Routing Area Working Group Internet-Draft Intended status: Standards Track Expires: August 18, 2014 S. Litkowski B. Decraene Orange P. Francois IMDEA Networks C. Filsfils Cisco Systems February 14, 2014

## Microloop prevention by introducing a local convergence delay draft-litkowski-rtgwg-uloop-delay-02

## Abstract

This document describes a mechanism for link-state routing protocols to prevent local transient forwarding loops in case of link failure. This mechanism Proposes a two-steps convergence by introducing a delay between the convergence of the node adjacent to the topology change and the network wide convergence.

As this mechanism delays the IGP convergence it may only be used for planned maintenance or when fast reroute protects the traffic between the link failure and the IGP convergence.

Simulations using real network topologies have been performed and show that local loops are a significant portion (>50%) of the total forwarding loops.

## Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [<u>RFC2119</u>].

#### Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <u>http://datatracker.ietf.org/drafts/current/</u>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on August 18, 2014.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>http://trustee.ietf.org/license-info</u>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Internet-Draft

# Table of Contents

$\underline{1}$ . Introduction	. <u>4</u>
$\underline{2}$ . Overview of the solution	. <u>4</u>
<u>3</u> . Specification	. <u>4</u>
<u>3.1</u> . Definitions	. <u>4</u>
<u>3.2</u> . Current IGP reactions	. <u>5</u>
<u>3.3</u> . Local events	. <u>5</u>
<u>3.4</u> . Local delay	. <u>6</u>
<u>3.4.1</u> . Link down event	. <u>6</u>
<u>3.4.2</u> . Link up event	. <u>6</u>
<u>4</u> . Applicability	· <u>7</u>
<u>4.1</u> . Applicable case : local loops	· <u>7</u>
<u>4.2</u> . Non applicable case : remote loops	· <u>7</u>
<u>5</u> . Simulations	. <u>8</u>
<u>6</u> . Deployment considerations	. <u>9</u>
<u>7</u> . Security Considerations	. <u>9</u>
8. Acknowledgements	. <u>10</u>
9. IANA Considerations	. <u>10</u>
<u>10</u> . References	. <u>10</u>
<u>10.1</u> . Normative References	. <u>10</u>
<u>10.2</u> . Informative References	. <u>10</u>
Authors' Addresses	. <u>11</u>

Internet-Draft

uloop-delay

# **1**. Introduction

In figure 1, upon link AD down event, for the destination A, if D updates its forwarding entry before C, a transient forwarding loop occurs between C and D. We have a similar loop for link up event, if C updates its forwarding entry A before D.

```
A ----- B

| |

| D-----C All the links have a metric of 1 except BC=5
```

Figure 1

# 2. Overview of the solution

This document defines a two-step convergence initiated by the router detecting the failure and advertising the topological changes in the IGP. This introduces a delay between the convergence of the local router and the network wide convergence. This delay is positive in case of "down" events and negative in case of "up" events.

This ordered convergence, is similar to the ordered FIB proposed defined in [I-D.ietf-rtgwg-ordered-fib], but limited to only one hop distance. As a consequence, it is simpler and becomes a local only feature not requiring interoperability; at the cost of only covering the transient forwarding loops involving this local router. The proposed mechanism also reuses some concept described in [I-D.ietf-rtgwg-microloop-analysis] with some limitation and improvements.

# $\underline{\mathbf{3}}$ . Specification

# 3.1. Definitions

This document will refer to the following existing IGP timers:

- o LSP\_GEN\_TIMER: to batch multiple local events in one single local LSP update. It is often associated with damping mechanism to slowdown reactions by incrementing the timer when multiple consecutive events are detected.
- o SPF\_TIMER: to batch multiple events in one single computation. It is often associated with damping mechanism to slowdown reactions by incrementing the timer when the IGP is instable.

o IGP\_LDP\_SYNC\_TIMER: defined in [<u>RFC5443</u>] to give LDP some time to establish the session and learn the MPLS labels before the link is used.

This document introduces the following two new timers :

- o ULOOP\_DELAY\_DOWN\_TIMER: slowdown the network wide IGP convergence in case of link down events.
- o ULOOP\_DELAY\_UP\_TIMER: slowdown the local node convergence in case of link up events.

#### 3.2. Current IGP reactions

Upon a change of status on an adjacency/link, the existing behavior of the router advertising the event is the following:

- 1. UP/Down event is notified to IGP.
- IGP processes the notification and postpones the reaction in LSP\_GEN\_TIMER msec.
- Upon LSP\_GEN\_TIMER expiration, IGP updates its LSP/LSA and floods it.
- 4. SPF is scheduled in SPF\_TIMER msec.
- 5. Upon SPF\_TIMER expiration, SPF is computed and RIB/FIB are updated.

## 3.3. Local events

In the next sections, we will use the concept of local events versus remote events. The notion of event we are using in this document is linked to IGP link state advertisements and not network events, as a single network event would create multiple IGP link state advertisement within the network.

A local event is a set of IGP link state advertisements describing only a change of a local component of the computing router (e.g. a link). As opposite to a remote event being a set of IGP link state advertisements describing any other type of changes.

Example :

+---- E ----++ | | | | A ----- B ------- C ------ D

Considering computing router is B, when B-C fails. B updates its local LSP describing the link B->C being down, C does exactly the same and starts flooding. During SPF\_TIMER, B and C LSPs would be taken into account. B and C LSPs are describing exactly the same event (B-C link down). For B point of view, both LSPs must be considered as a local event as they are describing the change of a local component of B (link B-C). If C node is failing, routers B,E and D are updating and flooding their LSPs. LSPs from E and D are considered as remote events for B as they are describing a change in a component that does not belong to B. Hence the local delay mechanism will be aborted. Hence this mechanism is not applicable to node failure.

## <u>3.4</u>. Local delay

#### 3.4.1. Link down event

Upon an adjacency/link down event, this document introduces a change in step 5 in order to delay the local convergence compared to the network wide convergence: the node SHOULD delay the forwarding entry updates by ULOOP\_DELAY\_DOWN\_TIMER. Such delay SHOULD only be introduced if all the LSDB modifications processed are only reporting down local events . Note that determining that all topological change are only local down events requires analyzing all modified LSP/LSA as a local link or node failure will typically be notified by multiple nodes. If a subsequent LSP/LSA is received/updated and a new SPF computation is triggered before the expiration of ULOOP\_DELAY\_DOWN\_TIMER, then the same evaluation SHOULD be performed.

As a result of this addition, routers local to the failure will converge slower than remote routers. Hence it SHOULD only be done for non urgent convergence, such as for administrative de-activation (maintenance) or when the traffic is Fast ReRouted.

## <u>3.4.2</u>. Link up event

Upon an adjacency/link up event, this document introduces the following change in step 3 where the node SHOULD:

- o Firstly build a LSP/LSA with the new adjacency but setting the metric to MAX\_METRIC . It SHOULD flood it but not compute the SPF at this time. This step is required to ensure the two way connectivity check on all nodes when computing SPF.
- o Then build the LSP/LSA with the target metric but SHOULD delay the flooding of this LSP/LSA by SPF\_TIMER + ULOOP\_DELAY\_UP\_TIMER. MAX\_METRIC is equal to MaxLinkMetric (0xFFFF) for OSPF and 2^24-2 (0xFFFFE) for IS-IS.

uloop-delay

o Then continue with next steps (SPF computation) without waiting for the expiration of the above timer. In other word, only the flooding of the LSA/LSP is delayed, not the local SPF computation.

As as result of this addition, routers local to the failure will converge faster than remote routers.

If this mechanism is used in cooperation with "LDP IGP Synchronization" as defined in [RFC5443] then the mechanism defined in RFC 5443 is applied first, followed by the mechanism defined in this document. More precisely, the procedure defined in this document is applied once the LDP session is considered "fully operational" as per [RFC5443].

#### **<u>4</u>**. Applicability

As previously stated, the mechanism only avoids the forwarding loops on the links between the node local to the failure and its neighbor. Forwarding loops may still occur on other links.

#### 4.1. Applicable case : local loops

А В	E								
/									
/									
GC	F	A11	the	links	have	а	metric	of	1

Figure 2

Let us consider the traffic from G to F. The primary path is G->D->C->E->F. When link CE fails, if C updates its forwarding entry for F before D, a transient loop occurs. This is sub-optimal as C has FRR enabled and it breaks the FRR forwarding while all upstream routers are still forwarding the traffic to itself.

By implementing the mechanism defined in this document on C, when the CE link fails, C delays the update of his forwarding entry to F, in order to let some time for D to converge. FRR keeps protecting the traffic during this period. When the timer expires on C, forwarding entry to F is updated. There is no transient forwarding loop on the link CD.

# **4.2**. Non applicable case : remote loops

A ----- B ---- E --- H | | | G---D-----C -----F --- J ---- K

All the links have a metric of 1 except BE=15

Figure 3

Let us consider the traffic from G to K. The primary path is G->D->C->F->J->K. When the CF link fails, if C updates its forwarding entry to K before D, a transient loop occurs between C and D.

By implementing the mechanism defined in this document on C, when the link CF fails, C delays the update of his forwarding entry to K, letting time for D to converge. When the timer expires on C, forwarding entry to F is updated. There is no transient forwarding loop between C and D. However, a transient forwarding loop may still occur between D and A. In this scenario, this mechanism is not enough to address all the possible forwarding loops. However, it does not create additional traffic loss. Besides, in some cases -such as when the nodes update their FIB in the following order C, A, D, for example because the router A is quicker than D to converge- the mechanism may still avoid the forwarding loop that was occuring.

#### 5. Simulations

Simulations have been run on multiple service provider topologies. So far, only link down event have been tested.

+   Topology	·+- 	Gain
+	+ -	+
T1	Ι	71%
T2	Ι	81%
ТЗ	Ι	62%
T4	Ι	50%
T5	Ι	70%
Тб	Ι	70%
T7	Ι	59%
Т8	T	77%
+	. + -	+

Table 1: Number of Repair/Dst that may loop

We evaluated the efficiency of the mechanism on eight different service provider topologies (different network size, design). The

benefit is displayed in the table above. The benefit is evaluated as follows:

- o We consider a tuple (link A-B, destination D, PLR S, backup nexthop N) as a loop if upon link A-B failure, the flow from a router S upstream from A (A could be considered as PLR also) to D may loop due to convergence time difference between S and one of his neighbor N.
- o We evaluate the number of potential loop tuples in normal conditions.
- o We evaluate the number of potential loop tuples using the same topological input but taking into account that S converges after N.
- o Gain is how much loops (remote and local) we succeed to suppress.

On topology 1, 71% of the transient forwarding loops created by the failure of any link are prevented by implementing the local delay. The analysis shows that all local loops are obviously solved and only remote loops are remaining.

#### <u>6</u>. Deployment considerations

Transient forwarding loops have the following drawbacks :

- o Limit FRR efficiency : even if FRR is activated in 50msec, as soon as PLR has converged, traffic may be affected by a transient loop.
- o It may impact traffic not directly concerned by the failure (due to link congestion).

This local delay proposal is a transient forwarding loop avoidance mechanism (like OFIB). Even if it only address local transient loops, , the efficiency versus complexity comparison of the mechanism makes it a good solution. It is also incrementally deployable with incremental benefits, which makes it an attractive option for both vendors to implement and Service Providers to deploy. Delaying convergence time is not an issue if we consider that the traffic is protected during the convergence.

## 7. Security Considerations

This document does not introduce change in term of IGP security. The operation is internal to the router. The local delay does not

uloop-delay

increase the attack vector as an attacker could only trigger this mechanism if he already has be ability to disable or enable an IGP link. The local delay does not increase the negative consequences as if an attacker has the ability to disable or enable an IGP link, it can already harm the network by creating instability and harm the traffic by creating forwarding packet loss and forwarding loss for the traffic crossing that link.

## 8. Acknowledgements

We wish to thanks the authors of [<u>I-D.ietf-rtgwg-ordered-fib</u>] for introducing the concept of ordered convergence: Mike Shand, Stewart Bryant, Stefano Previdi, and Olivier Bonaventure.

## 9. IANA Considerations

This document has no actions for IANA.

## **10**. References

#### <u>**10.1</u>**. Normative References</u>

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, March 1997.
- [RFC5443] Jork, M., Atlas, A., and L. Fang, "LDP IGP Synchronization", <u>RFC 5443</u>, March 2009.
- [RFC5715] Shand, M. and S. Bryant, "A Framework for Loop-Free Convergence", <u>RFC 5715</u>, January 2010.

# <u>10.2</u>. Informative References

[I-D.ietf-rtgwg-microloop-analysis]

Zinin, A., "Analysis and Minimization of Microloops in Link-state Routing Protocols", <u>draft-ietf-rtgwg-microloop-analysis-01</u> (work in progress), October 2005.

[I-D.ietf-rtgwg-ordered-fib]

Shand, M., Bryant, S., Previdi, S., Filsfils, C., Francois, P., and O. Bonaventure, "Framework for Loop-free convergence using oFIB", <u>draft-ietf-rtgwg-ordered-fib-12</u> (work in progress), May 2013.

- [I-D.ietf-rtgwg-remote-lfa]
  Bryant, S., Filsfils, C., Previdi, S., Shand, M., and S.
  Ning, "Remote LFA FRR", draft-ietf-rtgwg-remote-lfa-04
  (work in progress), November 2013.
- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", <u>RFC 3630</u>, September 2003.
- [RFC6571] Filsfils, C., Francois, P., Shand, M., Decraene, B., Uttaro, J., Leymann, N., and M. Horneffer, "Loop-Free Alternate (LFA) Applicability in Service Provider (SP) Networks", <u>RFC 6571</u>, June 2012.

Authors' Addresses

Stephane Litkowski Orange

Email: stephane.litkowski@orange.com

Bruno Decraene Orange

Email: bruno.decraene@orange.com

Pierre Francois IMDEA Networks

Email: pierre.francois@imdea.org

Clarence Fils Fils Cisco Systems

Email: cf@cisco.com