

Network Working Group  
Internet-Draft  
Intended status: Standards Track  
Expires: April 25, 2019

Y. Cai  
H. Ou  
Alibaba Group  
S. Vallepalli  
M. Mishra  
S. Venaas  
Cisco Systems, Inc.  
A. Green  
British Telecom  
October 22, 2018

**PIM Designated Router Load Balancing**  
**draft-ietf-pim-drlb-09**

Abstract

On a multi-access network, one of the PIM routers is elected as a Designated Router (DR). On the last hop LAN, the PIM DR is responsible for tracking local multicast listeners and forwarding traffic to these listeners if the group is operating in PIM-SM. This document specifies a modification to the PIM-SM protocol that allows more than one of these last hop routers to be selected, so that the forwarding load can be distributed among these routers.

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 April 25, 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 (<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 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

<a href="#">1.</a>	<a href="#">Introduction</a>	<a href="#">2</a>
<a href="#">2.</a>	<a href="#">Terminology</a>	<a href="#">5</a>
<a href="#">3.</a>	<a href="#">Applicability</a>	<a href="#">5</a>
<a href="#">4.</a>	<a href="#">Functional Overview</a>	<a href="#">6</a>
<a href="#">4.1.</a>	<a href="#">GDR Candidates</a>	<a href="#">6</a>
<a href="#">4.2.</a>	<a href="#">Hash Mask and Hash Algorithm</a>	<a href="#">7</a>
<a href="#">4.3.</a>	<a href="#">Modulo Hash Algorithm</a>	<a href="#">8</a>
<a href="#">4.3.1.</a>	<a href="#">Limitations</a>	<a href="#">9</a>
<a href="#">4.4.</a>	<a href="#">PIM Hello Options</a>	<a href="#">9</a>
<a href="#">5.</a>	<a href="#">Hello Option Formats</a>	<a href="#">10</a>
<a href="#">5.1.</a>	<a href="#">PIM DR Load Balancing Capability (DRLBC) Hello Option</a>	<a href="#">10</a>
<a href="#">5.2.</a>	<a href="#">PIM DR Load Balancing GDR (DRLBGDR) Hello Option</a>	<a href="#">10</a>
<a href="#">6.</a>	<a href="#">Protocol Specification</a>	<a href="#">11</a>
<a href="#">6.1.</a>	<a href="#">PIM DR Operation</a>	<a href="#">11</a>
<a href="#">6.2.</a>	<a href="#">PIM GDR Candidate Operation</a>	<a href="#">12</a>
<a href="#">6.2.1.</a>	<a href="#">Router Receives New DRLBGDR</a>	<a href="#">13</a>
<a href="#">6.2.2.</a>	<a href="#">Router Receives Updated DRLBGDR</a>	<a href="#">13</a>
<a href="#">6.3.</a>	<a href="#">PIM Assert Modification</a>	<a href="#">14</a>
<a href="#">7.</a>	<a href="#">Compatibility</a>	<a href="#">15</a>
<a href="#">8.</a>	<a href="#">Manageability Considerations</a>	<a href="#">16</a>
<a href="#">9.</a>	<a href="#">IANA Considerations</a>	<a href="#">16</a>
<a href="#">9.1.</a>	<a href="#">Initial registry</a>	<a href="#">16</a>
<a href="#">9.2.</a>	<a href="#">Assignment of new message types</a>	<a href="#">16</a>
<a href="#">10.</a>	<a href="#">Security Considerations</a>	<a href="#">16</a>
<a href="#">11.</a>	<a href="#">Acknowledgement</a>	<a href="#">17</a>
<a href="#">12.</a>	<a href="#">References</a>	<a href="#">17</a>
<a href="#">12.1.</a>	<a href="#">Normative References</a>	<a href="#">17</a>
<a href="#">12.2.</a>	<a href="#">Informative References</a>	<a href="#">17</a>
	<a href="#">Authors' Addresses</a>	<a href="#">17</a>

## [1.](#) Introduction

On a multi-access LAN such as an Ethernet, one of the PIM routers is elected as a DR. The PIM DR has two roles in the PIM-SM protocol. On the first hop LAN, the PIM DR is responsible for registering an active source with the Rendezvous Point (RP) if the group is



operating in PIM-SM. On the last hop LAN, the PIM DR is responsible for tracking local multicast listeners and forwarding to these listeners if the group is operating in PIM-SM.

Consider the following last hop LAN in Figure 1:

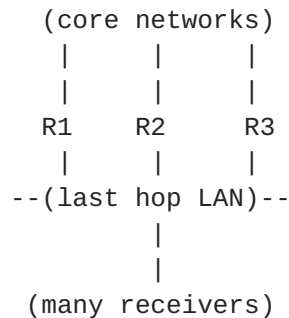


Figure 1: Last Hop LAN

Assume R1 is elected as the Designated Router. According to [\[RFC7761\]](#), R1 will be responsible for forwarding traffic to that LAN on behalf of any local members. In addition to keeping track of IGMP and MLD membership reports, R1 is also responsible for initiating the creation of source and/or shared trees towards the senders or the RPs.

Forcing sole data plane forwarding responsibility on the PIM DR uncovers a limitation in the protocol. In comparison, even though an OSPF DR or an IS-IS DIS handles additional duties while running the OSPF or IS-IS protocols, they are not required to be solely responsible for forwarding packets for the network. On the other hand, on a last hop LAN, only the PIM DR is asked to forward packets while the other routers handle only control traffic (and perhaps drop packets due to RPF failures). Hence the forwarding load of a last hop LAN is concentrated on a single router.

This leads to several issues. One of the issues is that the aggregated bandwidth will be limited to what R1 can handle towards this particular interface. It is very common that the last hop LAN consists of switches that run IGMP/MLD or PIM snooping. This allows the forwarding of multicast packets to be restricted only to segments leading to receivers who have indicated their interest in multicast groups using either IGMP or MLD. The emergence of the switched Ethernet allows the aggregated bandwidth to exceed, sometimes by a large number, that of a single link. For example, let us modify Figure 1 and introduce an Ethernet switch in Figure 2.



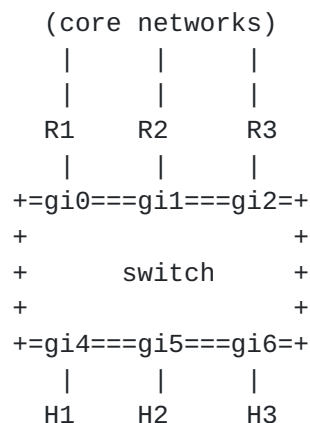


Figure 2: Last Hop Network with Ethernet Switch

Let us assume that each individual link is a Gigabit Ethernet. Each router, R1, R2 and R3, and the switch have enough forwarding capacity to handle hundreds of Gigabits of data.

Let us further assume that each of the hosts requests 500 Mbps of unique multicast data. This totals to 1.5 Gbps of data, which is less than what each switch or the combined uplink bandwidth across the routers can handle, even under failure of a single router.

On the other hand, the link between R1 and switch, via port gi0, can only handle a throughput of 1Gbps. And if R1 is the only DR (the PIM DR elected using the procedure defined by [RFC7761](#)) at least 500 Mbps worth of data will be lost because the only link that can be used to draw the traffic from the routers to the switch is via gi0. In other words, the entire network's throughput is limited by the single connection between the PIM DR and the switch (or the last hop LAN as in Figure 1).

Another important issue is related to failover. If R1 is the only forwarder on the last hop router for a shared LAN, when R1 goes out of service, multicast forwarding for the entire LAN has to be rebuilt by the newly elected PIM DR. However, if there was a way that allowed multiple routers to forward to the LAN for different groups, failure of one of the routers would only lead to disruption to a subset of the flows, therefore improving the overall resilience of the network.



There is a limitation in the hash algorithm used in this document, but this document provides the option to have different and more consistent hash algorithms in the future.

This document specifies a modification to the PIM-SM protocol that allows more than one of these routers, called Group Designated Routers (GDR) to be selected so that the forwarding load can be distributed among a number of routers.

## **2. Terminology**

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](#)].

With respect to PIM, this document follows the terminology that has been defined in [[RFC7761](#)].

This document also introduces the following new acronyms:

- o GDR: GDR stands for "Group Designated Router". For each multicast flow, either a (\*,G) for ASM, or an (S,G) for SSM, a hash algorithm (described below) is used to select one of the routers as a GDR. The GDR is responsible for initiating the forwarding tree building process for the corresponding multicast flow.
- o GDR Candidate: a last hop router that has the potential to become a GDR. A GDR Candidate must have the same DR priority and must run the same GDR election hash algorithm as the DR router. It must send and process new PIM Hello Options as defined in this document. There might be more than one GDR Candidate on a LAN, but only one can become GDR for a specific multicast flow.

## **3. Applicability**

The extension specified in this document applies to PIM-SM last hop routers only. It does not alter the behavior of a PIM DR on the first hop network. This is because the source tree is built using the IP address of the sender, not the IP address of the PIM DR that sends the registers towards the RP. The load balancing between first hop routers can be achieved naturally if an IGP provides equal cost multiple paths (which it usually does in practice). Also distributing the load to do registering does not justify the additional complexity required to support it.





## **4. Functional Overview**

In the PIM DR election as defined in [RFC7761], when multiple last hop routers are connected to a multi-access LAN (for example, an Ethernet), one of them is elected to act as PIM DR. The PIM DR is responsible for sending local Join/Prune messages towards the RP or source. In order to elect the PIM DR, each PIM router on the LAN examines the received PIM Hello messages and compares its own DR priority and IP address with those of its neighbors. The router with the highest DR priority is the PIM DR. If there are multiple such routers, their IP addresses are used as the tie-breaker, as described in [RFC7761].

In order to share forwarding load among last hop routers, besides the normal PIM DR election, the GDR is also elected on the last hop multi-access LAN. There is only one PIM DR on the multi-access LAN, but there might be multiple GDR Candidates.

For each multicast flow, that is, (\*,G) for ASM and (S,G) for SSM, a hash algorithm is used to select one of the routers to be the GDR. A new DR Load Balancing Capability (DRLBC) PIM Hello Option, which contains hash algorithm type, is announced by routers on interfaces where this specification is enabled. Last hop routers with the new DRLBC Option advertised in its Hello, and using the same GDR election hash algorithm and the same DR priority as the PIM DR, are considered as GDR Candidates.

Hash Masks are defined for Source, Group and RP separately, in order to handle PIM ASM/SSM. The masks, as well as a sorted list of GDR Candidate Addresses, are announced by the DR in a new DR Load Balancing GDR (DRLBGDR) PIM Hello Option.

A hash algorithm based on the announced Source, Group, or RP masks allows one GDR to be assigned to a corresponding multicast state. And that GDR is responsible for initiating the creation of the multicast forwarding tree for multicast traffic.

### **4.1. GDR Candidates**

GDR is the new concept introduced by this specification. GDR Candidates are routers eligible for GDR election on the LAN. To become a GDR Candidate, a router MUST support this specification, have the same DR priority and run the same GDR election hash algorithm as the DR on the LAN.

For example, assume there are 4 routers on the LAN: R1, R2, R3 and R4, which all support this specification. R1, R2 and R3 have the same DR priority while R4's DR priority is less preferred. In this



example, R4 will not be eligible for GDR election, because R4 will not become a PIM DR unless all of R1, R2 and R3 go out of service.

Furthermore, assume router R1 wins the PIM DR election, R1 and R2 run the same hash algorithm for GDR election, while R3 runs a different one. In this case, only R1 and R2 will be eligible for GDR election, while R3 will not.

As a DR, R1 will include its own Load Balancing Hash Masks and the identity of R1 and R2 (the GDR Candidates) in its DRLBGDR Hello Option.

#### **4.2. Hash Mask and Hash Algorithm**

A Hash Mask is used to extract a number of bits from the corresponding IP address field (32 for v4, 128 for v6) and calculate a hash value. A hash value is used to select a GDR from GDR Candidates advertised by PIM DR. For example, 0.0.255.0 defines a Hash Mask for an IPv4 address that masks the first, the second, and the fourth octets.

There are three Hash Masks defined:

- o RP Hash Mask
- o Source Hash Mask
- o Group Hash Mask

The hash masks need to be configured on the PIM routers that can potentially become a PIM DR, unless the implementation provides default Hash Mask values. An implementation SHOULD provide masks with default values 255.255.255.255 (IPv4) and FFFF:FFFF:FFFF:FFFF:FFFF:FFFF:FFFF:FFFF (IPv6).

- o If the group is in ASM mode and the RP Hash Mask announced by the PIM DR is not 0, calculate the value of hashvalue\_RP [[Section 4.3](#)] to determine GDR.
- o If the group is in ASM mode and the RP Hash Mask announced by the PIM DR is 0, obtain the value of hashvalue\_Group [[Section 4.3](#)] to determine GDR.
- o If the group is in SSM mode, use hashvalue\_SG [[Section 4.3](#)] to determine GDR.

A simple Modulo hash algorithm is defined in this document. However, to allow another hash algorithms to be used, a 1-octet "Hash



Algorithm Type" field is included in DRLBC Hello Option to specify the hash algorithm used by a last hop router.

If different hash algorithm types are advertised among last hop routers, only last hop routers running the same hash algorithm as the DR (and having the same DR priority as the DR) are eligible for GDR election.

#### [4.3.](#) Modulo Hash Algorithm

The Modulo hash algorithm is discussed here with a detailed description on `hashvalue_RP`. The same algorithm is described in brief for `hashvalue_Group` using the group address instead of the RP address for an ASM group with zero `RP_hashmask`, and also with `hashvalue_SG` for a the source address of an (S,G), instead of the RP address,

- o For ASM groups, with a non-zero `RP_Hash Mask`, hash value is calculated as:

$$\text{hashvalue\_RP} = (((\text{RP\_address} \& \text{RP\_hashmask}) \gg N) \& 0\text{xFFFF}) \% M$$

`RP_address` is the address of the RP defined for the group. `N` is the number of zeroes, counted from the least significant bit of the `RP_hashmask`. `M` is the number of GDR Candidates.

For example, Router X with IPv4 address 203.0.113.1 receives a DRLBGDR Hello Option from the DR, which announces RP Hash Mask 0.0.255.0 and a list of GDR Candidates, sorted by IP addresses from high to low: 203.0.113.3, 203.0.113.2 and 203.0.113.1. The ordinal number assigned to those addresses would be:

0 for 203.0.113.3; 1 for 203.0.113.2; 2 for 203.0.113.1 (Router X)

Assume there are 2 RPs: RP1 192.0.2.1 for Group1 and RP2 198.51.100.2 for Group2. Following the modulo hash algorithm:

`N` is 8 for 0.0.255.0, and `M` is 3 for the total number of GDR Candidates. The `hashvalue_RP` for RP1 192.0.2.1 is:

$$(((192.0.2.1 \& 0.0.255.0) \gg 8) \& 0\text{xFFFF} \% 3) = 2 \% 3 = 2$$

matches the ordinal number assigned to Router X. Router X will be the GDR for Group1, which uses 192.0.2.1 as the RP.

The `hashvalue_RP` for RP2 198.51.100.2 is:



$$(((198.51.100.2 \& 0.0.255.0) \gg 8) \& 0xFFFF \% 3) = 100 \% 3 = 1$$

which is different from Router X's ordinal number(2) hence,  
Router X will not be GDR for Group2.

- o If RP\_hashmask is 0, a hash value for an ASM group is calculated using the Group Hash Mask:

$$\text{hashvalue\_Group} = (((\text{Group\_address} \& \text{Group\_hashmask}) \gg N) \& 0xFFFF) \% M$$

Compare hashvalue\_Group with Ordinal number assigned to Router X, to decide if Router X is the GDR.

- o For SSM groups, a hash value is calculated using both the Source and Group Hash Mask:

$$\text{hashvalue\_SG} = (((\text{Source\_address} \& \text{Source\_hashmask}) \gg N\_S) \& 0xFFFF) \wedge (((\text{Group\_address} \& \text{Group\_hashmask}) \gg N\_G) \& 0xFFFF) \% M$$

#### **4.3.1. Limitations**

The Modulo Hash Algorithm has poor failover characteristics when a shared LAN has more than two GDRs. In the case of more than two GDRs on a LAN, when one GDR fails, all of the groups may be reassigned to a new GDR, even if they were not assigned to the failed GDR. However, many deployments use only two routers on a shared LAN for redundancy purposes. Future work may define new hash algorithms where only groups assigned to the failed GDR get reassigned.

#### **4.4. PIM Hello Options**

When a last hop PIM router sends a PIM Hello for an interface with this specification enabled, it includes a new option, called "Load Balancing Capability (DRLBC)".

Besides this DRLBC Hello Option, the elected PIM DR also includes a new "DR Load Balancing GDR (DRLBGDR) Hello Option". The DRLBGDR Hello Option consists of three Hash Masks as defined above and also a sorted list of GDR Candidate addresses on the last hop LAN.

The elected PIM DR uses DRLBC Hello Option advertised by all routers on the last hop LAN to compose the DRLBGDR Option. The GDR Candidates use the DRLBGDR Hello Option advertised by the PIM DR to calculate the hash value.





## 5. Hello Option Formats

### 5.1. PIM DR Load Balancing Capability (DRLBC) Hello Option

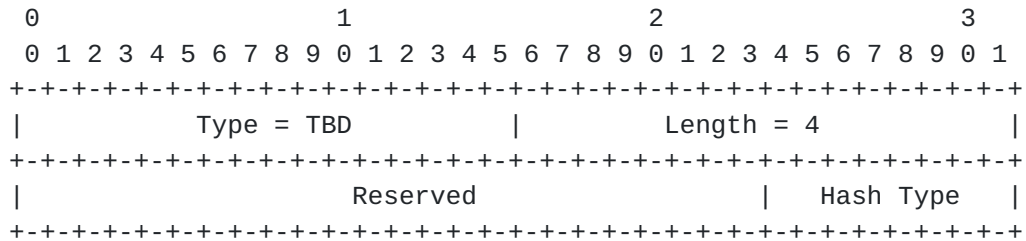


Figure 3: Capability Hello Option

Type: TBD

Length: 4

Hash Algorithm Type: 0 for Modulo hash algorithm

This DRLBC Hello Option **MUST** be advertised by last hop routers on interfaces with this specification enabled.

## 5.2. PIM DR Load Balancing GDR (DRLBGDR) Hello Option

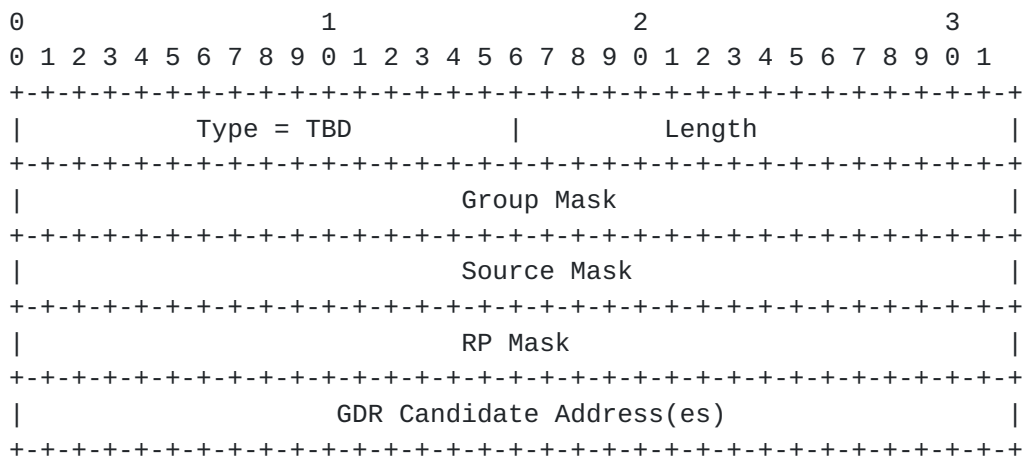


Figure 4: GDR Hello Option

Type: TBD

Length:  $(3 + n) \times (4 \text{ or } 16)$  where  $n$  is the number of GDR candidates.



Group Mask (32/128 bits): Mask

Source Mask (32/128 bits): Mask

RP Mask (32/128 bits): Mask

All masks MUST be in the same address family as the Hello IP header.

GDR Address (32/128 bits): Address(es) of GDR Candidate(s)

All addresses must be in the same address family as the Hello IP header. The addresses are sorted in descending order. The order is converted to the ordinal number associated with each GDR candidate in hash value calculation. For example, if addresses advertised are R3, R2, R1, the ordinal number assigned to R3 is 0, to R2 is 1 and to R1 is 2.

If the "Interface ID" option, as specified in [[RFC6395](#)], is present in a GDR Candidate's PIM Hello message, and the "Router ID" portion is non-zero:

- + For IPv4, the "GDR Candidate Address" will be set directly to the "Router ID".
- + For IPv6, the "GDR Candidate Address" will be set to the IPv4-IPv6 translated address of the "Router ID", as described in [[RFC4291](#)], that is the "Router-ID" is appended to the prefix of 96 bits of zeroes.

If the "Interface ID" option is not present in a GDR Candidate's PIM Hello message, or if the "Interface ID" option is present but the "Router ID" field is zero, the "GDR Candidate Address" will be the IPv4 or IPv6 source address of the PIM Hello message.

This DRLBGDR Hello Option MUST only be advertised by the elected PIM DR.

## **[6.](#) Protocol Specification**

### **[6.1.](#) PIM DR Operation**

The DR election process is still the same as defined in [[RFC7761](#)]. A DR that has this specification enabled on an interface advertises the new DRLBGDR Hello Option, which contains mask values from user



configuration, followed by a sorted list of GDR Candidate Addresses, from the highest value to the lowest value. Moreover, same as non-DR routers, the DR also advertises DRLBC Hello Option to indicate its capability of supporting this specification and the type of its GDR election hash algorithm.

If a PIM DR receives a PIM Hello with the DRLBGDR Option, the PIM DR SHOULD ignore the TLV.

If a PIM DR receives a neighbor DRLBC Hello Option, which contains the same hash algorithm type as the DR, and the neighbor has the same DR priority as the DR, PIM DR SHOULD consider the neighbor as a GDR Candidate and insert the GDR Candidate's Address into the sorted list of the DRLBGDR Option. However, the DR MAY have policies limiting which GDR Candidates, or the number of GDR Candidates to include.

## **6.2. PIM GDR Candidate Operation**

When an IGMP/MLD report is received, without this specification, only the PIM DR will handle the join and potentially run into the issues described earlier. Using this specification, a hash algorithm is used by the GDR Candidates to determine which router is going to be responsible for building forwarding trees on behalf of the host.

If this specification is enabled on an interface, the router MUST include the DRLBC Hello Option in its PIM Hello on the interface. Note that the presence of the DRLBC Option in PIM Hello does not guarantee that this router would be considered as a GDR candidate. Once DR election is done, the DRLBGDR Hello Option would be received from the current PIM DR on the link which would contain a list of GDRs selected by the PIM DR.

A router only acts as a GDR candidate if it is included in the GDR list of the DRLBGDR Hello Option.

A GDR Candidate may receive a DRLBGDR Hello Option from the PIM DR with different Hash Masks from those the candidate was configured with. The GDR Candidate MUST use the Hash Masks advertised by the PIM DR to calculate the hash value.

A GDR Candidate MUST ignore the DRLBGDR Hello Option if it is received from a PIM router which is not the DR.

If the PIM DR does not support this specification, GDR election will not take place, and only the PIM DR joins the multicast tree.



### **6.2.1. Router Receives New DRLBGDR**

The first time a router receives a DRLBGDR option from the PIM DR, it MUST process the option and check if it is in the GDR list.

1. If a router is not listed as a GDR candidate in DRLBGDR, no action is needed.
2. If a router is listed as a GDR candidate in DRLBGDR, then it MUST process each of the groups, or source and group pairs if SSM, in the IGMP/MLD reports. The masks are announced in the PIM Hello by the DR in the DRLBGDR Hello Option. For each group in the reports that is in ASM mode, and each source and group pair if the group is in SSM mode, it (PIM Router) needs to run the hash algorithm (described in [section 4.3](#)) based on the announced Source, Group or RP masks to determine if it is the GDR for specified group, or source and group pair. If the hash result is to be the GDR for the multicast flow, it does build the multicast forwarding tree. If it is not the GDR for the multicast flow, no action is needed.

### **6.2.2. Router Receives Updated DRLBGDR**

If a router (GDR or non GDR) receives an unchanged DRLBGDR from the current PIM DR, no action is needed.

If a router (GDR or non GDR) receives a new or modified DRLBGDR from the current PIM DR, it requires processing as described below:

1. If it was included in the previous GDR list, and still is included in the new GDR list: It needs to process each of the groups, or source and group pairs if the group is in SSM mode, and run the hash algorithm to check if it is still the GDR for the given group, or source and group pair if SSM.

If it was the GDR for a group, or source and group pair if SSM, and the new hash result chose it as the GDR, then no processing is required.

If it was the GDR for a group, or source and group pair if SSM, earlier and now it is no longer the GDR, then it sets its assert metric for the multicast flow to be (PIM\_ASSERT\_INFINITY - 1), as explained in [Section 6.3](#).

If it was not the GDR for a group, or source and group pair if SSM, earlier, and the new hash does not make it GDR, then no processing is required.





If it was not the GDR for an earlier group, or source and group pair if SSM, and now becomes the GDR, it starts building multicast forwarding tree for this flow.

2. If it was included in the previous GDR list, but is not included in the new GDR list: It needs to process each of the groups, or source and group pairs if the group is in SSM mode.

If it was the GDR for a group, or source and group pair if SSM, it sets its assert metric for the multicast flow to be (PIM\_ASSERT\_INFINITY - 1), as explained in [Section 6.3](#).

If it was not the GDR, then no processing is required.

3. If it was not included in the previous GDR list, but is included in the new GDR list, the router MUST run the hash algorithm for each of the groups, source and group pairs if SSM.

If it is not the GDR for a group, or source and group pair if SSM, no processing is required.

If it is hashed as the GDR, it needs to build a multicast forwarding tree.

### **[6.3](#). PIM Assert Modification**

It is possible that the identity of the GDR might change in the middle of an active flow. Examples when this could happen include:

When a new PIM router comes up

When a GDR restarts

When the GDR changes, existing traffic might be disrupted. Duplicates or packet loss might be observed. To illustrate the case, consider the following scenario where there are two flows G1 and G2. R1 is the GDR for G1, and R2 is the GDR for G2. When R3 comes up online, it is possible that R3 becomes GDR for both G1 and G2, hence R3 starts to build the forwarding tree for G1 and G2. If R1 and R2 stop forwarding before R3 completes the process, packet loss might occur. On the other hand, if R1 and R2 continue forwarding while R3 is building the forwarding trees, duplicates might occur.

This is not a typical deployment scenario but might still happen. Here we describe a mechanism to minimize the impact. We essentially want to minimize packet loss. Therefore, we would allow a small amount of duplicates and depend on PIM Assert to minimize the duplication.



When the role of GDR changes as above, instead of immediately stopping forwarding, R1 and R2 continue forwarding to G1 and G2 respectively, while, at the same time, R3 build forwarding trees for G1 and G2. This will lead to PIM Asserts.

With the introduction of GDR, the following modification to the Assert packet MUST be done: if a router enables this specification on its downstream interface, but it is not a GDR (before network event it was GDR), it would adjust its Assert metric to (PIM\_ASSERT\_INFINITY - 1).

Using the above example, for G1, assume R1 and R3 agree on the new GDR, which is R3. R1 will set its Assert metric as (PIM\_ASSERT\_INFINITY - 1). That will make R3, which has normal metric in its Assert as the Assert winner.

For G2, assume it takes a slightly longer time for R2 to find out that R3 is the new GDR and still considers itself being the GDR while R3 already has assumed the role of GDR. Since both R2 and R3 think they are GDRs, they further compare their metric and IP addresses. If R3 has the better routing metric, or the same metric but a better tie-breaker, the result will be consistent during GDR selection. If unfortunately, R2 has the better metric or the same metric but a better tie-breaker, R2 will become the Assert winner and continues to forward traffic. This will continue until:

The next PIM Hello Option from DR selects R3 as the GDR. R3 will then build the forwarding tree and send an Assert.

The process continues until R2 agrees to the selection of R3 as the GDR, and sets its own Assert metric to (PIM\_ASSERT\_INFINITY - 1), which will make R3 the Assert winner. During the process, we will see intermittent duplication of traffic but packet loss will be minimized. In the unlikely case that R2 never relinquishes its role as GDR (while every other router thinks otherwise), the proposed mechanism also helps to keep the duplication to a minimum until manual intervention takes place to remedy the situation.

## **7. Compatibility**

In the case of a hybrid Ethernet shared LAN (where some PIM routers enable the specification defined in this document, and some do not)

- o If a router which does not support this specification becomes the DR on the LAN, then it is the only router acting as a DR, and there will be no load-balancing.



- o If a router which does not support this specification becomes a non-DR on link, then it acts as non-DR defined in [[RFC7761](#)], and it will not take part in any load-balancing.

## **8. Manageability Considerations**

Only the routers announcing the same Hash Algorithm as the DR would be considered as GDR candidates. Network administrators need to make sure that the desired set of routers announce the same algorithm. Migration between different algorithm types is not considered in this document.

## **9. IANA Considerations**

IANA has temporarily assigned type 34 for the PIM DR Load Balancing Capability (DRLBC) Hello Option, and type 35 for the PIM DR Load Balancing GDR (DRLBGDR) Hello Option. IANA is requested to make these assignments permanent when this document is published as an RFC. The string TBD should be replaced by the assigned values accordingly. This document requests IANA to create a DRLB hash type registry. This should be placed in the "Protocol Independent Multicast (PIM)" branch of the tree.

### **9.1. Initial registry**

The initial content of the registry should be as follows.

Type	Name	Reference
-----	-----	-----
0	Hash algorithm modulo	This document
1-255	Unassigned	

### **9.2. Assignment of new message types**

Assignment of new message types is done according to the "IETF Review" model, see [[RFC5226](#)].

## **10. Security Considerations**

Security of the new DR Load Balancing PIM Hello Options is only guaranteed by the security of PIM Hello messages, so the security considerations for PIM Hello messages as described in PIM-SM [[RFC7761](#)] apply here.



## **11. Acknowledgement**

The authors would like to thank Steve Simlo, Taki Millonis for helping with the original idea, Bill Atwood, Bharat Joshi for review comments, Toerless Eckert and Rishabh Parekh for helpful conversation on the document.

Special thanks to Anish Kachinthaya, Anvitha Kachinthaya and Jake Holland for reviewing the document and providing comments.

## **12. References**

### **12.1. 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, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC4291] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", [RFC 4291](#), DOI 10.17487/RFC4291, February 2006, <<https://www.rfc-editor.org/info/rfc4291>>.
- [RFC6395] Gulrajani, S. and S. Venaas, "An Interface Identifier (ID) Hello Option for PIM", [RFC 6395](#), DOI 10.17487/RFC6395, October 2011, <<https://www.rfc-editor.org/info/rfc6395>>.
- [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, <<https://www.rfc-editor.org/info/rfc7761>>.

### **12.2. Informative References**

- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", [RFC 5226](#), DOI 10.17487/RFC5226, May 2008, <<https://www.rfc-editor.org/info/rfc5226>>.

#### Authors' Addresses

Yiqun Cai  
Alibaba Group

Email: [yiqun.cai@alibaba-inc.com](mailto:yiqun.cai@alibaba-inc.com)





Heidi Ou  
Alibaba Group

Sri Vallepalli  
Cisco Systems, Inc.  
3625 Cisco Way  
San Jose CA 95134  
USA

Email: svallepa@cisco.com

Mankamana Mishra  
Cisco Systems, Inc.  
821 Alder Drive,  
Milpitas CA 95035  
USA

Email: mankamis@cisco.com

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

Email: stig@cisco.com

Andy Green  
British Telecom  
Adastral Park  
Ipswich IP5 2RE  
United Kingdom

Email: andy.da.green@bt.com

