

Internet Engineering Task Force (IETF)
Request for Comments: 6198
Category: Informational
ISSN: 2070-1721

B. Decraene
France Telecom
P. Francois
UCL
C. Pelsser
IIJ
Z. Ahmad
Orange Business Services
A.J. Elizondo Armengol
Telefonica I+D
T. Takeda
NTT
April 2011

Requirements for the Graceful Shutdown of BGP Sessions

Abstract

The Border Gateway Protocol (BGP) is heavily used in Service Provider networks for both Internet and BGP/MPLS VPN services. For resiliency purposes, redundant routers and BGP sessions can be deployed to reduce the consequences of an Autonomous System Border Router (ASBR) or BGP session breakdown on customers' or peers' traffic. However, simply taking down or even bringing up a BGP session for maintenance purposes may still induce connectivity losses during the BGP convergence. This is no longer satisfactory for new applications (e.g., voice over IP, online gaming, VPN). Therefore, a solution is required for the graceful shutdown of a (set of) BGP session(s) in order to limit the amount of traffic loss during a planned shutdown. This document expresses requirements for such a solution.

Status of This Memo

This document is not an Internet Standards Track specification; it is published for informational purposes.

This document is a product of the Internet Engineering Task Force (IETF). It represents the consensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Not all documents approved by the IESG are a candidate for any level of Internet Standard; see Section 2 of RFC 5741.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <http://www.rfc-editor.org/info/rfc6198>.

Copyright Notice

Copyright (c) 2011 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	2
2. Conventions Used in This Document	3
3. Problem Statement	4
3.1. Example of Undesirable BGP Routing Behavior	4
3.2. Causes of Packet Loss	5
4. Terminology	6
5. Goals and Requirements	7
6. Security Considerations	10
7. References	10
7.1. Normative References	10
7.2. Informative References	10
Acknowledgments	11
Appendix A. Reference BGP Topologies	12
A.1. EBGp Topologies	12
A.2. IBGP Topologies	15
A.3. Routing Decisions	19

1. Introduction

The Border Gateway Protocol (BGP) [RFC4271] is heavily used in Service Provider networks for both Internet and BGP/MPLS VPN services [RFC4364]. For resiliency purposes, redundant routers and BGP sessions can be deployed to reduce the consequences of an Autonomous System Border Router (ASBR) or BGP session breakdown on customers' or peers' traffic.

We place ourselves in the context where a Service Provider performs a maintenance operation and needs to shut down one or multiple BGP peering link(s) or a whole ASBR. If an alternate path is available within the Autonomous System (AS), the requirement is to avoid or reduce customer or peer traffic loss during the BGP convergence.

Indeed, as an alternate path is available in the AS, it should be made possible to reroute the customer or peer traffic on this backup path before the BGP session(s) is/are torn down, the nominal path withdrawn, and the forwarding stopped.

The requirements also cover the subsequent re-establishment of the BGP session as even this "UP" case can currently trigger route loss, and thus traffic loss, at some routers.

BGP [RFC4271] and MP-BGP [RFC4760] do not currently have a mechanism to gracefully migrate traffic from one BGP next-hop to another without interrupting the flow of traffic. When a BGP session is taken down, BGP behaves as if there were a sudden link or router failure and withdraws the prefixes learned over that session, which may trigger traffic loss. While still being advertised as reachable, there is no mechanism to advertise to its BGP peers that the prefix will soon be unreachable. When applicable, such mechanism would reduce or prevent traffic loss. It would typically be applicable in case of a maintenance operation requiring the shutdown of a forwarding resource. Typical examples would be a link or line card maintenance, replacement, or upgrade. It may also be applicable for a software upgrade, as it may involve a firmware reset on the line cards and hence forwarding interruption.

The introduction of route reflectors (RRs) as per [RFC4456] to solve scalability issues bound to Internal BGP (IBGP) full-meshes has worsened the duration of routing convergence as some route reflectors may hide the backup path. Thus, depending on RR topology, more IBGP hops may be involved in the IBGP convergence.

Note that these planned maintenance operations cannot be addressed by Graceful Restart (GR) extensions [RFC4724] as GR only applies when the forwarding is preserved during the control plane restart. On the contrary, graceful shutdown applies when the forwarding is interrupted.

Also, note that some protocols are already considering such a graceful shutdown procedure (e.g., GMPLS in [RFC5817]).

A metric of success is the degree to which such a mechanism eliminates traffic loss during maintenance operations.

2. 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].

3. Problem Statement

As per [RFC4271], when one (or many) BGP session(s) are shut down, a BGP NOTIFICATION message is sent to the peer and the session is then closed. A protocol convergence is then triggered both by the local router and by the peer. Alternate paths to the destination are selected, if known. If those alternate paths are not known prior to the BGP session shutdown, additional BGP convergence steps are required in each AS to search for an alternate path.

This behavior is not satisfactory in a maintenance situation because the traffic that was directed towards the removed next-hops may be lost until the end of the BGP convergence. As it is a planned operation, a make-before-break solution should be made possible.

As maintenance operations are frequent in large networks [Reliable], the global availability of the network is significantly impaired by this BGP maintenance issue.

3.1. Example of Undesirable BGP Routing Behavior

To illustrate these problems, let us consider the following simple example where one customer router "CUST" is dual-attached to two Service Providers' routers, "ASBR1" and "ASBR2".

ASBR1 and ASBR2 are in the same AS and are owned by the same Service Provider. Both are IBGP clients of the route reflector R1.

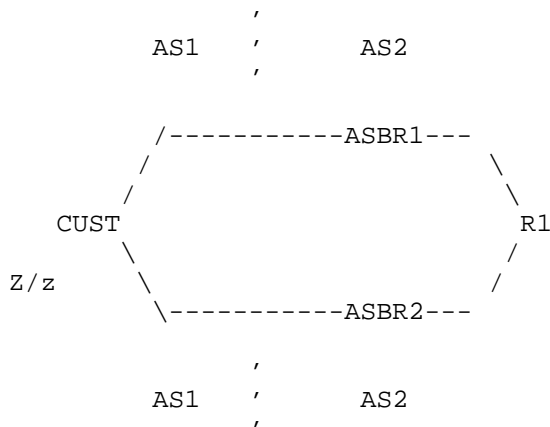


Figure 1. Dual-Attached Customer

Before the maintenance, packets for destination Z/z use the ASBR1-CUST link because R1 selects ASBR1's route based on the IGP cost.

Let's assume the Service Provider wants to shut down the ASBR1-CUST link for maintenance purposes. Currently, when the shutdown is performed on ASBR1, the following steps are performed:

1. ASBR1 withdraws its prefix Z/z to its route reflector, R1.
2. R1 runs its decision process, selects the route from ASBR2, and advertises the new path to ASBR1.
3. ASBR1 runs its decision process and recovers the reachability of Z/z.

Traffic is lost at step 1 when ASBR1 loses its route until step 3 when it discovers a new path.

Note that this is a simplified description for illustrative purposes. In a bigger AS, multiple steps of BGP convergence may be required to find and select the best alternate path (e.g., ASBR1 may be chosen based on a higher LOCAL_PREF, hierarchical route reflectors may be used, etc.). When multiple BGP routers are involved and plenty of prefixes are affected, the recovery process can take longer than application requirements.

3.2. Causes of Packet Loss

The loss of packets during maintenance has two main causes:

- lack of an alternate path on some routers, and
- transient routing inconsistency.

Some routers may lack an alternate path because another router is hiding the backup path. This router can be:

- a route reflector only propagating its best path.
- the backup ASBR not advertising the backup path because it prefers the nominal path.

This lack of knowledge regarding the alternate path is the first target of this requirements document.

Transient routing inconsistencies happen during IBGP convergence because routers do not simultaneously update their Routing Information Bases (RIBs) and hence do not simultaneously update their

Forwarding Information Bases (FIBs) entries. This can lead to forwarding loops, which result in both link congestion and packet drops. The duration of these transient micro-loops is dependent on the IBGP topology (e.g., number of route reflectors between ingress and egress ASBR), implementation differences among router platforms (which result in differences in the time taken to update specific prefix in the FIB), and forwarding mode (hop-by-hop IP forwarding versus tunneling).

Note that when an IP lookup is only performed on entry to the AS, for example, prior to entry into a tunnel across the AS, micro-loops will not occur. An example of this is when BGP is being used as the routing protocol for MPLS VPN as defined in [RFC4364].

Note that [RFC5715] defines a framework for loop-free convergence. It has been written in the context of IP fast reroute for link state IGP [RFC5714], but some concepts are also of interest for BGP convergence.

4. Terminology

g-shut: Graceful shutdown. A method for explicitly notifying the BGP routers that a BGP session (and hence the prefixes learned over that session) is going to be disabled.

g-noshut: Graceful no shutdown. A method for explicitly notifying the BGP routers that a BGP session (and hence the prefixes learned over that session) is going to be enabled.

g-shut initiator: the router on which the session(s) shutdown(s) is (are) performed for maintenance.

g-shut neighbor: a router that peers with the g-shut initiator via (one of) the session(s) undergoing maintenance.

affected prefixes: a prefix initially reached via the peering link(s) undergoing maintenance.

affected router: a router reaching an affected prefix via a peering link undergoing maintenance.

initiator AS: the autonomous system of the g-shut initiator router.

neighbor AS(es): the autonomous system(s) of the g-shut neighbor router(s).

5. Goals and Requirements

Currently, when a BGP session of the router under maintenance is shut down, the router removes the routes and then triggers the BGP convergence on its BGP peers by withdrawing its route.

The goal of BGP graceful shutdown of a (set of) BGP session(s) is to minimize traffic loss during a planned shutdown. Ideally, a solution should reduce this traffic loss to zero.

Another goal is to minimize and, preferably, to eliminate packet loss when the BGP session is re-established following the maintenance.

As the event is known in advance, a make-before-break solution can be used in order to initiate the BGP convergence, find and install the alternate paths before the nominal paths are removed. As a result, before the nominal BGP session is shut down, all affected routers learn and use the alternate paths. Those alternate paths are computed by BGP, taking into account the known status of the network, which includes known failures that the network is processing concurrently with the BGP session graceful shutdown and possibly other known graceful shutdowns under way. Therefore, multiple BGP graceful shutdowns overlapping within a short time frame are gracefully handled. Indeed, a given graceful shutdown takes into account all previous ones.

As a result, provided an alternate path with enough remaining capacity is available, the packets are rerouted before the BGP session termination and fewer packets (possibly none) are lost during the BGP convergence process since, at any time, all routers have a valid path.

From the above goals, we can derive the following requirements:

- a) A mechanism to advertise the maintenance action to all affected routers is REQUIRED. Such a mechanism may be either implicit or explicit. Note that affected routers can be located both in the local AS and in neighboring ASes. Note also that the maintenance action can either be the shutdown of a BGP session or the establishment of a BGP session.

The mechanism SHOULD allow BGP routers to minimize and, preferably, eliminate packet loss when a path is removed or advertised. In particular, it SHOULD be ensured that the old path is not removed from the routing tables of the affected routers before the new path is known.

The solution mechanism MUST significantly reduce and, ideally, eliminate packet loss. A trade-off may be made between the degree of packet loss and the simplicity of the solution.

- b) An Internet-wide convergence is OPTIONAL. However, if the initiator AS and the neighbor AS(es) have a backup path, they SHOULD be able to gracefully converge before the nominal path is shut down.
- c) The proposed solution SHOULD be applicable to any kind of BGP sessions (External BGP (EBGP), IBGP, IBGP route reflector client, EBGP confederations, EBGP multi hop, MultiProtocol BGP extension, etc.) and any address family. If a BGP implementation allows the closing or enabling of a subset of Address Family Identifiers (AFIs) carried in an MP-BGP session, this mechanism MAY be applicable to this subset of AFIs.

Depending on the kind of session, there may be some variations in the proposed solution in order to fulfill the requirements.

The following cases should be handled in priority:

- The shutdown of an inter-AS link and therefore the shutdown of an EBGP session;
- The shutdown of an ASBR and therefore the shutdown of all its BGP sessions.

Service Providers and platforms implementing a graceful shutdown solution should note that in BGP/MPLS VPN as per [RFC4364], the Provider Edge - Customer Edge (PE-CE) routing can be performed by protocols other than BGP (e.g., static routes, RIPv2, OSPF, IS-IS). This is out of scope of this document.

- d) The proposed solution SHOULD NOT change the BGP convergence behavior for the ASes exterior to the maintenance process, namely, ASes other than the initiator AS and its neighbor AS(es).
- e) An incremental deployment on a per-AS or per-BGP session basis MUST be made possible. In case of partial deployment, the proposed solution SHOULD incrementally improve the maintenance process. It should be noted that in an inter-domain relation, one AS may have more incentive to use graceful shutdown than the other. Similarly, in a BGP/MPLS VPN environment, it's much easier to upgrade the PE routers than the CE ones, mainly because there is at least an order of magnitude more CE and CE locations than PE and PE locations. As a consequence, when

splitting the cost of the solution between the g-shut initiator and the g-shut neighbor, the solution SHOULD favor a low-cost solution on the neighbor AS side in order to reduce the impact on the g-shut neighbor. Impact should be understood as a generic term that includes first hardware, then software, then configuration upgrade.

- f) Redistribution or advertisement of (static) IP routes into BGP SHOULD also be covered.
- g) The proposed solution MAY be designed in order to avoid transient forwarding loops. Indeed, forwarding loops increase packet transit-delay and may lead to link saturation.
- h) The specific procedure SHOULD end when the BGP session is closed following the g-shut and once the BGP session is gracefully opened following the g-noshut. In the end, once the planned maintenance is finished, the nominal BGP routing MUST be re-established. The duration of the g-shut procedure, and hence the time before the BGP session is safely closed, SHOULD be discussed by the solution document. Examples of possible solutions are the use of a pre-configured timer, the use of a message to signal the end of the BGP convergence, or the monitoring of the traffic on the g-shut interface.
- i) The solution SHOULD be simple and simple to operate. Hence, it MAY only cover a subset of the cases. As a consequence, most of the above requirements are expressed as "SHOULD" rather than "MUST".

The metrics to evaluate and compare the proposed solutions are:

- The duration of the remaining loss of connectivity when the BGP session is brought down or up;
- The applicability to a wide range of BGP and network topologies;
- The simplicity;
- The duration of transient forwarding loops;
- The additional load introduced in BGP (e.g., BGP messages sent to peer routers, peer ASes, the Internet).

6. Security Considerations

At the requirements stage, this graceful shutdown mechanism is not expected to affect the security of the BGP protocol, especially if it can be kept simple. No new sessions are required and the additional ability to signal the graceful shutdown is not expected to bring additional attack vectors, as BGP neighbors already have the ability to send incorrect or misleading information or even shut down the session.

Security considerations MUST be addressed by the proposed solutions. In particular, they SHOULD address the issues of bogus g-shut messages and how they would affect the network(s), as well as the impact of hiding a g-shut message so that g-shut is not performed.

The solution SHOULD NOT increase the ability of one AS to selectively influence routing decision in the peer AS (inbound Traffic Engineering) outside of the case of the BGP session shutdown. Otherwise, the peer AS SHOULD have means to detect such behavior.

7. References

7.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, January 2006.
- [RFC4760] Bates, T., Chandra, R., Katz, D., and Y. Rekhter, "Multiprotocol Extensions for BGP-4", RFC 4760, January 2007.
- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", RFC 4456, April 2006.
- [RFC4364] Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", RFC 4364, February 2006.

7.2. Informative References

- [RFC5817] Ali, Z., Vasseur, JP., Zamfir, A., and J. Newton, "Graceful Shutdown in MPLS and Generalized MPLS Traffic Engineering Networks", RFC 5817, April 2010.

- [RFC5715] Shand, M. and S. Bryant, "A Framework for Loop-Free Convergence", RFC 5715, January 2010.
- [RFC5714] Shand, M. and S. Bryant, "IP Fast Reroute Framework", RFC 5714, January 2010.
- [RFC4724] Sangli, S., Chen, E., Fernando, R., Scudder, J., and Y. Rekhter, "Graceful Restart Mechanism for BGP", RFC 4724, January 2007.
- [Reliable] Network Strategy Partners, LLC. "Reliable IP Nodes: A prerequisite to profitable IP services", November 2002.
http://www.nspllc.com/NewPages/Reliable_IP_Nodes.pdf

Acknowledgments

The authors would like to thank Nicolas Dubois, Benoit Fondeviole, Christian Jacquenet, Olivier Bonaventure, Steve Uhlig, Xavier Vinet, Vincent Gillet, Jean-Louis le Roux, Pierre Alain Coste, and Ronald Bonica for their useful discussions on this subject, review, and comments.

This document has been partly sponsored by the European project IST AGAVE.

Appendix A. Reference BGP Topologies

This section describes some frequent BGP topologies used both within the AS (IBGP) and between ASes (EBGP). Solutions should be applicable to the following topologies and their combinations.

A.1. EBGP Topologies

This section describes some frequent BGP topologies used between ASes. In each figure, a line represents a BGP session.

A.1.1. One ASBR in AS1 Connected to Two ASBRs in the Neighboring AS2

In this topology, we have an asymmetric protection scheme between AS1 and AS2:

- On the AS2 side, two different routers are used to connect to AS1.
- On the AS1 side, one single router with two BGP sessions is used.



Figure 2. EBGP Topology with Redundant ASBR in One of the ASes

BGP graceful shutdown is expected to be applicable for the maintenance of:

- one of the routers of AS2;
- one link between AS1 and AS2, performed either on an AS1 or AS2 router.

A.2. IBGP Topologies

This section describes some frequent BGP topologies used within an AS. In each figure, a line represents a BGP session.

A.2.1. IBGP Full-Mesh

In this topology, we have a full-mesh of IBGP sessions:

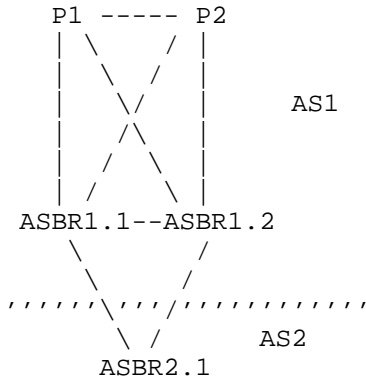


Figure 5. IBGP Full-Mesh

When the session between ASBR1.1 and ASBR2.1 is gracefully shut down, it is required that all affected routers of AS1 reroute traffic to ASBR1.2 before the session between ASBR1.1 and ASBR2.1 is shut down.

Similarly, when the session between ASBR1.1 and ASBR2.1 is gracefully brought up, all affected routers of AS1 preferring ASBR1.1 over ASBR1.2 need to reroute traffic to ASBR1.1 before the less preferred path through ASBR1.2 is possibly withdrawn.

A.2.2. Route Reflector

In this topology, route reflectors are used to limit the number of IBGP sessions. There is a single level of route reflectors and the route reflectors are fully meshed.

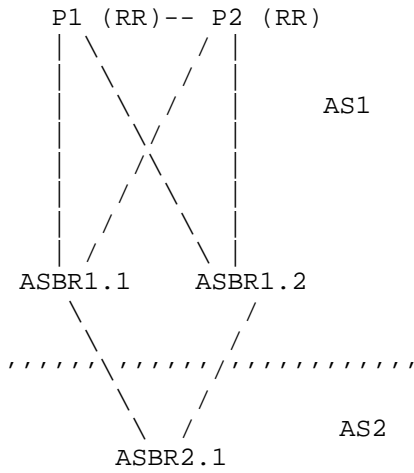


Figure 6. Route Reflector

When the session between ASBR1.1 and ASBR2.1 is gracefully shut down, all BGP routers of AS1 need to reroute traffic to ASBR1.2 before the session between ASBR1.1 and ASBR2.1 is shut down.

Similarly, when the session between ASBR1.1 and ASBR2.1 is gracefully brought up, all affected routers of AS1 preferring ASBR1.1 over ASBR1.2 need to reroute traffic to ASBR1.1 before the less preferred path through ASBR1.2 is possibly withdrawn.

A.2.3. Hierarchical Route Reflector

In this topology, hierarchical route reflectors are used to limit the number of IBGP sessions. There could be more than two levels of route reflectors and the top-level route reflectors are fully meshed.

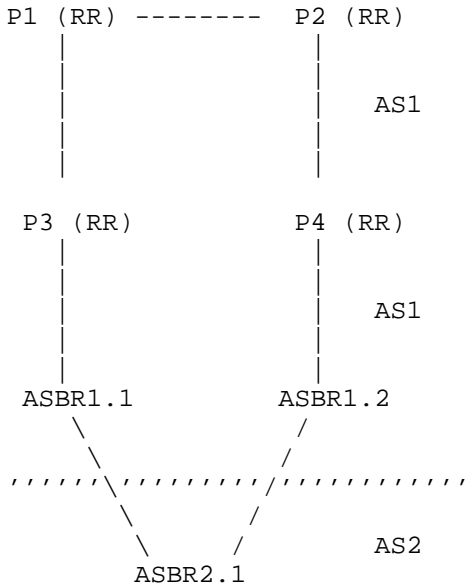


Figure 7. Hierarchical Route Reflector

When the session between ASBR1.1 and ASBR2.1 is gracefully shut down, all BGP routers of AS1 need to reroute traffic to ASBR1.2 before the session between ASBR1.1 and ASBR2.1 is shut down.

Similarly, when the session between ASBR1.1 and ASBR2.1 is gracefully brought up, all affected routers of AS1 preferring ASBR1.1 over ASBR1.2 need to reroute traffic to ASBR1.1 before the less preferred path through ASBR1.2 is possibly withdrawn.

A.2.4. Confederations

In this topology, a confederation of ASes is used to limit the number of IBGP sessions. Moreover, RRs may be present in the member ASes of the confederation.

Confederations may be run with different sub-options. Regarding the IGP, each member AS can run its own IGP or they can all share the same IGP. Regarding BGP, LOCAL_PREF may or may not cross the member AS boundaries.

A solution should support the graceful shutdown and graceful bringing up of EBGP sessions between member ASes in the confederation in addition to the graceful shutdown and graceful bringing up of EBGP sessions between a member-AS and an AS outside of the confederation.

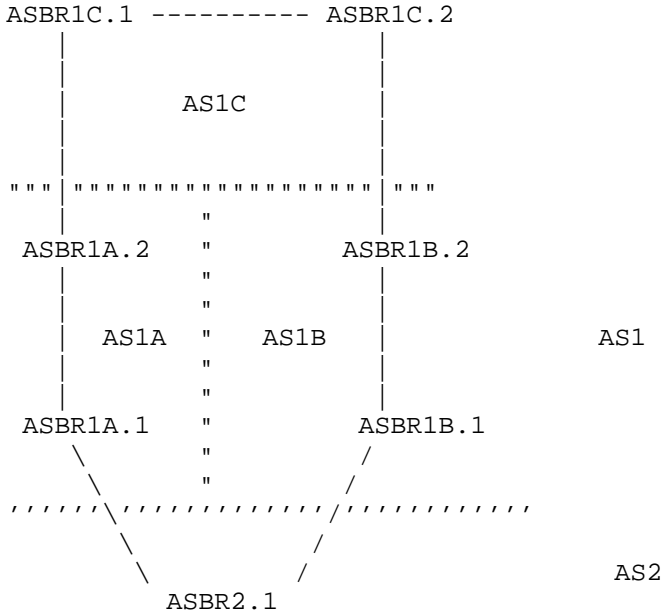


Figure 8. Confederation

In the above figure, member ASes AS1A, AS1B, and AS1C belong to a confederation of ASes in AS1. AS1A and AS1B are connected to AS2.

In normal operation, for the traffic toward AS2:

- AS1A sends the traffic directly to AS2 through ASBR1A.1.
- AS1B sends the traffic directly to AS2 through ASBR1B.1.
- AS1C load balances the traffic between AS1A and AS1B.

When the session between ASBR1A.1 and ASBR2.1 is gracefully shut down, all BGP routers of AS1 need to reroute traffic to ASBR1B.1 before the session between ASBR1A.1 and ASBR2.1 is shut down.

Similarly, when the session between ASBR1A.1 and ASBR2.1 is gracefully brought up, all affected routers of AS1 preferring ASBR1A.1 over ASBR1B.1 need to reroute traffic to ASBR1A.1 before the less preferred path through ASBR1B.1 is possibly withdrawn.

A.3. Routing Decisions

Here we describe some routing engineering choices that are frequently used in ASes and that should be supported by the solution.

A.3.1. Hot Potato (IGP Cost)

The ingress router selects the nominal egress ASBR (AS exit point) based on the IGP cost to reach the BGP next-hop.

A.3.2. Cold Potato (BGP LOCAL_PREF)

The ingress router selects the nominal egress ASBR based on the BGP LOCAL_PREF value set and advertised by the exit point.

A.3.3. Cold Potato (BGP Preference Set on Ingress)

The ingress router selects the nominal egress ASBR based on preconfigured policy information. (Typically, this is done by locally setting the BGP LOCAL_PREF based on the BGP communities attached on the routes).

As per [RFC4271], note that if tunnels are not used to forward packets between the ingress and egress ASBR; this can lead to persistent forwarding loops.

Authors' Addresses

Bruno Decraene
France Telecom
38-40 rue du General Leclerc
92794 Issy Moulineaux cedex 9
France
EMail: bruno.decraene@orange-ftgroup.com

Pierre Francois
Universite catholique de Louvain
Place Ste Barbe, 2
Louvain-la-Neuve 1348
BE
EMail: francois@info.ucl.ac.be

Cristel Pelsser
Internet Initiative Japan
Jinbocho Mitsui Building
1-105 Kanda jinbo-cho
Chiyoda-ku, Tokyo 101-0051
Japan
EMail: cristel@iij.ad.jp

Zubair Ahmad
Orange Business Services
13775 McLearen Road, Oak Hill VA 20171
USA
EMail: zubair.ahmad@orange-ftgroup.com

Antonio Jose Elizondo Armengol
Division de Analisis Tecnologicos
Technology Analysis Division
Telefonica I+D
C/ Emilio Vargas 6
28043, Madrid
EMail: ajea@tid.es

Tomonori Takeda
NTT Corporation
9-11, Midori-Cho 3 Chrome
Musashino-Shi, Tokyo 180-8585
Japan
EMail: takeda.tomonori@lab.ntt.co.jp

