

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 28, 2009

P. Marques
R. Fernando
Juniper Networks
E. Chen
P. Mohapatra
Cisco Systems
July 27, 2008

**Advertisement of the best-external route to IBGP
draft-marques-idr-best-external-00.txt**

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on January 28, 2009.

Abstract

This document makes a case and provides the rules for a border router to advertise its best external route towards its IBGP peers when its overall best is a route received from an IBGP peer.

The best external route may be different from the overall best route installed in the Loc-Rib. Advertising the best-external route (when different from the overall best route) into an IBGP helps in speeding up routing convergence, has positive effects in reducing inter-domain churn and in some limited scenarios could help avoid permanent IBGP

route oscillation.

The document also extends this mechanism to route reflectors and confederation border routers to advertise a best route that is external to the cluster/domain.

Table of Contents

| | | |
|----------------------|--|--------------------|
| 1. | Introduction | 3 |
| 1.1. | Requirements Language | 3 |
| 2. | Consistency between routing and forwarding | 3 |
| 3. | Algorithm for selection of best-external route | 5 |
| 4. | Route Reflection | 6 |
| 5. | Confederations | 6 |
| 6. | Applications | 7 |
| 6.1. | Fast Connectivity Restoration | 7 |
| 6.2. | Inter-Domain Churn Reduction | 7 |
| 6.3. | Reducing Persistent IBGP oscillation | 7 |
| 7. | Acknowledgments | 8 |
| 8. | IANA Considerations | 8 |
| 9. | Security Considerations | 8 |
| 10. | Normative References | 8 |
| | Authors' Addresses | 8 |
| | Intellectual Property and Copyright Statements | 10 |

1. Introduction

The term best-external route describes the most preferred route among the routes received by a router from its EBGPeers. The best-external route might differ from the overall route installed in the Loc-RIB in the case when the overall best route happens to be an internal route. Advertising the best-external route, when different from the overall best, presents additional information into an IBGP mesh which may be of value for several purposes including:

- o Faster restoration of connectivity, by providing additional paths, that may be used to fail over in case the primary path becomes invalid or is withdrawn.
- o Reducing inter-domain churn and traffic blackholing due to the readily available alternate path.
- o Reducing the potential for situations of permanent IBGP route oscillation, as discussed in some scenarios [[RFC3345](#)].
- o Improving selection of lower MED routes from the same neighboring AS.

In current networks, BGP is typically deployed in topologies that include the use of route reflectors [[RFC4456](#)] and/or confederations [[RFC5065](#)]. It is straightforward to extend the concept of "external" route to a cluster or confed sub-AS. A route is considered "external" if it has not been received from the cluster/sub-AS which is being considered for advertisement.

1.1. 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 [RFC 2119](#) [[RFC2119](#)].

2. Consistency between routing and forwarding

The BGP protocol, as defined in [[RFC1771](#)], specifies that a BGP speaker shall advertise to its internal peers the route with the highest degree of preference among routes to the same destination received from external neighbors.

This section discusses problems present with the approach described in [[RFC1771](#)] and the next section offers an alternative algorithm to select a best external route which can be advertised to an IBGP mesh.

The internal update advertisement rules contained in the original BGP-4 specification [[RFC1771](#)] can lead to situations where traffic is forwarded through a route other than the route advertised by BGP.

Inconsistencies between forwarding and routing are highly undesirable. Service providers use BGP with the dual objective of learning reachability information and expressing policy over network resources. The latter assumes that forwarding follows routing information.

Consider the Autonomous system presented in figure 1, where r1 ... r4 are members of a single IBGP mesh and routes a, b, and c are received from external peers.

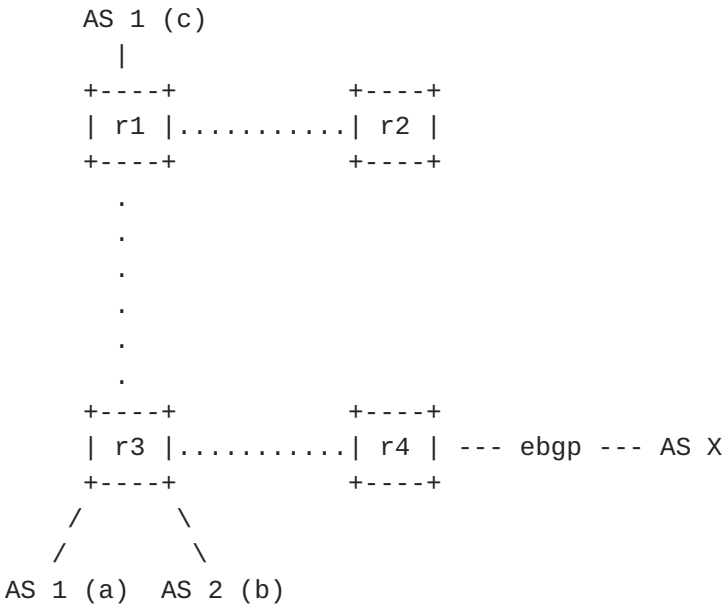


Figure 1: Inconsistency in Routing

| Path | AS | MED | rtr_id |
|------|----|-----|--------|
| a | 1 | 10 | 1 |
| b | 2 | 5 | 10 |
| c | 1 | 5 | 5 |

Figure 2: Path Attribute Table

Following the rules as specified in [[RFC1771](#)], router r3 will select path (b) received from AS 2 as its overall best to install in the

Loc-Rib, since path (b) is preferable to path (c), the lowest MED route from AS 1. However for the purposes of Internal Update route selection, it will ignore the presence of path (c), and elect (a) as its advertisement, via the router-id tie-breaking rule.

In this scenario, router r4 will receive (c) from r1 and (a) from r3. It will pick the lowest MED route (c) and advertise it out via ebgp to AS X. However at this point routing is inconsistent with forwarding as traffic received from AS X will be forwarded towards AS 2, while the ebgp advertisement is being made for an AS 1 path.

Routing policies are typically specified in terms of neighboring ASes. In the situation above, assuming that AS 1 is network for which this AS provides transit services while AS 2 and AS X are peer networks, one can easily see how the inconsistency between routing and forwarding would lead to transit being inadvertently provided between AS X and AS 2. This could lead to persistent forwarding loops.

Inconsistency between routing and forwarding may happen, whenever a bgp speaker chooses to advertise an external route into IBGP that is different from the overall best route and its overall best is external.

3. Algorithm for selection of best-external route

Given that the intent in advertising an external route, when the overall best for the same destination is an internal route, is to provide additional information into the IBGP mesh into which a route is participating, it is desirable to take into account the routes received from interior neighbors in the selection process.

We propose a route selection algorithm that selects a global order between routes and which selects the same overall best route as the one currently specified [[RFC4271](#)].

In order to achieve this we need to introduce the concept of route group. A route group is a set of routes to the same destination received from the same neighboring AS and which is equal in terms of route selection prior to the MED comparison step.

Routes are ordered within a group via MED or subsequent route selection rules.

The order of all routes for the same destination is determined by the order of the best route in each group.

As an example, the following set of received routes:

| Path | AS | MED | rtr_id |
|------|----|-----|--------|
| a | 1 | 10 | 10 |
| b | 2 | 5 | 1 |
| c | 1 | 5 | 5 |
| d | 2 | 20 | 20 |
| e | 2 | 30 | 30 |
| f | 3 | 10 | 20 |

Figure 3: Path Attribute Table - 2

Would yield the following order (from the most to the least preferred):

$b < d < e < c < a < f$

In this example, comparison of the best route within each group provides the sequence ($b < c < f$). The remaining routes are ordered in relation to their respective group best.

The route to be advertised to the IBGP mesh or a given cluster/sub-AS is selected by choosing the most preferred route that is external to that particular domain. Note that whenever the overall best route is external it will automatically be selected by this algorithm.

4. Route Reflection

A route reflector that chooses to implement this algorithm, will advertise to its non-client IBGP peers, the most preferred path received from its clients. This is referred to as the best intra-cluster route. It will advertise to its client peers the most preferred path received from a neighbor outside the cluster. This is referred to as the best inter-cluster route.

In order for a reflector to be able to advertise the best of its inter-cluster routes into a cluster it is necessary that client-to-client reflection be disabled, since its advertisement may otherwise

5. Confederations

When a BGP speaker is configured as a confederation border router, it shall consider the best-external route as follows:

- o When advertising into its sub-AS, it should select the most preferred route not received from within its sub-AS.
- o o When advertising into confed ebgp, it should select the most preferred route not received from the neighboring sub-AS.

6. Applications

6.1. Fast Connectivity Restoration

When two exits are available to reach a particular destination and one is preferred over the other, the availability of an alternate path provides fast connectivity restoration when the primary path fails.

Restoration can be quick since the alternate path is already at hand. The border router could precompute the backup route and preinstall it in FIB ready to be switched when the primary goes away. Note that this requires the border router that's the backup to also preinstall the secondary path and switch to it on failure.

6.2. Inter-Domain Churn Reduction

Within an AS, the non availability of backup best leads to a border router sending a withdraw upstream when the primary fails. This leads to inter-domain churn and packet loss for the time the network takes to converge to the alternate path. Having the alternate path will reduces the churn and eliminates packet loss.

6.3. Reducing Persistent IBGP oscillation

Advertising the best-external route, according to the algorithm described in this document will reduce the possibility of route oscillation by introducing additional information into the IBGP system.

For a permanent oscillation condition to occur, it is necessary that a circular dependency between paths occurs such that the selection of a new best path by a router, in response to a received IBGP advertisement, causes the withdrawal of information that another router depends on in order to generate the original event.

In vanilla BGP, when only the best overall route is advertised, as in most implementations, oscillation can occur whenever there are 2 or clusters/sub-ASes such that at least one cluster has more than one path that can potentially contribute to the dependency.

7. Acknowledgments

This document greatly benefits from the comments of Yakov Rekhter, John Scudder and Jenny Yuan.

8. IANA Considerations

This document has no actions for IANA.

9. Security Considerations

There are no additional security risks introduced by this design.

10. Normative References

- [RFC1771] Rekhter, Y. and T. Li, "A Border Gateway Protocol 4 (BGP-4)", [RFC 1771](#), March 1995.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC3345] McPherson, D., Gill, V., Walton, D., and A. Retana, "Border Gateway Protocol (BGP) Persistent Route Oscillation Condition", [RFC 3345](#), August 2002.
- [RFC4271] Rekhter, Y., Li, T., and S. Hares, "A Border Gateway Protocol 4 (BGP-4)", [RFC 4271](#), January 2006.
- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", [RFC 4456](#), April 2006.
- [RFC5065] Traina, P., McPherson, D., and J. Scudder, "Autonomous System Confederations for BGP", [RFC 5065](#), August 2007.

Authors' Addresses

Pedro Marques
Juniper Networks
1194 N. Mathilda Ave
Sunnyvale, CA 94089
USA

Phone:
Email: roque@juniper.net

Rex Fernando
Juniper Networks
1194 N. Mathilda Ave
Sunnyvale, CA 94089
USA

Phone:
Email: rex@juniper.net

Enke Chen
Cisco Systems
170 W. Tasman Drive
San Jose, CA 95134
USA

Phone:
Email: enkechen@cisco.com

Pradosh Mohapatra
Cisco Systems
170 W. Tasman Drive
San Jose, CA 95134
USA

Phone:
Email: pmohapat@cisco.com

Full Copyright Statement

Copyright (C) The IETF Trust (2008).

This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY, THE IETF TRUST AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Intellectual Property

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in [BCP 78](#) and [BCP 79](#).

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at ietf-ipr@ietf.org.

