

Network Working Group  
Internet-Draft  
Intended status: Informational  
Expires: April 5, 2013

A. Karan  
C. Filsfils  
D. Farinacci  
IJ. Wijnands, Ed.  
Cisco Systems, Inc.  
B. Decraene  
France Telecom  
U. Joorde  
Deutsche Telekom  
W. Henderickx  
Alcatel-Lucent  
October 02, 2012

**Multicast only Fast Re-Route  
draft-ietf-rtgwg-mofrr-00**

**Abstract**

As IPTV deployments grow in number and size, service providers are looking for solutions that minimize the service disruption due to faults in the IP network carrying the packets for these services. This draft describes a mechanism for minimizing packet loss in a network when node or link failures occur. Multicast only Fast Re-Route (MoFRR) works by making simple enhancements to multicast routing protocols such as PIM and mLDP.

**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 April 5, 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

<a href="#">1.</a>	Introduction . . . . .	<a href="#">3</a>
<a href="#">1.1.</a>	Conventions used in this document . . . . .	<a href="#">3</a>
<a href="#">1.2.</a>	Terminology . . . . .	<a href="#">3</a>
<a href="#">2.</a>	Basic Overview . . . . .	<a href="#">4</a>
<a href="#">3.</a>	Upstream Multicast Hop Selection . . . . .	<a href="#">4</a>
<a href="#">3.1.</a>	PIM . . . . .	<a href="#">4</a>
<a href="#">3.2.</a>	mLDP . . . . .	<a href="#">5</a>
<a href="#">4.</a>	Topologies for MoFRR . . . . .	<a href="#">5</a>
<a href="#">4.1.</a>	Dual-Plane Topology . . . . .	<a href="#">5</a>
<a href="#">5.</a>	Detecting Failures . . . . .	<a href="#">8</a>
<a href="#">6.</a>	ECMP-mode MoFRR . . . . .	<a href="#">9</a>
<a href="#">7.</a>	Non-ECMP-mode MoFRR . . . . .	<a href="#">9</a>
<a href="#">7.1.</a>	Variation . . . . .	<a href="#">11</a>
<a href="#">8.</a>	Keep It Simple Principle . . . . .	<a href="#">11</a>
<a href="#">9.</a>	Capacity Planning for MoFRR . . . . .	<a href="#">11</a>
<a href="#">10.</a>	Other Applications . . . . .	<a href="#">12</a>
<a href="#">11.</a>	Security Considerations . . . . .	<a href="#">12</a>
<a href="#">12.</a>	Acknowledgments . . . . .	<a href="#">13</a>
<a href="#">13.</a>	Contributing authors . . . . .	<a href="#">13</a>
<a href="#">14.</a>	References . . . . .	<a href="#">13</a>
<a href="#">14.1.</a>	Normative References . . . . .	<a href="#">13</a>
<a href="#">14.2.</a>	Informative References . . . . .	<a href="#">13</a>
	Authors' Addresses . . . . .	<a href="#">13</a>



## **1. Introduction**

Multiple techniques have been developed and deployed to improve service guarantees, both for multicast video traffic and Video on Demand traffic. Most existing solutions are geared towards finding an alternate path around one or more failed network elements (link, node, path failures).

This draft describes a mechanism for minimizing packet loss in a network when node or link failures occur. Multicast only Fast Re-Route (MoFRR) works by making simple changes to the way selected routers use multicast protocols such as PIM and mLDp. No changes to the protocols themselves are required. With MoFRR, in many cases, multicast routing protocols don't necessarily have to depend on or have to wait on unicast routing protocols to detect network failures.

On a merge point MoFRR logic determines a primary Upstream Multicast Hop (UMH) and a secondary UMH and joins the tree via both simultaneously. Data packets are received over the primary and secondary paths. Only the packets from the primary UMH are accepted and forwarded down the tree, the packets from the secondary UMH are discarded. The UMH determination is different for PIM and mLDp and explained later in this document. When a failure is detected on the path to the primary UMH, the repair occurs by changing the secondary UMH into the primary and the primary into the secondary. Since the repair is local, it is fast - greatly improving convergence times in the event of node or link failures on the path to the primary UMH.

### **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**

MoFRR : Multicast only Fast Re-Route.

ECMP : Equal Cost Multi-Path.

mLDp : Multi-point Label Distribution Protocol.

PIM : Protocol Independent Multicast.

UMH : Upstream Multicast Hop, a candidate next-hop that can be used to reach the root of the tree.



tree : Either a PIM (S,G)/(\*,G) tree or a mLDP P2MP or MP2MP LSP.

OIF : Outgoing InterFace, an interface used to forward multicast packets down the tree towards the receivers. Either a PIM (S,G)/(\*,G) tree or a mLDP P2MP or MP2MP LSP.

## **2. Basic Overview**

The basic idea of MoFRR is for a merge point router to join a multicast tree via two divergent upstream paths in order to get maximum redundancy. The two divergent paths SHOULD never merge upstream, otherwise the maximal redundancy is compromised. Sometimes the topology guarantees maximal redundancy, other times additional configuration or techniques are needed to enforce it. See later in this document.

A merge point router should only accept and forward on one of the upstream paths at the time in order to avoid duplicate packet forwarding. The selection of the primary and secondary UMH is done by the MoFRR logic and normally based on unicast routing to find loop free candidates.

Note, the impact of additional amount of data on the network is mitigated when tree membership is densely populated. When a part of the network has redundant data flowing, join latency for new joining members is reduced because its likely a tree merge point is not far away.

## **3. Upstream Multicast Hop Selection**

An Upstream Multicast Hop (UMH) is a candidate next-hop that can be used to reach the root of the tree. This is normally based on unicast routing to find loop free candidate(s). With MoFRR procedures we select a primary and a backup UMH. The procedures for determining the UMH are different for PIM and mLDP. See below;

### **3.1. PIM**

The UMH selection in PIM is also known as the Reverse Path Forwarding (RPF) procedure. Based on a unicast route lookup on either the Source address or Rendezvous Point (RP) [[RFC4601](#)], an upstream interface is selected for sending the PIM Joins/Prunes AND accepting the multicast packets. The interface the packets are received on is used to pass or fail the RPF check. If packets are received on an interface that was not selected by the RPF procedure, or not the primary, the packets are discarded.



### **[3.2.](#) mLDP**

The UMH selection in mLDP also depends on unicast routing, but the difference with PIM is that the acceptance of multicast packets is based on MPLS labels and independent on the interface the packet is received on. Using the procedures as defined in [\[RFC6388\]](#) an upstream Label Switched Router (LSR) is elected. The upstream LSR that was elected for a Label Switched Path (LSP) gets a unique local MPLS Label allocated. Multicast packets are only forwarded if the MPLS label matches the MPLS label that was allocated for that LSPs (primary) upstream LSR.

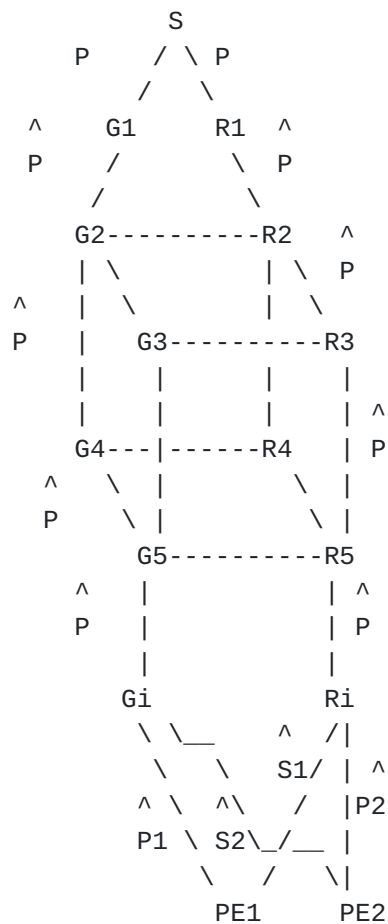
## **[4.](#) Topologies for MoFRR**

MoFRR works best in topologies illustrated in the figure below. MoFRR may be enabled on any router in the network. In the figures below, MoFRR is shown enabled on the Provider Edge (PE) routers to illustrate one way in which the technology may be deployed.

### **[4.1.](#) Dual-Plane Topology**







P = Primary path  
S = Secondary path

FIG1. Two-Plane Network Design

The topology has two planes, a primary plane and a secondary plane that are fully disjoint from each other all the way into the POPs. This two plane design is common in service provider networks as it eliminates single point of failures in their core network. The links marked PJ indicate the normal path of how the PIM joins flow from the POPs towards the source of the network. Multicast streams, especially for the densely watched channels, typically flow along both the planes in the network anyways.

The only change MoFRR adds to this is on the links marked S where the PE routers join a secondary path to their secondary ECMP UMH. As a result of this, each PE router receives two copies of the same stream, one from the primary plane and the other from the secondary plane. As a result of normal UMH behavior, the multicast stream received over the primary path is accepted and forwarded to the downstream receivers. The copy of the stream received from the



secondary UNH is discarded.

When a router detects a routing failure on the path to its primary UMH, it will switch to the secondary UMH and accept packets for that stream. If the failure is repaired the router may switch back. The primary and secondary UMHs have only local context and not end-to-end context.

As one can see, MoFRR achieves the faster convergence by pre-building the secondary multicast tree and receiving the traffic on that secondary path. The example discussed above is a simple case where there are two ECMP paths from each PE device towards the source, one along the primary plane and one along the secondary. In cases where the topology is asymmetric or is a ring, this ECMP nature does not hold, and additional rules have to be taken into account to choose when and where to join the secondary path.

MoFRR is appealing in such topologies for the following reasons:

1. Ease of deployment and simplicity: the functionality is only required on the PE devices although it may be configured on all routers in the topology. Furthermore, each PE device can be enabled separately. PEs not enabled for MoFRR do not see any change or degradation. Inter-operability testing is not required as there are no PIM or mLDp protocol change.
2. End-to-end failure detection and recovery: any failure along the path from the source to the PE can be detected and repaired with the secondary disjoint stream.
3. Capacity Efficiency: as illustrated in the previous example, the Multicast trees corresponding to IPTV channels cover the backbone and distribution topology in a very dense manner. As a consequence, the secondary path graft into the normal Multicast trees (ie. trees signaled by PIM or mLDp without MoFRR extension) at the aggregation level and hence do not demand any extra capacity either on the distribution links or in the backbone. They simply use the capacity that is normally used, without any duplication. This is different from conventional FRR mechanisms which often duplicate the capacity requirements (the backup path crosses links/nodes which already carry the primary/normal tree and hence twice as much capacity is required).
4. Loop free: the secondary path join is sent on an ECMP disjoint path. By definition, the neighbor receiving this request is closer to the source and hence will not cause a loop.

The topology we just analyzed is very frequent and can be modeled as



per Fig2. The PE has two ECMP disjoint paths to the source. Each ECMP path uses a disjoint plane of the network.

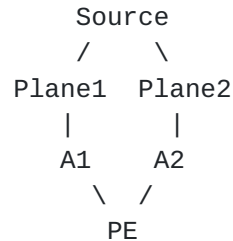


FIG2. PE is dual-homed to Dual-Plane Backbone

Another frequent topology is described in Fig 3. PEs are grouped by pairs. In each pair, each PE is connected to a different plane. Each PE has one single shortest-path to a source (via its connected plane). There is no ECMP like in Fig 2. However, there is clearly a way to provide MoFRR benefits as each PE can offer a disjoint secondary path to the other plane PE (via the disjoint path).

MoFRR secondary neighbor selection process needs to be extended in this case as one cannot simply rely on using an ECMP path as secondary neighbor. This extension is referred to as non-ecmp extension and is described later in the document.

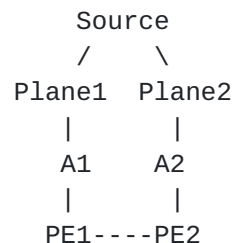


FIG3. PEs are connected in pairs to Dual-Plane Backbone

## 5. Detecting Failures

Once the two paths are established, the next step is detecting a failure on the primary path to know when to switch to the backup path.

A first option consists of comparing the packets received on the



primary and secondary streams but only forwarding one of them -- the first one received, no matter which interface it is received on. Zero packet loss is possible for RTP-based streams.

A second option assumes a minimum known packet rate for a given data stream. If a packet is not received on the primary RPF within this time frame, the router assumes primary path failure and switches to the secondary RPF interface. 50msec switchover is possible.

A third option leverages the significant improvements of the IGP convergence speed. When the primary path to the source is withdrawn by the IGP, the MoFRR-enabled router switches over to the backup path, the UMH is changed to the secondary UMH. Since the secondary path is already in place, and assuming it is disjoint from the primary path, convergence times would not include the time required to build a new tree and hence are smaller. Realistic availability requirements (sub-second to sub-200msec) should be possible.

A fourth option consists in leveraging connected link failure. This option makes sense when MoFRR is deployed across the network (not only at PE).

## **6. ECMP-mode MoFRR**

If the IGP installs two ECMP paths to the source and if the Multicast tree is enabled for ECMP-Mode MoFRR, the router installs them as primary and secondary UMH. Only packets received from the primary UMH path are processed. Packets received from the secondary UMH are dropped.

The selected primary UMH should be the same as if MoFRR extension was not enabled.

If more than two ECMP paths exist, two are selected as primary and secondary UMH. Information from the IGP link-state topology could be leveraged to optimize this selection.

Note, MoFRR does not restrict the number of UMH paths that are joined. Implementations may use as many paths as are configured.

## **7. Non-ECMP-mode MoFRR**





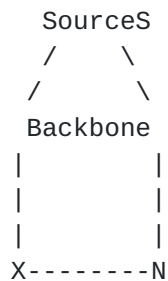


Fig5. Non-ECMP-Mode MoFRR

X is configured for MoFRR for a Multicast tree

$R(X)$  is  $X$ 's UMH to S

N is a neighbor of X

$R(N)$  is N's UMH to S

$x_s$  represents the IGP metric from X to S

$n_s$  represents the IGP metric from N to S

$x_n$  represents the IGP metric from X to N

A router X configured for non-ECMP-mode MoFRR for a Multicast tree joins a primary path to its primary UMH  $R(X)$  and a secondary path to UMH N if the following three conditions are met.

C1:  $x_s < x_n + n_s$

C2:  $n_s < n_x + x_s$

C3: X cannot join the secondary path N if N is the only member of the OIF list

The first condition ensures that N is not on the primary branch from X to S.

The second condition ensures that X is not on the primary branch from N to S.

These two conditions ensure that at least locally the two paths are disjoint.

The third condition is required to break control-plane loops which could occur in some scenarios.

For example in FIG3, if PE1 and PE2 have received an igmp request for a Multicast tree, they will both join the primary path on their plane and a secondary path to the neighbor PE. If their receivers would leave at the same time, it could be possible for the Multicast tree on PE1 and PE2 to never get deleted as each PE refresh each other via the secondary path joins (remember that a secondary path join is not distinguishable from a primary join. MoFRR does not require any PIM or mLDp protocol modification).



A control-plane loop occurs when two nodes keep a state forever due to joining the secondary path to each other. This forever condition is not acceptable as no real receiver is connected to the nodes (directly via IGMP or indirectly via PIM). Rule 3 prevents this case as it prevents the mutual refresh of secondary joins and it applies it in the specific case where there is no real receiver connected.

### **7.1. Variation**

Rule R3 can be removed if Rule 2 is restricted as follows:

R2p:  $ns < xs$

This ensures that X will only join the secondary path to a neighbor N who is strictly closer to the source than X is. By reciprocity, N will thus never join the secondary path for the same Multicast tree via X. The strictly smaller than is key here.

Note that this non-ECMP-mode MoFRR variation does not support the square topology and hence is less preferred.

## **8. Keep It Simple Principle**

Many Service Providers devise their topology such that PEs have disjoint paths to the multicast sources. MoFRR leverages the existence of these disjoint paths without any PIM or mLDp protocol modification. Interoperability testing is thus not required. In such topologies, MoFRR only needs to be deployed on the PE devices. Each PE device can be enabled one by one. PEs not enabled for MoFRR do not see any change or degradation.

Multicast streams with Tight SLA requirements are often characterized by a continuous high packet rate (SD video has a continuous interpacket gap of  $\sim 3\text{msec}$ ). MoFRR simply leverages the stream characteristic to detect any failures along the primary branch and switch-over on the secondary branch in a few 10s of msec.

## **9. Capacity Planning for MoFRR**

As for LFA FRR ([draft-ietf-rtgwg-lfa-applicability-00](#)), MoFRR applicability is topology dependent.

In this document, we have described two very frequent designs (Fig 2 and Fig 3) which provide maximum MoFRR benefits.

Designers with topologies different than Fig2 and 3 can still benefit



from MoFRR benefits thanks to the use of capacity planning tools.

Such tools are able to simulate the ability of each PE to build two disjoint branches of the same tree. This for hundreds of PEs and hundreds of sources.

This allows to assess the MoFRR protection coverage of a given network, for a set of sources.

If the protection coverage is deemed insufficient, the designer can use such tool to optimize the topology (add links, change igp metrics).

## **10. Other Applications**

While all the examples in this document show the MoFRR applicability on PE devices, it is clear that MoFRR could be enabled on aggregation or core routers.

MoFRR can be popular in Data Center network configurations. With the advent of lower cost ethernet and increasing port density in routers, there is more meshed connectivity than ever before. When using a 3-level access, distribution, and core layers in a Data Center, there is a lot of inexpensive bandwidth connecting the layers. This will lend itself to more opportunities for ECMP paths at multiple layers. This allows for multiple layers of redundancy protecting link and node failure at each layer with minimal redundancy cost.

Redundancy costs are reduced because only one packet is forwarded at every link along the primary and secondary data paths so there is no duplication of data on any link thereby providing make-before-break protection at a very small cost.

Alternate methods to detect failures such as MPLS-OAM or BFD may be considered.

The MoFRR principle may be applied to MVPNs.

## **11. Security Considerations**

There are no security considerations for this design other than what is already in the main PIM specification [[RFC4601](#)] and mLDp specification [[RFC6388](#)] .



## **12. Acknowledgments**

The authors would like to thank John Zwiebel, Greg Shepherd and Dave Oran for their review of the draft.

## **13. Contributing authors**

Below is a list of other contributing authors in alphabetical order:

Nicolai Leymann  
Deutsche Telekom  
Winterfeldtstrasse 21  
Berlin 10781  
DE  
Email: N.Leymann@telekom.de

## **14. References**

### **14.1. Normative References**

- [RFC5036] Andersson, L., Minei, I., and B. Thomas, "LDP Specification", [RFC 5036](#), October 2007.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

### **14.2. Informative References**

- [RFC4601] Fenner, B., Handley, M., Holbrook, H., and I. Kouvelas, "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)", [RFC 4601](#), August 2006.
- [RFC6388] Wijnands, IJ., Minei, I., Kompella, K., and B. Thomas, "Label Distribution Protocol Extensions for Point-to-Multipoint and Multipoint-to-Multipoint Label Switched Paths", [RFC 6388](#), November 2011.





Authors' Addresses

Apoorva Karan  
Cisco Systems, Inc.  
3750 Cisco Way  
San Jose CA, 95134  
USA

Email: [apoorva@cisco.com](mailto:apoorva@cisco.com)

Clarence Filsfils  
Cisco Systems, Inc.  
De kleetlaan 6a  
Diegem BRABANT 1831  
Belgium

Email: [cfilsfil@cisco.com](mailto:cfilsfil@cisco.com)

Dino Farinacci  
Cisco Systems, Inc.  
425 East Tasman Drive  
San Jose CA, 95134  
USA

Email: [dino@cisco.com](mailto:dino@cisco.com)

IJsbrand Wijnands (editor)  
Cisco Systems, Inc.  
De Kleetlaan 6a  
Diegem 1831  
BE

Email: [ice@cisco.com](mailto:ice@cisco.com)

Bruno Decraene  
France Telecom  
38-40 rue du General Leclerc  
Issy Moulineaux cedex 9, 92794  
FR

Email: [bruno.decraene@orange-ftgroup.com](mailto:bruno.decraene@orange-ftgroup.com)



Uwe Joorde  
Deutsche Telekom  
Hammer Str. 216-226  
Muenster D-48153  
DE

Email: Uwe.Joorde@telekom.de

Wim Henderickx  
Alcatel-Lucent  
Copernicuslaan 50  
Antwerp 2018  
Belgium

Email: wim.henderickx@alcatel-lucent.com

