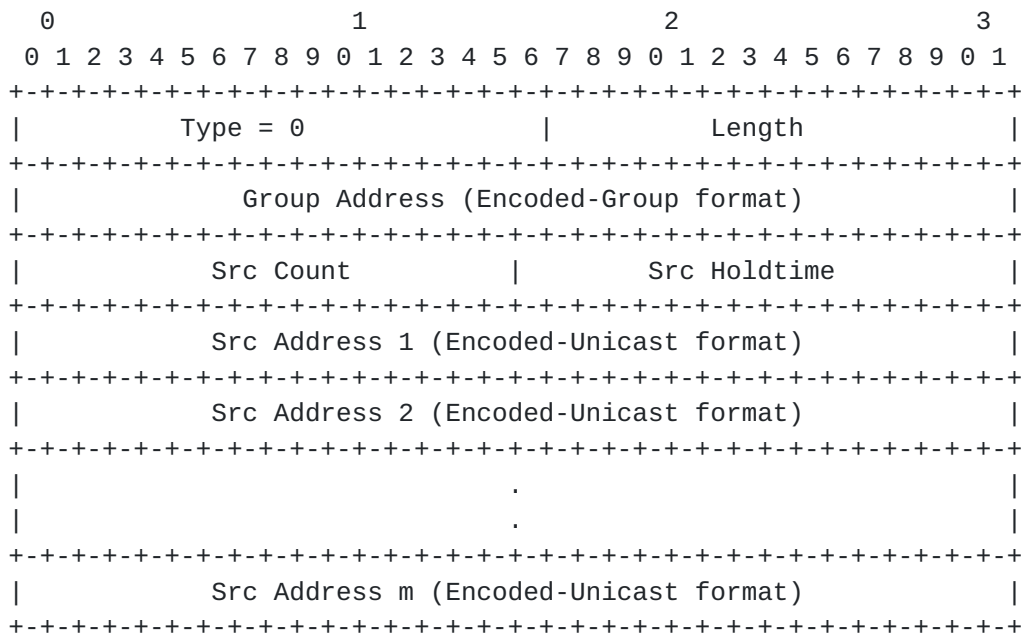
When the message is received, the initial checks above must be performed. If it passes the checks, we then for each included TLV perform processing according to the specification for that TLV.

After processing, we forward the message. Unless otherwise specified by the type specification, the TLVs in the forwarded message are identical to the TLVs in the received message. However, if the most significant bit in the type field is set (the type value is larger than 32767) and we do not support the type, then that particular type should be omitted from the forwarded messages. The message is forwarded out of all interfaces with PIM neighbors (including the interface it was received on).

4. Distributing Source to Group Mappings

The generic flooding mechanism (PFM) defined in the previous section can be used for distributing source to group mappings about active multicast sources throughout a PIM domain. A Group Source Holdtime (GSH) TLV is defined for this purpose.

4.1. Group Source Holdtime TLV



Type: This TLV has type 0.

Length: The length of the value.

Group Address: The group we are announcing sources for. The format for this address is given in the Encoded-Group format in [\[RFC7761\]](#).

Src Count: How many unicast encoded sources address encodings follow.

Src Holdtime: The Holdtime (in seconds) for the corresponding source(s).

Src Address: The source address for the corresponding group. The format for these addresses is given in the Encoded-Unicast address in [\[RFC7761\]](#).

[4.2. Originating PFM messages](#)

A PFM message MAY contain one or more Group Source Holdtime (GSH) TLVs. This is used to flood information about active multicast sources. Each FHR that is directly connected to an active multicast source originates PFM messages containing GSH TLVs. How a multicast router discovers the source of the multicast packet and when it considers itself the FHR follows the same procedures as the registering process described in [\[RFC7761\]](#). When a FHR has decided that a register needs to be sent per [\[RFC7761\]](#), the SG is not registered via the PIM SM register procedures, but the SG mapping is

included in an GSH TLV in a PFM message. Note, only the SG mapping is distributed in the message, not the entire packet as would have been done with a PIM register. The router originating the PFM messages includes one of its own addresses in the originator field. Note that this address SHOULD be routeable due to RPF checking. The PFM messages containing the GSH TLV are periodically sent for as long as the multicast source is active, similar to how PIM registers are periodically sent. The default announcement period is 60 seconds, which means that as long as the source is active, it is included in a PFM message originated every 60 seconds. The holdtime for the source is by default 210 seconds. Other values MAY be configured, but the holdtime MUST be either zero, or larger than the announcement period. It is RECOMMENDED to be 3.5 times the announcement period. A source MAY be announced with a holdtime of zero to indicate that the source is no longer active.

If an implementation supports originating GSH TLVs with different holdtimes for different sources, it can if needed send multiple TLVs with the same group address. Due to the format, all the sources in the same TLV have the same holdtime.

4.3. Processing GSH TLVs

A router that receives a PFM message containing GSH TLVs SHOULD parse the message and store each of the GSH TLVs as SG mappings with a holdtimer started with the advertised holdtime. For each group that has directly connected receivers, this router SHOULD send PIM (S,G) joins for all the SG mappings advertised in the message for the group. The SG mappings are kept alive for as long as the holdtimer for the source is running. Once the holdtimer expires a PIM router MAY send a PIM (S,G) prune to remove itself from the tree. However, when this happens, there should be no more packets sent by the source, so it may be desirable to allow the state to time out rather than sending a prune.

Note that a holdtime of zero has a special meaning. It is to be treated as if the source just expired, and state to be removed. Source information MUST NOT be removed due to the source being omitted in a message. For instance, if there is a large number of sources for a group, there may be multiple PFM messages, each message containing a different list of sources for the group.

4.4. The first packets and bursty sources

The PIM register procedure is designed to deliver Multicast packets to the RP in the absence of a Shortest Path Tree (SPT) from the RP to the source. The register packets received on the RP are decapsulated and forwarded down the shared tree to the LHRs. As soon as an SPT is

built, multicast packets would flow natively over the SPT to the RP or LHR and the register process would stop. The PIM register process ensures packet delivery until an SPT is in place reaching the FHR. If the packets were not unicast encapsulated to the RP they would be dropped by the FHR until the SPT is setup. This functionality is important for applications where the initial packet(s) must be received for the application to work correctly. Another reason would be for bursty sources. If the application sends out a multicast packet every 4 minutes (or longer), the SPT is torn down (typically after 3:30 minutes of inactivity) before the next packet is forwarded down the tree. This will cause no multicast packet to ever be forwarded. A well behaved application should be able to deal with packet loss since IP is a best effort based packet delivery system. But in reality this is not always the case.

With the procedures defined in this document the packet(s) received by the FHR will be dropped until the LHR has learned about the source and the SPT is built. That means for bursty sources or applications sensitive for the delivery of the first packet this solution would not be very applicable. This solution is mostly useful for applications that don't have strong dependency on the initial packet(s) and have a fairly constant data rate, like video distribution for example. For applications with strong dependency on the initial packet(s) we recommend using PIM Bidir [[RFC5015](#)] or SSM [[RFC4607](#)]. The protocol operations are much simpler compared to PIM SM, it will cause less churn in the network and both guarantee best effort delivery for the initial packet(s).

Another solution to address the problems described above is documented in [[I-D.ietf-magma-msnip](#)]. This proposal allows for a host to tell the FHR its willingness to act as Source for a certain Group before sending the data packets. LHRs have time to join the SPT before the host starts sending which would avoid packet loss. The SG mappings announced by [[I-D.ietf-magma-msnip](#)] can be advertised directly in SG messages, allowing a nice integration of both proposals. The life time of the SPT is not driven by the liveliness of Multicast data packets (which is the case with PIM SM), but by the announcements driven via [[I-D.ietf-magma-msnip](#)]. This will also prevent packet loss due to bursty sources.

4.5. Resiliency to network partitioning

In a PIM SM deployment where the network becomes partitioned, due to link or node failure, it is possible that the RP becomes unreachable to a certain part of the network. New sources that become active in that partition will not be able to register to the RP and receivers within that partition are not able to receive the traffic. Ideally you would want to have a candidate RP in each partition, but you

never know in advance which routers will form a partitioned network. In order to be fully resilient, each router in the network may end up being a candidate RP. This would increase the operational complexity of the network.

The solution described in this document does not suffer from that problem. If a network becomes partitioned and new sources become active, the receivers in that partitioned will receive the SG Mappings and join the source tree. Each partition works independently of the other partition(s) and will continue to have access to sources within that partition. As soon as the network heals, the SG Mappings are re-flooded into the other partition(s) and other receivers can join to the newly learned sources.

5. Security Considerations

The security considerations are mainly similar to what is documented in [[RFC5059](#)]. It is a concern that rogue devices can inject packets that are flooded throughout a domain. PFM packets must only be accepted from a PIM neighbor. Deployments may use mechanisms for authenticating PIM neighbors. For PFM-SA it is an issue that injected packets from a rogue device could send SG mappings for a large number of source addresses, causing routers to use memory storing these mappings, and also if they have interest in the groups, build Shortest Path Trees for sources that are not actually active.

6. IANA considerations

This document requires the assignment of a new PIM message type for the PIM Flooding Mechanism (PFM). IANA is also requested to create a registry for PFM TLVs, with type 0 assigned to the "Source Group Holdtime" TLV. Values in the range 1-65535 are "Unassigned". Assignments for the registry are to be made according to the policy "IETF Review" as defined in [[RFC5226](#)].

7. Acknowledgments

The authors would like to thank Arjen Boers for contributing to the initial idea, and Yiqun Cai and Dino Farinacci for their comments on the draft.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), DOI 10.17487/RFC2119, March 1997, <<http://www.rfc-editor.org/info/rfc2119>>.
- [RFC5059] Bhaskar, N., Gall, A., Lingard, J., and S. Venaas, "Bootstrap Router (BSR) Mechanism for Protocol Independent Multicast (PIM)", [RFC 5059](#), DOI 10.17487/RFC5059, January 2008, <<http://www.rfc-editor.org/info/rfc5059>>.
- [RFC7761] Fenner, B., Handley, M., Holbrook, H., Kouvelas, I., Parekh, R., Zhang, Z., and L. Zheng, "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)", STD 83, [RFC 7761](#), DOI 10.17487/RFC7761, March 2016, <<http://www.rfc-editor.org/info/rfc7761>>.

8.2. Informative References

- [RFC4607] Holbrook, H. and B. Cain, "Source-Specific Multicast for IP", [RFC 4607](#), DOI 10.17487/RFC4607, August 2006, <<http://www.rfc-editor.org/info/rfc4607>>.
- [RFC5015] Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano, "Bidirectional Protocol Independent Multicast (BIDIR-PIM)", [RFC 5015](#), DOI 10.17487/RFC5015, October 2007, <<http://www.rfc-editor.org/info/rfc5015>>.
- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", [BCP 26](#), [RFC 5226](#), DOI 10.17487/RFC5226, May 2008, <<http://www.rfc-editor.org/info/rfc5226>>.
- [I-D.ietf-magma-msnip] Fenner, B., Haberman, B., Holbrook, H., Kouvelas, I., and S. Venaas, "Multicast Source Notification of Interest Protocol (MSNIP)", [draft-ietf-magma-msnip-06](#) (work in progress), March 2011.

Authors' Addresses

IJsbrand Wijnands
Cisco Systems, Inc.
De kleetlaan 6a
Diegem 1831
Belgium

Email: ice@cisco.com

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

Email: stig@cisco.com

Michael Brig
Aegis BMD Program Office
17211 Avenue D, Suite 160
Dahlgren VA 22448-5148
USA

Email: michael.brig@mda.mil

Anders Jonasson
Swedish Defence Material Administration (FMV)
Loennvaegen 4
Vaexjoe 35243
Sweden

Email: anders@jomac.se

