

idr
Internet-Draft
Expires: August 28, 2008

B. Dickson
Afilias Canada, Inc
February 25, 2008

**Enhanced BGP Capabilities for Exchanging Second-best and Back-up Paths
draft-dickson-idr-second-best-backup-01**

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 August 28, 2008.

Copyright Notice

Copyright (C) The IETF Trust (2008).

Abstract

This Internet Draft describes an enhanced way to exchange prefix information, to permit multiple copies of a prefix with different paths to be announced and withdrawn.

This negotiated capability provides faster local (inter-AS) and global (intra-AS) convergence, reduces path-hunting, improves route-reflector behaviour, including eliminating both persistent oscillations and BGP "wedgies".

Additional prefix instances have new optional BGP attributes, to control path selection.

Withdrawal of prefixes will require new attributes to disambiguate prefix instances.

Benefits are seen both when deployed intra-AS, and on inter-AS peering.

Author's Note

This Internet Draft is intended to result in this draft or a related draft(s) being placed on the Standards Track for idr.

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 [5].

Intended Status: Proposed Standard.

Table of Contents

- [1. Background](#) [4](#)
- [1.1. Localized Information](#) [4](#)
- [1.2. The Withdrawl Problem](#) [5](#)
- [1.3. The Uniqueness Problem](#) [5](#)
- [2. Proposed Changes](#) [6](#)
- [2.1. New Negotiated Option: USE_SECOND_BEST_AND_BACKUP](#) [6](#)
- [2.2. New Optional Path Attribute: SECOND_BEST](#) [6](#)
- [2.3. New Optional Path Attribute: BACKUP_ONLY](#) [6](#)
- [2.4. New Optional Path Attribute: BACKUP_ONLY_SECOND_BEST](#) [6](#)
- [2.5. New Update Format](#) [7](#)
- [2.6. New Withdraw Format](#) [7](#)
- [3. Modifications to BGP Behavior](#) [10](#)
- [3.1. Changes to Path Selection Rules](#) [10](#)
- [3.2. Second Best - Basic Method](#) [11](#)
- [3.3. Second Best - Route Reflector](#) [11](#)
- [3.4. Second Best - Inter-AS Hybrid Method](#) [11](#)
- [3.5. Backup Only - Basic Method](#) [11](#)
- [3.6. Backup Only - Route Reflector](#) [12](#)
- [3.7. IBGP vs EBGp](#) [12](#)
- [4. Implementation Guidelines](#) [13](#)
- [5. Security Considerations](#) [15](#)
- [6. IANA Considerations](#) [16](#)
- [7. Acknowledgements](#) [17](#)
- [8. References](#) [18](#)
- [8.1. Normative References](#) [18](#)
- [8.2. Informative References](#) [18](#)
- [Appendix A. Path-Hunting Examples](#) [19](#)
- [Appendix B. Persistent Oscillation Examples](#) [20](#)
- [Appendix C. BGP Wedgie Examples](#) [22](#)
- [Author's Address](#) [25](#)
- [Intellectual Property and Copyright Statements](#) [26](#)

1. Background

Even when all the best current practises are observed, operational problems may be experienced when running a BGP network.

These include slow convergence due to "path-hunting", persistent oscillations [[1](#)], and BGP "wedgies" [[2](#)].

Standardization of MRAI timers helps this, as well as [RFC 5004](#) [[4](#)].

These RFCs identify the above issues as needing further work.

1.1. Localized Information

The problems listed above occur as a result of additional information not being available (either on a transient basis, or permanently.)

In the case of "path hunting", the information needed for achieving a stable final state is eventually received, but until it is, sub-optimal forwarding will occur, and possibly even transient routing loops.

The "problem" mechanisms involved are:

- o the suppression of announcement of "second-best" paths, because of IBGP-received "best" paths;
- o the suppression by route-reflectors, of IBGP non-best paths (i.e. those normally seen directly by IBGP peers)
- o the suppression of announcement of "second-best" paths, because of EBGP-received "best" paths.
- o the lack of explicit global mechanism for expressing de-prefering announcements via "back-up" providers.

When a prefix+path received is better than the local "best", the new "best" is normally sent.

However, once a new "best" is received, the side-effect is to force the speaker to WITHDRAW the previous best path within the same "regime" (IBGP mesh or EBGP peers).

When we consider the extra (e.g. suppressed) information, with special rules on what to send and how to treat it, the specified problems may go away, or be reduced in scope, duration, or likelihood.

1.2. The Withdrawl Problem

When a prefix (plus path) is withdrawn, the desired stable state is for the next-best path for that prefix (if one exists) to be chosen at each BGP speaker per its local policy.

If that second-best path is already on hand, the delay and intermediate states can be reduced or entirely avoided. This is especially true for both intra-AS and inter-AS "path hunting".

To avoid inconsistent behavior, routing loops, and routing-information loops, the second-best path received from a neighbor, should never be selected as a best path locally.

The second-best path from a neighbor **MUST ONLY** be considered as a candidate for best path, when the previous best path from that neighbor is withdrawn. When this occurs, the path in question is promoted to "best" status.

1.3. The Uniqueness Problem

Currently, for each prefix, only one path for that prefix is ever announced from one peer to another (except in the instance of Route Reflectors). Because of this property, uniqueness, a withdrawl on a prefix does not require path information. This also means that a change of best path is accomplished via an update for a prefix with the new path information.

If, however, more than one path for a given prefix was sent, then any attempt to withdraw a prefix+path would require that the specific path for the prefix being withdrawn be supplied in the withdrawl update message.

In an environment where multiple paths per prefix are possible, but only one path per prefix is maintained, then two steps would be involved in changing the "best" path. In no particular order, that would be the withdrawl of the old prefix+path, and the announcement of the new prefix+path.

2. Proposed Changes

2.1. New Negotiated Option: USE_SECOND_BEST_AND_BACKUP

This is a new BGP Capabilities value, which can be optionally included in the capabilities negotiation. The specific value is a code-point to be assigned by IANA.

When negotiated:

- o Update messages MUST be in the new format
- o Updates without any of the new optional attributes are considered BEST
- o For each prefix, at most one of each type (BEST, SECOND_BEST, BACKUP_ONLY, BACKUP_ONLY_SECOND_BEST) may be sent

2.2. New Optional Path Attribute: SECOND_BEST

This is a new BGP Path Attribute type. It MAY be used only if the USE_SECOND_BEST_AND_BACKUP capability has been negotiated. The type value is a new code point to be assigned by IANA.

This is an Optional, Non-Transitive, Non-Extended, Non-Partial attribute. All the "attr flag bits" (from BGP [3]) are zero. The length is 1, and the value is 1.

2.3. New Optional Path Attribute: BACKUP_ONLY

This is a new BGP Path Attribute type. The type value is a new code point to be assigned by IANA. This is an Optional, Transitive, Non-Extended, Non-Partial attribute, with the "attr flag bits" (from BGP [3]) set to appropriate values. The length is 1, and the value is 1.

2.4. New Optional Path Attribute: BACKUP_ONLY_SECOND_BEST

This is a new BGP Path Attribute type. It MAY be used only if the USE_SECOND_BEST_AND_BACKUP capability has been negotiated. The type value is a new code point to be assigned by IANA.

This is an Optional, Non-Transitive, Non-Extended, Non-Partial attribute. All the "attr flag bits" (from BGP [3]) are zero. The length is 1, and the value is 1.

2.5. New Update Format

Update messages are identical to existing format, with the exception of the new Withdrawl format, and the new optional Path Attributes (SECOND_BEST ,BACKUP_ONLY, and.BACKUP_ONLY_SECOND_BEST). If BGP capability USE_SECOND_BEST_AND_BACKUP has been negotiated, any Update MAY have a Path Attribute(s) which include SECOND_BEST, BACKUP_ONLY, and/or BACKUP_ONLY_SECOND_BEST. More than one instance of a given prefix, with distinct values of Path Attributes, MAY be sent between BGP speakers.

At most four instances may be sent, specifically one of each combination of with/without SECOND_BEST and BACKUP_ONLY and BACKUP_ONLY_SECOND_BEST: One with neither, one with SECOND_BEST only, one with just BACKUP_ONLY, and one with BACKUP_ONLY_SECOND_BEST. both SECOND_BEST and BACKUP_ONLY.

Two prefix paths are considered identical if they differ only in the presence or absence of any of the new attributes. An Update which contains a path which differs by either or both of these, will result in the path information for the prefix being modified.

2.6. New Withdraw Format

Since it is no longer possible to identify which instance of an prefix is affected by an update containing a withdrawl, a new format for Withdrawls is needed. For simplicity of implementations, this consists of four Withdrawl sections, one for each of the types (BEST, SECOND_BEST, BACKUP_ONLY, BACKUP_ONLY_SECOND_BEST). They occur in REVERSE order, to simplify state transitions if/when a "BEST" path is withdrawn. Each Withdrawl section has the same format as the original Withdrawl section.

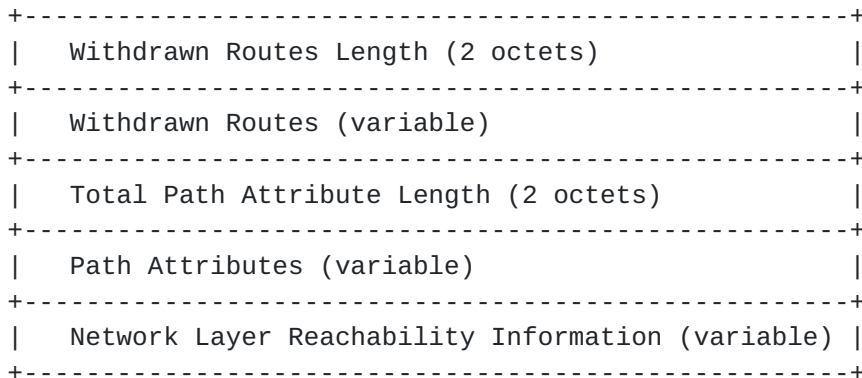


Figure 1

Withdrawn Routes Length: This 2-octets unsigned integer indicates the total length of the Withdrawn Routes field in octets. Its value allows the length of the Network Layer Reachability Information field to be determined, as specified below.

A value of 0 indicates that no routes are being withdrawn from service, and that the WITHDRAWN ROUTES field is not present in this UPDATE message.

Withdrawn Routes Field: This field now consists of four sub-fields and their respective lengths. The value for Withdrawn Routes Length above, must be the sum of the four lengths, plus 8 (the sum of the lengths of the Subfield Lengths).

The format and sequence of the subfields is as follows:

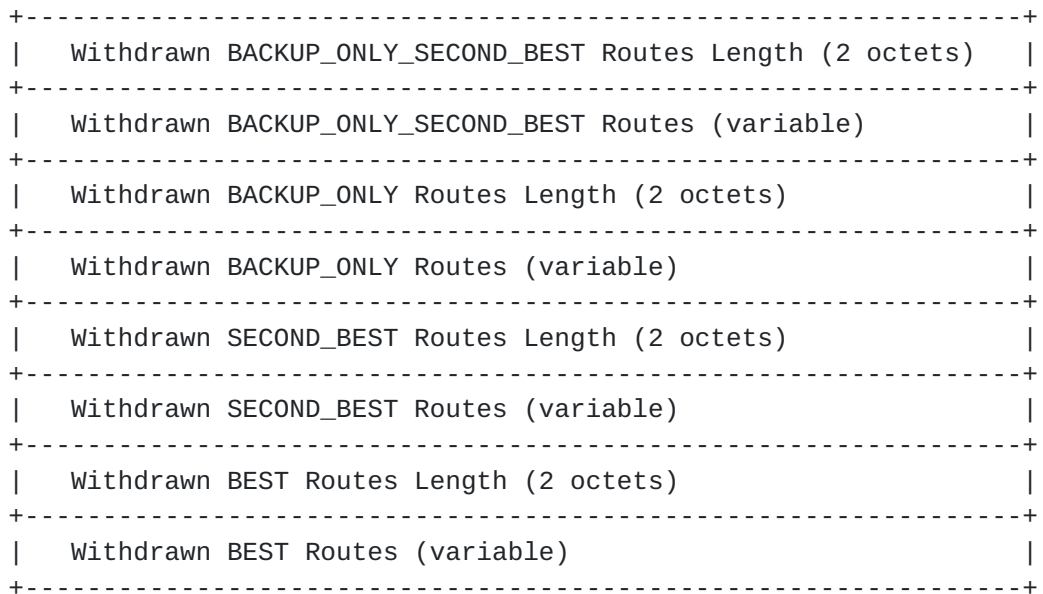


Figure 2

Withdrawn Routes Subfield Lengths These 2-octets unsigned integers indicates the total length of their respective Withdrawn Routes subfields in octets.

Withdrawn Routes Subfields: Each of these is a variable-length field that contains a list of IP address prefixes for the routes that are being withdrawn from service. Each IP address prefix is encoded as a 2-tuple of the form <length, prefix>, whose fields are described below:


```
+-----+  
| Length (1 octet) |  
+-----+  
| Prefix (variable) |  
+-----+
```


3. Modifications to BGP Behavior

3.1. Changes to Path Selection Rules

The path selection rules for BGP ([section 9.1.2.2](#) of BGP4 [3]) are changed as follows:

- o The following rule is placed before step (a): If paths with and without BACKUP_ONLY (or BACKUP_ONLY_SECOND_BEST) are both available, those with BACKUP_ONLY/BACKUP_ONLY_SECOND_BEST are eliminated
- o The following rule is a modification to step (c): Step (c) is first performed INCLUDING paths with SECOND_BEST. If, at the end of the first attempt at step (c), only paths with SECOND_BEST remain, re-run step (c), this time EXCLUDING the paths with SECOND_BEST. After this modified version of step (c), the remaining paths MUST NOT have the SECOND_BEST attribute. In other words, Step (c) MUST remove any SECOND_BEST paths.
- o The remainder of the usual BGP path selection rules are applied as normal
- o If the final path selected has the BACKUP_ONLY/BACKUP_ONLY_SECOND_BEST attribute, the attribute BACKUP_ONLY MUST be set.

The path selection rules for "Second Best" path are as follows:

- o The already-selected "best" path is removed from the set of paths to compare
- o The same rules are applied as for the "best" path
- o The selected path is advertised with the attribute SECOND_BEST applied
- o If the selected path had the BACKUP_ONLY attribute, the attribute BACKUP_ONLY_SECOND_BEST must be set.

The prefix instances for consideration of second-best path are the REMAINDER of non-SECOND_BEST instances, and the SECOND_BEST instance received on the in-RIB from which the best path was selected (if one exists). Only one SECOND_BEST instance received may be considered for the local (and out-RIB) SECOND_BEST path.

3.2. Second Best - Basic Method

Once the capability for doing so has been negotiated between a pair of BGP speakers, each sends the best two paths for each prefix. The path information will include the additional SECOND_BEST attribute on the second best path.

When the current "best" path is withdrawn, the withdrawal MAY be propagated without having to perform a full BGP table path selection. The current "second best" path in the local-RIB is promoted to "best". This is because the alternate candidates have already been evaluated and "second-best" has already been selected.

Whenever an AS consists of a mesh of BGP speakers who have negotiated this capability, the withdrawal will propagate through the entire AS. This will either have no effect, or with a change in "best" without requiring non-local information to choose the new "best" path.

3.3. Second Best - Route Reflector

The "best" and "second best" are reflected. The same mechanism is used for determining both best and second-best per prefix. Updates must be reflected whenever the choice of either or both of the "best" or "second best" change. Withdrawals may be propagated immediately.

3.4. Second Best - Inter-AS Hybrid Method

When a withdrawal of the current best path is received from a peer doing USE_SECOND_BEST_AND_BACKUP, and the rules for sending updates require that an update for this prefix be sent to a peer who does not support USE_SECOND_BEST_AND_BACKUP, the current second-best instance of the prefix is sent to that peer in an Update. The neighbor does not need the withdrawal, since the new path replaces the old path.

When the selection of best path results in the selection of a path with BACKUP_ONLY, the path is sent as the best path. This is the only time where a BACKUP_ONLY path is sent as BEST, without preserving the BACKUP_ONLY attribute.

3.5. Backup Only - Basic Method

The main reason for establishing the BACKUP_ONLY attribute is to permit the global implementation of actual "backup only" announcements. It is not to facilitate change of policies, or to circumvent local policies, instead it is to make possible the implementation of policies where those have been negotiated by two or more parties.

Currently, there are several documented scenarios in the "Wedgies" RFC [2] where the mutually desired policy is either unable to be implemented, or does not deterministically reach the desired state.

Use of the BACKUP_ONLY attribute on announcements sent to a backup provider, permit these problems to be resolved.

The same prefix is announced to both the primary and backup provider. When announced to the primary provider, the BACKUP_ONLY attribute is NOT set. When announced to the backup provider, the BACKUP_ONLY attribute IS set.

The propagation of the BACKUP_ONLY instance will be limited by the availability of multiple paths and the use of SECOND_BEST peerings.

In Figure 10 (of [Appendix C](#)), the BACKUP_ONLY instance will be seen by the backup provider, and be passed with both SECOND_BEST and BACKUP_ONLY to the backup provider's transit provider. The latter will prefer any other instance without BACKUP_ONLY, even if it has applied a LOCAL_PREFERENCE to the received prefix instance. Should the other instance be withdrawn, the BACKUP_ONLY will be selected and subsequently propagated. The withdrawal will also eventually result in an Update with the BACKUP_ONLY attribute but WITHOUT the SECOND_BEST attribute (since the prefix will now only be reachable via the backup provider.)

[3.6.](#) Backup Only - Route Reflector

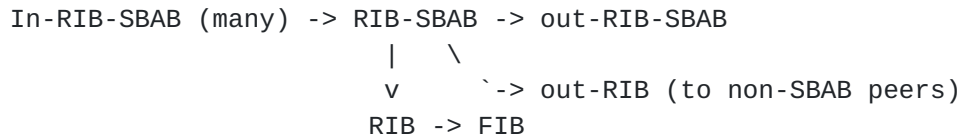
Route Reflectors operate the same as always. The BACKUP_ONLY attribute MUST be preserved during reflection. Thus, if "Second Best" is in operation, then the BACKUP_ONLY attribute of both best and second-best MUST be preserved on both instances. And, if "Second Best" is not in use, then the selected "best" prefix, if it has BACKUP_ONLY set, must be reflected with BACKUP_ONLY as well.

[3.7.](#) IBGP vs EBG

The same rules apply for EBG->EBG, EBG->IBG, IBG->EBG, and IBG->IBG. If a particular peering has had USE_SECOND_BEST_AND_BACKUP negotiated, then any update for a particular prefix that results in new selection of either or both of best and second-best, the new selections (and possible withdrawal of old selections) is sent to the appropriate peers. Additionally, updates which have BACKUP_ONLY MAY be sent.

4. Implementation Guidelines

In order to encourage effective implementation schemes, and to demonstrate some of the benefits of deployment, here are some suggestions for facilitating fast propagation of path changes, which are anticipated as improving behavior. This applies in particular to Path Hunting issues.



```

+-----+-----+-----+-----+
| PREFIX | IN-SBAB | OUT-SBAB | *PATH-info-ptr |
+-----+-----+-----+-----+

```

Figure 4

Where IN-SBAB and OUT-SBAB are 4-bit fields indicating what the SECOND_BEST_AND_BACKUP (SBAB) are attributes (BEST, SECOND_BEST, BACKUP_ONLY, SECOND_BEST_BACKUP_ONLY). IN-SBAB are the attributes received from a peer, and for ONLY those prefixes selected for inclusion into the RIB-SBAB, what the corresponding attributes are.

For example, if all external peers have NOT negotiated SBAB, those prefixes would have SBAB binary values of 1000. Each In-RIB-SBAB would have at most one instance. And for each prefix, at most one In-RIB-SBAB would be selected as best, and have its corresponding OUT-SBAB set to binary value 1000.

This forward-chaining allows for processing of SBAB updates to determine whether withdrawals need to be flooded to peers, and if so, what SBAB attribute to apply to the withdrawals that are flooded. This flooding MAY be performed in parallel to normal BGP table update processing.

For clarity, it should be pointed out that:

- o The process for the step RIB-SBAB to RIB is "select prefixes marked 'best'".
- o The process for the step RIB-SBAB to out-RIB is also "select prefixes marked 'best'".

- o The process for the step RIB-SBAB to out-RIB-SBAB is the same as ordinary RIB to out-RIB, except for preservation of SBAB attributes (if any).

5. Security Considerations

No additional security considerations beyond those already present in BGP are introduced.

6. IANA Considerations

IANA will need to assign new code points for BGP Capabilities for USE_SECOND_BEST_AND_BACKUP. IANA will need to assign new code points for BGP Attribute Types for SECOND_BEST, BACKUP_ONLY and BACKUP_ONLY_SECOND_BEST.

[7.](#) Acknowledgements

The author wishes to acknowledge the helpful guidance of Joe Abley, Tony Li, and Yakhov Rehkter. The author also wishes to acknowledge the insight gained from his Scottish Deerhound, Skylar, winning a Reserve Best-in-Show. (The selection method of "second best" comes from the Reserve system used at the group and best-in-show levels of dog shows).

8. References

8.1. Normative References

- [1] McPherson, D., Gill, V., Walton, D., and A. Retana, "Border Gateway Protocol (BGP) Persistent Route Oscillation Condition", [RFC 3345](#), August 2002.
- [2] Griffin, T. and G. Huston, "BGP Wedgies", [RFC 4264](#), November 2005.
- [3] Rekhter, Y., Li, T., and S. Hares, "A Border Gateway Protocol 4 (BGP-4)", [RFC 4271](#), January 2006.
- [4] Chen, E. and S. Sangli, "Avoid BGP Best Path Transitions from One External to Another", [RFC 5004](#), September 2007.

8.2. Informative References

- [5] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

[Appendix A](#). Path-Hunting Examples

(These will be included in a subsequent version of this ID.)

Appendix B. Persistent Oscillation Examples

Consider the example in Figure 5 where

- o R1, R2, R3, R4, and R5 belong to one AS.
- o R1 is a route reflector with R2 and R3 as its clients.
- o R4 is a route reflector with R5 as its client.
- o The IGP metrics are as listed.
- o External paths (a), (b), and (c) are as described in Figure 6.

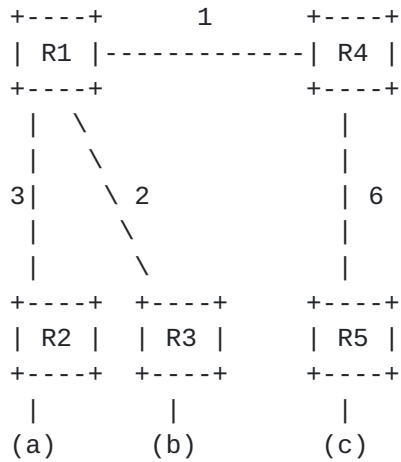


Figure 5

Path	AS_PATH	MED
a	1 3	10
b	2 3	1
c	2 3	0

Figure 6

With the addition of "second best", we have:

R1 has the following:

Path	AS_PATH	MED	IGP-metric
a	1 3	10	3 (received:best) (best)
b	2 3	1	2 (received:best)
c	2 3	0	7 (received:best) (second_best - not sent)

R4 has the following:

Path	AS_PATH	MED	IGP-metric
a	1 3	10	4 (received:best) (best - not sent)
c	2 3	0	6 (received: best) (second_best)

This results in R1 having:

Path	AS_PATH	MED	IGP-metric
a	1 3	10	3 (received:best) (best)
b	2 3	1	2 (received:best)
c	2 3	0	7 (received:second_best) (second_best - not sent)

By including the second_best in the best path calculation, the persistent oscillation problem is resolved.

Appendix C. BGP Wedgie Examples

The following examples from [RFC 4264 \[2\]](#) show the effects of the proposed changes, in resolving "wedgie" issues.

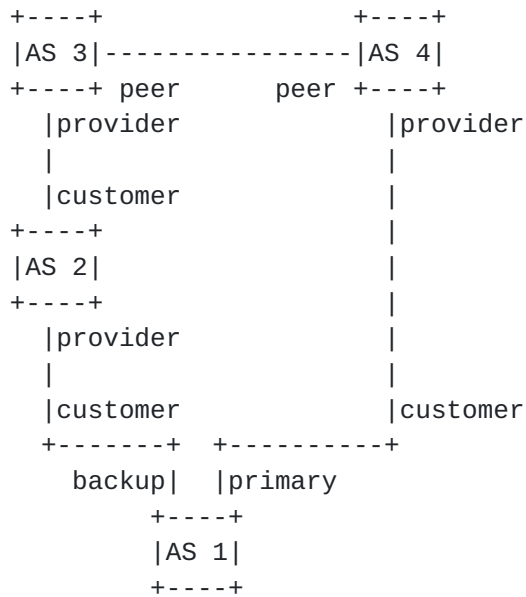


Figure 10

In Figure 10 above, the announcement via the backup link is sent with BACKUP_ONLY.

- o AS 4 sends the "best" (the direct link to AS 1) to AS 3.
- o AS 2 sends its "best", which is the BACKUP_ONLY path from AS 1, to AS 3, also with BACKUP_ONLY (since it is a transitive attribute).
- o AS 3 and AS 4 exchange their respective "best" paths.
- o AS 3 prefers the path "4 1" over "2 1" because "2 1" is BACKUP_ONLY.
- o AS 3 sends a revised BACKUP_ONLY update to AS 4 as SECOND_BEST.
- o AS 3 sends the new "best" to AS 2.
- o AS 2 sends a revised BACKUP_ONLY update to AS 3 as SECOND_BEST.

This state will be reached regardless of sequence of disconnects and reconnects.

Link failures will also result in propagation of withdrawls of "best"

and the SECOND_BEST promotions will result in immediate correct behavior.

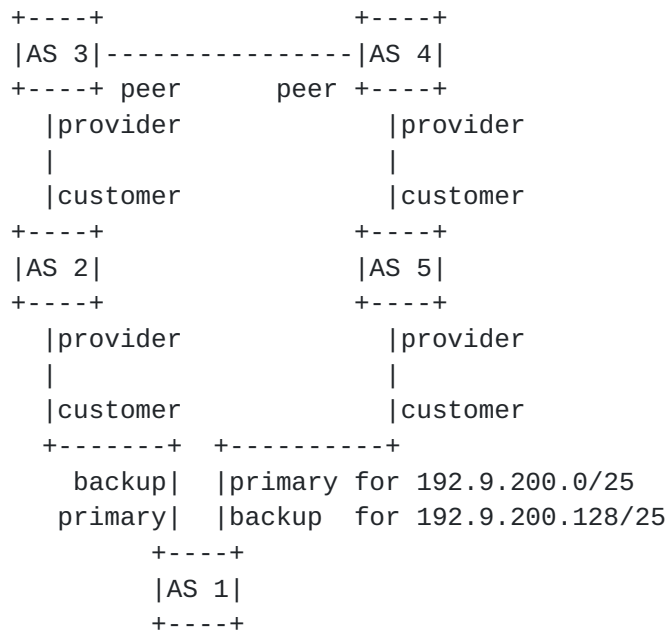


Figure 11

In Figure 11 above, the announcements via the backup links will work the same as in Example 1.

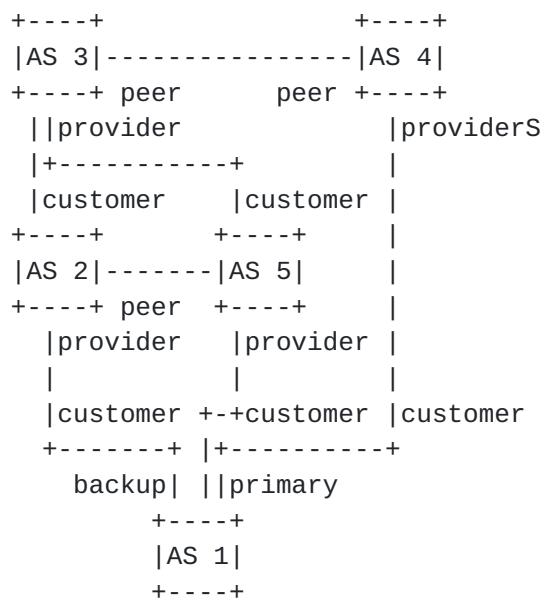


Figure 12

In Figure 12 above, the announcements via both backup links will

result in:

- o AS 2 selecting its best path via "3 4 1" (the only path it hears from AS 3)
- o AS 2 hearing two paths from AS 5:
 - * its "second best" path "5 3 4 1"
 - * another path marked SECOND_BEST and BACKUP_ONLY with path "5 1"
- o AS 2 hearing a BACKUP_ONLY directly from AS 1

Any announcement that AS 3 hears from AS 2 or AS 5 will always be marked BACKUP_ONLY. Thus, any combination of break/restore on any links in any order, will always result in the desired state being reached.

Author's Address

Brian Dickson
Afilias Canada, Inc
4141 Yonge St,
Suite 204
North York, ON M2P 2A8
Canada

Email: brian.peter.dickson@gmail.com

URI: www.afilias.info

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.

Acknowledgment

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).

