

RTGWG
Internet-Draft
Intended status: Standards Track
Expires: April 5, 2013

P. Thubert, Ed.
P. Bellagamba
Cisco Systems
October 2, 2012

Available Routing Constructs
draft-thubert-rtgwg-arc-00

Abstract

This draft introduces the concept of ARC, a two-edged routing construct that forms its own fault isolation and recovery domain. The new paradigm can be leveraged to improve the network utilization and resiliency for unicast and multicast traffic in multiple environments.

Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

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

1. Introduction	3
2. Terminology	4
3. ARC Set representations	7
4. Applicability	11
4.1. Load Balancing	11
4.1.1. Routing Hierarchies	11
5. Lowest ARC First	11
5.1. Init	12
5.2. Growing Trees	12
5.3. Being Safe	12
5.4. Bending An ARC	13
5.5. Orienting Links	14
5.6. Looping or recursing	14
6. Forwarding Along An ARC Set	15
6.1. Control Plane Recovery	16
6.2. Data Plane Recovery	16
7. Manageability	17
8. IANA Considerations	17
9. Security Considerations	17
10. Acknowledgements	17
11. References	18
11.1. Normative References	18
11.2. Informative References	18
Authors' Addresses	19

1. Introduction

Traditional routing and forwarding uses the concept of path as the basic routing paradigm to get a packet from a source to a destination by following an ordered sequence of arrows between intermediate nodes. In this serial design, a path is broken as soon as a single arrow is, and getting around a breakage can require path recomputation, network reconvergence, and incur delays to till service is restored.

Multiple paths can be bound together for instance to form a Directed Acyclic Graph (DAG) to a destination, but that technique can be difficult to balance and cannot provide a full path redundancy even in the case of a biconnected graph. For instance, if the node that is closest to the DAG destination has only one link to that destination, then it does not have a alternate path to get to that destination.

It is also possible to compute an alternate routing topology for fast rerouting to a given destination, in which case some signalling, tagging or labelling can be put in place to indicate whether a packet follows the normal path or was rerouted over an alternate topology. Once a packet is rerouted, it is bound to the alternate topology so only one breakage can be handled with looplessness guarantees in most practical situations.

This draft introduces the concept of an Available Routing Construct (ARC) as a routing construct made of a sequence of nodes and links with 2 outgoing edges, so that, upon a single breakage, each lively node in along ARC can still reach one of the outgoing edges. As a result, an ARC is this resilient to one breakage as opposed to an arrow that has only one outgoing edge, and an ARC topology is resilient to one breakage per ARC.

The routing graph to reach a certain destination is expressed as a cascade of ARCs, each ARC providing its own independent domain of fault isolation and recovery. Unicast traffic may enter an ARC via any node but it may only leave the ARC through one of its two edges. One node along the ARC is designated as the cursor. In normal unicast operations, the traffic inside an ARC flows away from the cursor towards an edge. Upon a failure, packets may bounce on the breakage point and flow the other way along the ARC to take the other exit.

As a result an ARC is resilient to any single failure, and the recovery can be driven either from the data plane or the control plane. A second failure occurring within a same ARC will isolate an ARC segment. This can be further corrected from the control plane by

reversing all the incoming Edges in a process that might recurse till an exit is found. When ARC reversal is applied, an ARC topology is resilient to some cases of Shared Risk Link Group (SRLG) failures.

This draft presents the concept and provides an intuition of how ARCs can simplify the operation and improve the network utilization and resiliency for all sorts of traffic in multiple environments, but defers to further documents to elaborate on the algorithms and optimizations in the different application domains.

For instance, ARCs can also be used in datacenters for the purpose of fast-reroute, or within a service provider network to simplify load balancing operations or leverage optimally the ring topologies [RFC5921]. An ARC topology can be flooded over itself and serve as a backbone for reliable multicasting operations.

2. Terminology

The draft uses the following terminology:

ARC: Available Routing Construct. An ARC is a loopless ordered set of nodes and links whereby traffic may enter via any node in the ARC but may only leave the ARC through either one of the ARC edges.

Comb: An ARC generalization: a Comb is a n -edged loopless set of nodes and links with $n \geq 2$; traffic may enter via any node in the Comb but may only exit the Comb through one of its n edges. A Comb comes with a walk operation that enables to attempt to exit via every edge and to discover when all have been tried.

Cursor: A virtual point along an ARC that can be located on a node or on a link between 2 nodes. In normal operations, the traffic along the ARC flows away from its Cursor. If the cursor is a node, then traffic can be distributed on both sides. The Cursor may be moved to change the way traffic is load balanced along an ARC. It may also be placed at the location of a failure to direct traffic away from that point.

ARC Node: A Node that belongs to an ARC.

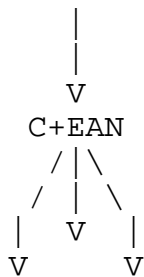
Edge ARC Node: An ARC Node at an edge of its ARC. An Edge ARC Node is a node via which traffic can exit the ARC.

Edge Link: A directed link outgoing from an Edge ARC Node. Traffic can only exit from an ARC via an Edge Link. An Edge Link does not accept traffic into an ARC.

Intermediate ARC Node: A node that is not at an edge of an ARC. A Intermediate ARC Node node that can receive traffic and forward traffic between its adjacent nodes.

Intermediate Link: A link between two Intermediate ARC Nodes. An Intermediate Link is reversible, meaning that traffic is allowed in both directions though an individual packet is constrained in the way its direction is reversed. For stable links such as wired links, the typical constraint is that the direction of a packet may be reversed at most once along a given ARC.

Collapsed ARC: An ARC that is formed of a single node. This node is altogether the cursor and both Edge Nodes. This implies that the node has at least 2 outgoing links to 2 different Safe Nodes.



E: Edge ARC Node	-	collapsed in a single node
C: Cursor	-	

Figure 1: Collapsed ARC

Infrastructure ARC: An ARC that is formed of more than one node, which also means that the Edge Nodes are different nodes.

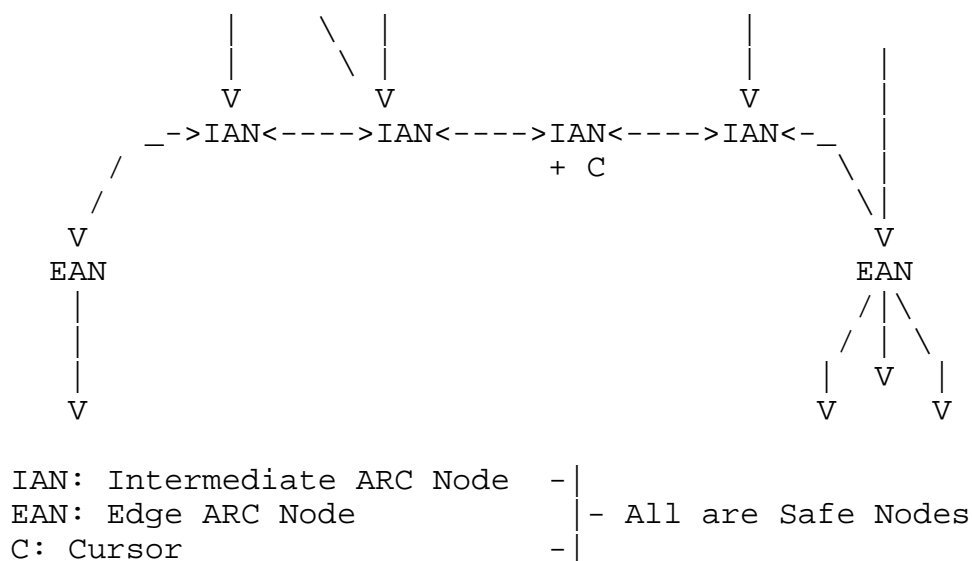


Figure 2: Infrastructure ARC

DAG: Directed Acyclic Graph.

ARC Set (or Cascade): A DAG with ARCs as vertices. In the DAG, an edge between ARC A and ARC B corresponds to a link from an Edge ARC Node in ARC A and an arbitrary ARC Node in ARC B. Note that by definition, an ARC has at least 2 outgoing Edge Links, one per Edge Node, and maybe more if an Edge Node has multiple outgoing Edge Links. All vertices in the DAG have 2 forwarding solutions, even the ARC closest to the destination.

Omega: the abstract destination (== root) of an ARC Set.

ARC Height: An arbitrary distance from Omega that is associated to an ARC. The Height of an ARC must be more than the Height of any of the ARCs it terminates into. The order of ARC formation by a given algorithm can be used as a Height whereby an ARC is always strictly higher or lower than another.

Buttressing ARC: A split ARC that is merged into another ARC at one edge. An ARC and one or more Buttressing ARCs form a Comb construct that is resilient to additional breakages. A Buttressing ARC may be applied to an ARC or a Comb iff traffic outgoing the Buttressing ARC Edge always reaches in an ARC that is lower than this ARC, or Omega.

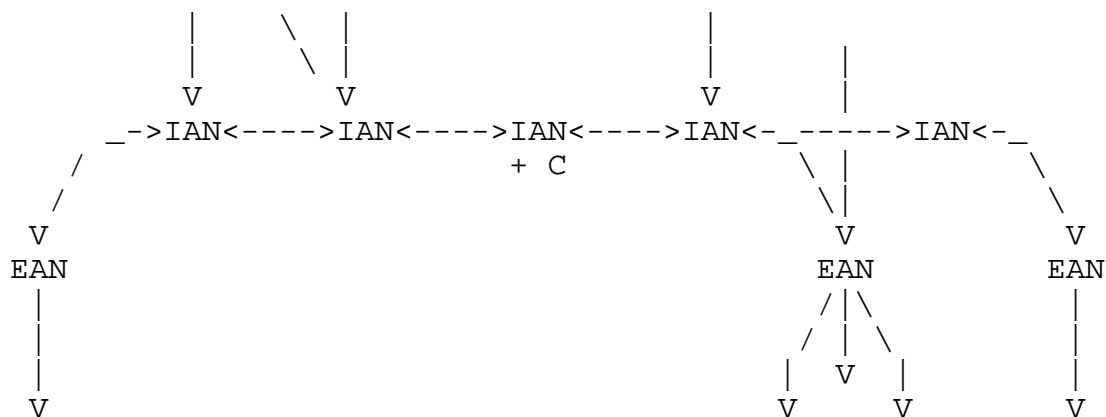


Figure 3: Comb with Buttrressing ARC

Safe Node: A node is Safe if there is no single point of failure - apart from the node itself - on its way to Omega. From this definition, a node is Safe if it has at least two non-congruent paths to two different other Safe Nodes. It results that a Safe node that is not Omega has at least two completely disjunct paths to Omega. When an ARC has been successfully constructed, all its nodes become safe with respect to the Omega for which the ARC was constructed. By extension for a collapsed path Omega is deemed to be Safe, that is any node that pertains in Omega is a Safe Node.

?-S: A node N is deemed dependent on a node S or S-dependent (denoted as ?-S) if S is the last single point of failure along N's shortest path to Omega.

3. ARC Set representations

An ARC Set can be represented in a number of fashions:

Graph View:

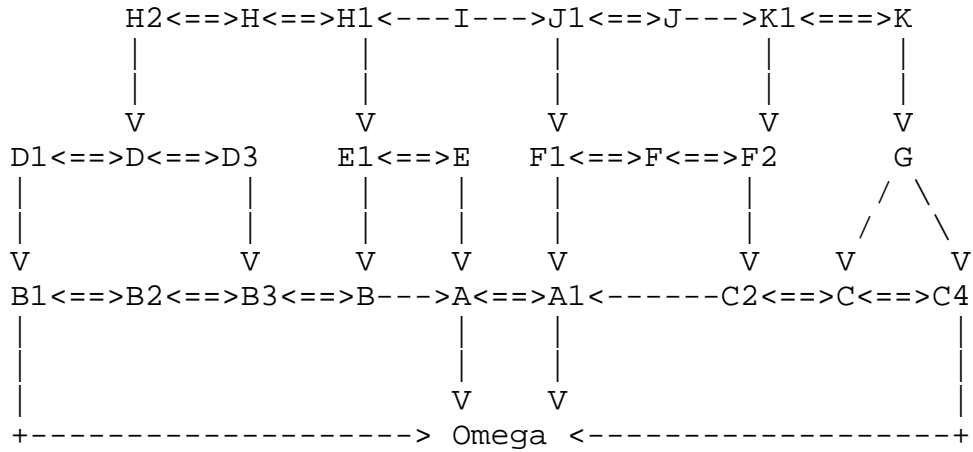


Figure 4: Routing Graph View

This representation is similar to a classical routing graph with the peculiarity that some Links are marked reversible. An ARC is represented as a sequence of reversible links. The node that holds the cursor is also indicated somehow.

ARC View:

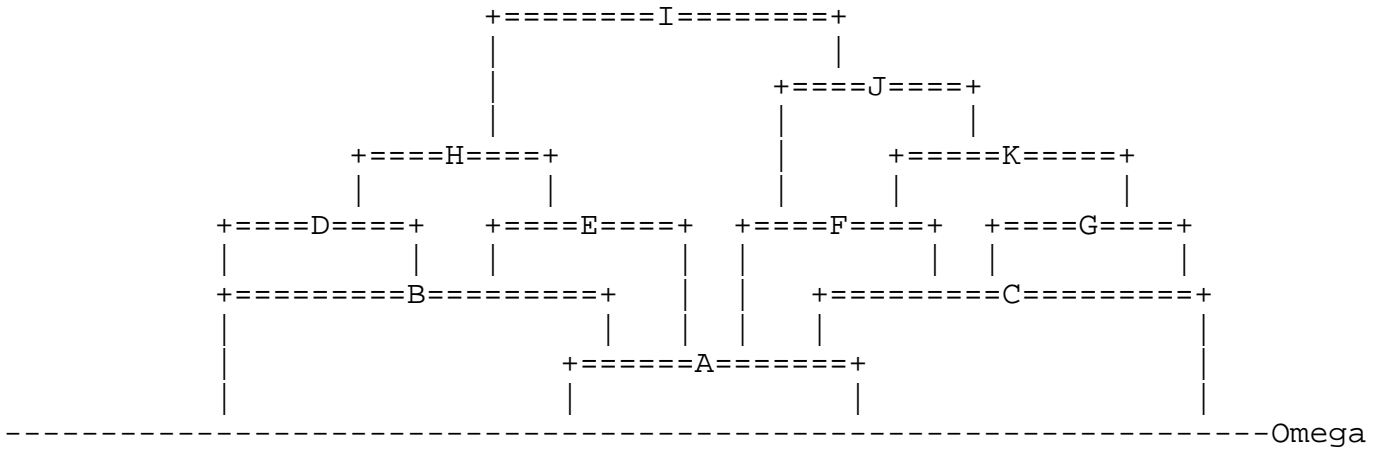


Figure 5: ARC Representation

This representation is similar to a classical routing graph with the peculiarity that some Links are marked reversible. An ARC is represented as a sequence of reversible links.

Collapsed DAG view:

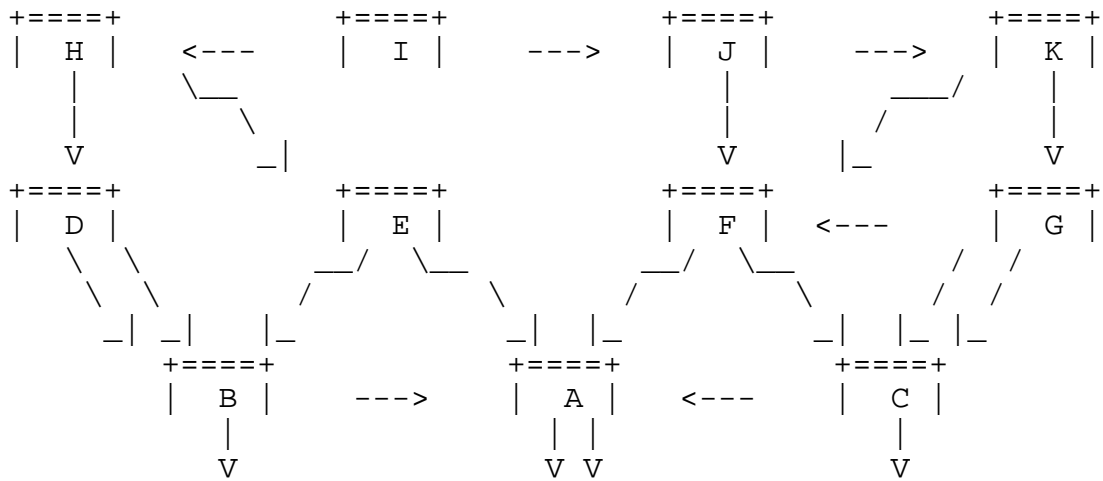


Figure 6: ARC DAG

A DAG representation whereby an ARC is abstracted as a vertice and links between ARCs are shown as directed edges. This way, the reversible links are omitted and the graph is simplified. It can be noted that even the vertice closest to Omega has 2 non-congruent forwarding solutions, that is Heir Links to Omega.

4. Applicability

This section has to be refined. ARCs probably apply to both unicast and multicast and the authors expect further documents to explain how that is done. The examples below are provided as an indication but is not limiting the field of applications.

4.1. Load Balancing

In normal conditions, only the cursor may distribute its traffic between the two Edge Nodes. If an Edge Node is still congested after the cursor forwards all its traffic towards the other Edge Node, then the cursor can be moved towards the congested Edge in order to derive even more traffic towards the other Edge. If both Edges are congested, then a backpressure can be applied on the incoming ARCs so that they move their own traffic towards their own alternate Edge. The process may recurse.

4.1.1. Routing Hierarchies

The ARC methods may be used to build and/or leverage routing hierarchies, allowing high availability at multiple hierarchical levels. In one hand, the view of an ARC Set can be simplified by abstracting an ARC as a node in a DAG. The view of the routing topology is thus simplified, as illustrated in Figure 6. In the other hand, ARCs may be used inside a subtopology, such as a ring, to enable forwarding inside a ring towards a next ring. Then, abstracting a full ring as a node, ARCs can be applied to a graph of rings, providing another level of redundancy and an abstract end to end path computation that is represented as a cascade of ARCs of rings.

5. Lowest ARC First

The open Lowest ARC First (oLAF) algorithm is presented below in such a way as to help the reader figure how an ARC Set can be obtained but not in a computer-optimized fashion that is left to be determined. oLAF is based on Dijkstra's algorithm for Shortest Path First (SPF) computation, and is designed in such a fashion that the reverse SPF tree towards a destination is conserved and preferred for forwarding along the resulting ARC Set.

We make the computation on behalf of Omega, that is an abstraction, but could represent the node or the set of nodes that we want to reach with an ARC Set. If Omega is instantiated as an actual destination node, then that node may be a fine location for an ARC Computing Engine.

5.1. Init

So we start with an proverbial Initial Set of Nodes that are interconnected by Links, and Omega that is the destination that we want to reach with an ARC Set.

If there is no Heir, we're done. If there is a single Heir then the graph is monoconnected, so we restart the computation taking that Heir off the Set of Nodes and making it Omega.

Else, if Omega is a single Node, or if Omega is composed of multiple nodes but we are willing to accept that both ends of an ARC terminate in a same node in Omega, then we create virtual Omega nodes, a minimum of two and at most one per Heir, and we make them the new Omega. Note: we need at least two destinations because both ends of an ARC cannot terminate in a same node.

Now we can start building an ARC Set towards the resulting Omega.

In this process, we create so-called Dependent Sets of nodes, each owned by a Safe Node S , $DSet(S)$. $DSet(S)$ contains nodes that are not determined to be Safe at the current stage of the computation and for which S , the owner Safe Node, is the last single point of failure on the shortest path tree to Omega. It results that a given node can be at most in one $DSet$, and that a Safe Node belongs to its own $DSet$.

For each node S in Omega we create a $DSet(S)$ in which we place S .

5.2. Growing Trees

And then the process goes like this:

We select the node in the Set of Nodes that is closest to Omega using the cost towards Omega as if we were building a traditional reverse SPF tree and we place the selected node in the same Dependent Set as its parent in the reverse SPF tree. Note that for a Heir, the parent might be a real node in Omega, or a virtual Omega node.

If we kept it at that, we would be building subtrees that are hanging off a Safe Node and together would represent the reverse shortest path tree towards Omega, each subtree being grown separately inside $DSet(S)$ where S is the (virtual) Safe node that is the root of the subtree.

5.3. Being Safe

But once we have placed the selected node in a $DSet$, we consider its neighbors one by one. If at least one of the neighbors is already in

a different DSet than this node, we select the neighbor that provides the shortest alternate path to Omega for the selected node.

Doing so, we have isolated two paths:

- o one along its own shortest path that is contained within its own Dependent Set and that leads to the owner Safe Node of this set.
- o and one via the selected neighbor, along its own shortest path within the selected neighbor's Dependent Set and that leads to the owner Safe Node of that other set.

Because the two sets are different and have no intersection, these paths are non-congruent. And because the two non-congruent paths lead to two different Safe Nodes, this node is Safe.

It might happen that:

- o the selected node's parent is already a Safe Node, in which case the selected node is the Edge AN on its shortest path side.
- o It might also happen that the selected neighbor is already a Safe Node, in which case selected node is the Edge AN on its alternate side.

If both conditions are met for a same AN, then that AN forms a collapsed ARC by itself.

5.4. Bending An ARC

Now we form an ARC as follows:

- o A height is attributed to this ARC that must be strictly more than that of the ARCs it terminates into, if any. The order in which the ARCs are built may be used in some cases.
- o The ARC terminates in the two Safe Nodes that are the owners of the two DSets. The normal behaviour is to make a Edge Link the link to the Safe Node.
- o If the Safe Node at one end forms a collapsed ARC by itself, it may be absorbed in the ARC in order to build a multi-edged ARC.
- o If one of the two Safe Nodes pertains in a ARC or a Comb construct that is higher than the other end, then this ARC may be merged at the Safe Node with its original ARC, in order to form a Comb construct whereby this ARC is a Buttrussing ARC of the Comb. The resulting Comb conserves the height on the original ARC or Comb

that it extends.

- o The ARC is built by adjoining the two non-congruent paths that we isolated for the selected node.
- o The selected node is the node farthest from Omega in the resulting ARC, so we make it the cursor.
- o The link between the selected node and the selected neighbor would not have been used in a classical reverse SPF tree. Here, we have determined that this link is in fact critical to connect 2 zones of the network (the DSets) that can act as a back up for one another in case of the failure of their respective single points of failure (the Safe Nodes).
- o Because the ARC can be used in both directions, each AN along the ARC has two non-congruent paths to the Safe Nodes that the ARC terminates into. So it is a Safe Node. We create individual DSets for each AN and we move the AN to its own DSet.

5.5. Orienting Links

For each ARC Node along the ARC:

- o any link (there can be zero for a collapsed ARC, one for an Edge AN or two of them for a Intermediate AN) between this AN and a next AN along this ARC is made an Intermediate Link, that is, reversible. The normal direction, away from the cursor, preserves the shortest path.
- o If this AN is an Edge AN for this ARC, than all links off this node that terminate in a Safe Node are made Edge Links, that is, outgoing but not reversible.
- o All the other links left undertermined.

The nodes left in the Dependent Sets but the owner Safe Node are still not Safe. They are moved back to the original Set of Nodes to enable forming additional ARCs which might depend on this ARC in the ARC Set.

5.6. Looping or recursing

We are done processing the particular node we had picked in the original Set of Nodes. If the Set of Nodes as it stands now is not empty, we continue from Section 5.2.

If the Set of Nodes went empty, we are done with this pass and we

consider the Dependent Sets that we have put together. In a biconnected graph, there should be one set per node and one node per set, denoting that every node is a Safe Node.

If some portion of the graph is monoconnected, then each monoconnected portion forms the Dependent Set of the Safe Node that is its single point of failure. In order to be maximally redundant, we need to form the ARCs again, within the Dependent Set.

To do so, we remove the Safe Node from the Dependent set and make it Omega. We make the resulting DSet our Set of Nodes and run the algorithm again.

This may recurse a number of times if the graph has monoconnected zones within others.

6. Forwarding Along An ARC Set

Under normal conditions, the traffic flows away from the cursor of the current ARC and cascades into the next ARC on that side of the cursor, with the Height of the current ARC decreasing monotonically from ARC to ARC till Omega is reached.

The same goes for a generic Comb construct. When Buttrussing ARCs are applied on a main ARC or other Buttrussing ARCs, the final construct assumes the shape of a tree. The tree may be walked in different manners but the shortest path requires to start going down the current ARC or Buttrussing ARC to its Edge.

In case of Label forwarding, the same recursivity is applied and a multiple ARC label path is constructed. Each ARC has its own set of label path per Omega, each ARC Set label path being merged into the lower ARC label set, thus at the interconnection point. At minimum, ARC label path should be built from the cursor toward each edge, but this would require label path recompilation upon cursor move, the proposed approach is then to build for the normal flow to an Omega one pair of label path from edge to edge.

As this label construct maps the ARC topology with local significant label, the Label Distribution Protocol (LDP) could be reused to announce label association to neighbors on the ARC.

Upon a breakage inside an ARC, until a corrective action takes place, some traffic will be lost. The corrective action might be either operated at the control plane or the data plane, if immediate action and near-zero packet loss is required.

6.1. Control Plane Recovery

Upon a first breakage in an ARC, the cursor is moved to the breakage point, either a node or a link, so that traffic flows away from the cursor again.

Upon a second breakage within a same ARC, a segment of the ARC is now isolated. Both breakage points become sinks till an additional corrective action, such as modifying the ARC Set, takes place. All incoming links in the isolated segment are blocked, causing the traffic to exit at the other end of the incoming ARCs.

Blocking an Edge Link in the incoming ARC may create an isolated segment in the incoming ARC as well if it is a second breakage there too, or if both edges of the incoming ARC terminate in the broken segment. In that case the process recurses and the broken zone can be determined as the collection of the isolated segments.

If a segment of an ARC is getting isolated by a dual failure but that ARC segment has incoming Edges then the ARC can be reversed. This reversal is done by reversing of all the incoming Edges, which become outgoing. The segment that was isolated now benefits from multiple exits in a loop free fashion. This process might in turn isolate a segment of an ARC that was incoming and the process recurses and some links flap. If a real exit exists the process will stabilize, but a count to infinity must be put in place to avoid a permanent flapping when a whole ARC Subset is physically isolated. One may consider that this process is in fact the classical link reversal technique, as applied to the DAG of ARCs.

6.2. Data Plane Recovery

Upon a breakage inside an ARC, it is possible in the data plane to reverse the direction of -to turn- a given packet once along the ARC so the packets exits over the other Edge Link. But in order to avoid loops, it is undesirable to reverse the direction of a given packet a second time.

Note that once a given packet leaves an ARC to enter the next, it is free to bounce again in the next ARC. In other words, the domain that is impacted by a turn is limited to the current ARC itself; the ARC forms the event horizon wherein the notion that a turn happened may cause a loop.

So a local strategy must be put in place inside an ARC to allow a given packet to bounce once upon a breakage, and get dropped upon a second breakage.

In the case of IP packet forwarding, a packet can be tagged when it bounces inside an ARC, or when it passes the cursor, for instance by reserving a TOS bit for that purpose. When the packet bounces, the bit is set and when the packet leaves the ARC, the bit is reset and may be used again in the next ARC. In the generic case of a Comb, a strategy must be put in place to walk the structure and drop a packet that tries all the Edges. it attempts to pass the cursor twice in a same direction, meaning that more than a full walk was already accomplished.

In the case of MPLS forwarding, the same result can be achieved with either 3 or 4 Labels Switched Paths (LSPs) along the ARC.

3-Labels method: In this case we lay a primary LSP from the cursor to the Edge in each direction, and a backup LSP Edge to Edge in each direction. So a node along the way has three labels, one primary and two backup, one in each direction. Should the primary path fail, the packet can be placed along the backup LSP in the other direction. We'll note that this method constrains the location of the cursor. Should the cursor move, The primary LSPs have to be recomputed, at a minimum between the old and the new location of the cursor where the direction is reversed.

4-Labels method: In this case we have two primary and two backup LSPs Edge to Edge in each direction. The labels are independent of the location of the cursor, so the cursor can be moved in control plane with no impact on labels.

7. Manageability

This specification describes a generic model. Protocols and management will come later

8. IANA Considerations

This specification does not require IANA action.

9. Security Considerations

This specification is not found to introduce new security threat.

10. Acknowledgements

The authors wishes to thank Dirk Anteunis, Stewart Bryant, IJsbrand

Wijnands, George Swallow, Eric Osborne, Clarence Filsfils and Eric Levy-Abegnoli for their participation and continuous support to the work presented here.

11. References

11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4291] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", RFC 4291, February 2006.
- [RFC4861] Narten, T., Nordmark, E., Simpson, W., and H. Soliman, "Neighbor Discovery for IP version 6 (IPv6)", RFC 4861, September 2007.
- [RFC4862] Thomson, S., Narten, T., and T. Jinmei, "IPv6 Stateless Address Autoconfiguration", RFC 4862, September 2007.

11.2. Informative References

- [I-D.phinney-roll-rpl-industrial-applicability] Phinney, T., Thubert, P., and R. Assimiti, "RPL applicability in industrial networks", draft-phinney-roll-rpl-industrial-applicability-00 (work in progress), October 2011.
- [I-D.thubert-lowpan-backbone-router] Thubert, P., "LoWPAN Backbone Router", draft-thubert-lowpan-backbone-router-00 (work in progress), November 2007.
- [RFC4191] Draves, R. and D. Thaler, "Default Router Preferences and More-Specific Routes", RFC 4191, November 2005.
- [RFC4541] Christensen, M., Kimball, K., and F. Solensky, "Considerations for Internet Group Management Protocol (IGMP) and Multicast Listener Discovery (MLD) Snooping Switches", RFC 4541, May 2006.
- [RFC5213] Gundavelli, S., Leung, K., Devarapalli, V., Chowdhury, K., and B. Patil, "Proxy Mobile IPv6", RFC 5213, August 2008.
- [RFC5415] Calhoun, P., Montemurro, M., and D. Stanley, "Control And Provisioning of Wireless Access Points (CAPWAP) Protocol

Specification", RFC 5415, March 2009.

- [RFC5865] Baker, F., Polk, J., and M. Dolly, "A Differentiated Services Code Point (DSCP) for Capacity-Admitted Traffic", RFC 5865, May 2010.
- [RFC5921] Bocci, M., Bryant, S., Frost, D., Levrau, L., and L. Berger, "A Framework for MPLS in Transport Networks", RFC 5921, July 2010.
- [RFC6275] Perkins, C., Johnson, D., and J. Arkko, "Mobility Support in IPv6", RFC 6275, July 2011.
- [RFC6550] Winter, T., Thubert, P., Brandt, A., Hui, J., Kelsey, R., Levis, P., Pister, K., Struik, R., Vasseur, JP., and R. Alexander, "RPL: IPv6 Routing Protocol for Low-Power and Lossy Networks", RFC 6550, March 2012.

Authors' Addresses

Pascal Thubert (editor)
Cisco Systems
Village d'Entreprises Green Side
400, Avenue de Roumanille
Batiment T3
Biot - Sophia Antipolis 06410
FRANCE

Phone: +33 497 23 26 34
Email: pthubert@cisco.com

Patrice Bellagamba
Cisco Systems
214 Avenue des fleurs
Saint-Raphael 83700
FRANCE

Phone: +33.6.1998.4346
Email: pbellaga@cisco.com