

GROW
Internet-Draft
Expires: September 27, 2004

T. Griffin
University of Cambridge
G. Huston
APNIC
March 29, 2004

BGP Wedgies
draft-ietf-grow-bgp-wedgies-01.txt

Status of this Memo

This document is an Internet-Draft and is subject to all provisions of [section 3 of RFC 3667](#). 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 become aware will be disclosed, in accordance with [RFC 3668](#).

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/lid-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 September 27, 2004.

Copyright Notice

Copyright (C) The Internet Society (2004).

Abstract

It has commonly been assumed that the Border Gateway Protocol (BGP) is a tool for distributing reachability information in a manner that creates forwarding paths in a deterministic manner. In this memo we will describe a class of BGP configurations for which there is more than one potential outcome, and where forwarding states other than the intended state are equally stable, and that the stable state

Internet-Draft

BGP Wedgies

March 2004

where BGP converges may be selected by BGP in a non-deterministic manner. These stable, but unintended, BGP states are termed here "BGP Wedgies".

1. Introduction

It has commonly been assumed that the Border Gateway Protocol (BGP) [[RFC1771](#)] is a tool for distributing reachability information in a manner that creates forwarding paths in a deterministic manner. This is a 'problem statement' memo that describes a class of BGP configurations for which there is more than one stable forwarding state. In this class of configurations forwarding states other than the intended state are equally stable, and the stable state where BGP converges may be selected by BGP in a non-deterministic manner.

These stable, but unintended, BGP states are termed here "BGP Wedgies".

2. Describing BGP Routing Policy

BGP routing policies generally reflect each network administrator's objective to optimize their position with respect to their network's cost, performance and reliability.

With respect to cost optimization, the local network's default routing policy often reflects a local preference to prefer routes learned from a customer to routes learned from some form of peering exchange. In the same vein the local network is often configured to prefer routes learned from a peer or a customer over those learned from a directly connected upstream transit provider. These preferences may be expressed via a local preference configuration setting, where the local preference overrides the AS path length metric of the base BGP operation.

In terms of engineering reliability in the inter-domain routing environment it is commonly the case that a service provider may enter into arrangements with two or more upstream transit providers, passing routes to both providers, and receiving traffic from both sources. If the path to one upstream fails the traffic will switch to other links, and once the path is recovered, the traffic should switch back.

In such situations of multiple upstream providers it is also commonplace to place a relative preference on the providers, so that one connection is regarded as a preferred, or "primary" connection, and other connections are regarded as less preferred, or "backup" connections. The intent is typically that the backup connections will be used for traffic only for the duration of a failure in the

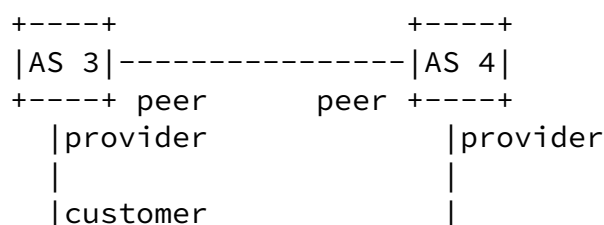
primary connection.

It is possible to express this primary / backup policy using local AS path prepending, where the AS path is artificially lengthened towards the backup providers, using additional instances of the local AS. This is not a deterministic selection algorithm, as the selected primary provider may in turn be using AS path prepending to its backup upstream provider, and in certain cases the path through the backup provider may still be selected as the shortest AS path length.

An alternative approach to routing policy specification uses BGP communities [[RFC1997](#)]. In this case the provider publishes a set of community values that allows the client to select the provider's local preference setting. The client can use a community to mark a route as "backup only" towards the backup provider, and "primary preferred" to the primary provider, assuming both providers support community values with such semantics. In this case the local preference overrides the AS path length metric, so that if the route is marked "backup only", the route will be selected only when there is no other source of the route.

3. BGP Wedgies

The richness of local policy expression through the use of communities, when coupled with the behavior of a distance vector protocol like BGP leads to the observation that certain configurations have more than one "solution", or more than one stable BGP state. An example of such a situation is indicated in Figure 1.



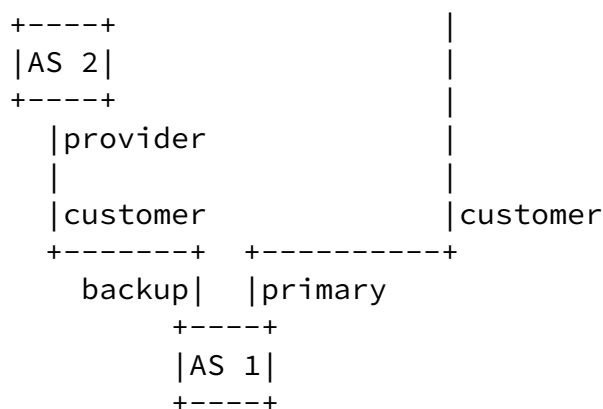


Figure 1

In this case AS1 has marked its advertisement of prefixes to AS2 as "backup only", and its advertisement of prefixes to AS4 as "primary". AS3 will hear AS4's advertisement across the peering link, and pick of AS1's prefixes with the path "AS4, AS1". AS3 will advertise this to AS2. AS2 will hear two paths to AS1, the first is by the direct connection to AS1, and the second is via the path "AS3, AS4, AS1". AS2 will prefer the longer path as the directly connected routes are marked "backup only", and AS2's local preference decision will prefer the AS3 advertisement over the AS1 advertisement.

This is the intended outcome of AS1's policy settings, where no traffic passes from AS2 to AS1, and AS2, reaches AS1 via a path that transits AS3 and AS4.

This intended outcome is achieved as long as AS1 announces its routes on the primary path, to AS4, before announcing its backup routes to AS2.

If the AS1 - AS4 path is broken, causing aBGP sesssion failure between AS1 and AS4, then AS4 will withdraw its advertisement of AS1's routes to AS3, who, in turn will send a withdrawal to AS2. As2, will then select the backup path to AS1. AS2 will advertise this path to AS3, and AS3 will advertise this path to AS4. Again, this is part of the intended operation of the primary / backup policy setting.

When connectivity between AS4 and AS1 is restored the BGP state will not revert to the original state. AS4 will learn the primary path to

AS1, and readvertise this to AS3 using the path "AS4, AS1". AS3, using a default preference of preferring customer-advertised routes over peer routes will continue to prefer the "AS2, AS1" path. AS3 will not pass any updates to AS2. After the restoration of the circuit traffic from AS3 to AS1 and from AS2 to AS1 will be presented to AS1 via the backup path, even through the primary path via AS4 is in service.

The intended forwarding state can only be restored by AS1 deliberately bringing down its eBGP session with AS2, even though it is carrying traffic. This will cause the BGP state to revert to the intended configuration.

It is often the case that an AS will attempt to balance incoming traffic across multiple providers, again using the primary / backup mechanism. For some prefixes one link is configured as the primary link, and the others as the backup link, while for other prefixes another link is selected as the primary link. An example is shown in Figure 2.

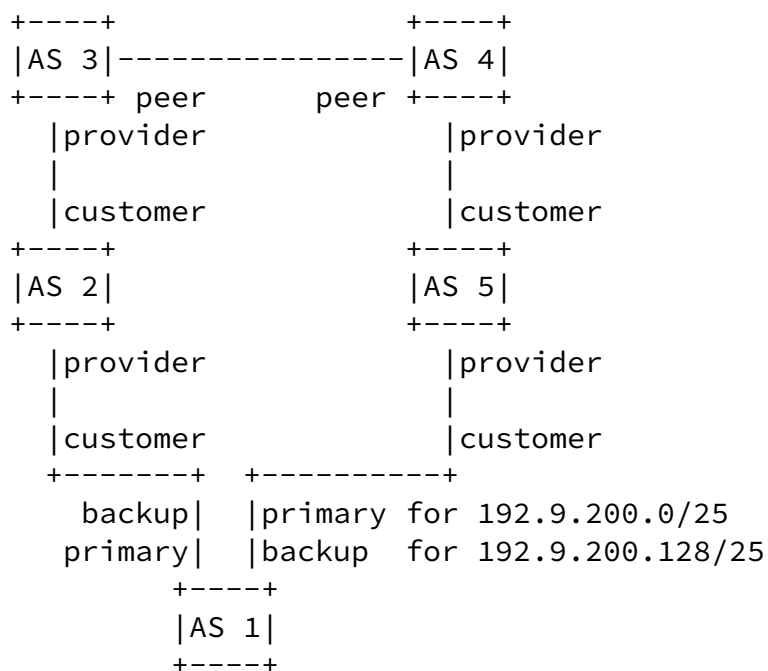


Figure 2

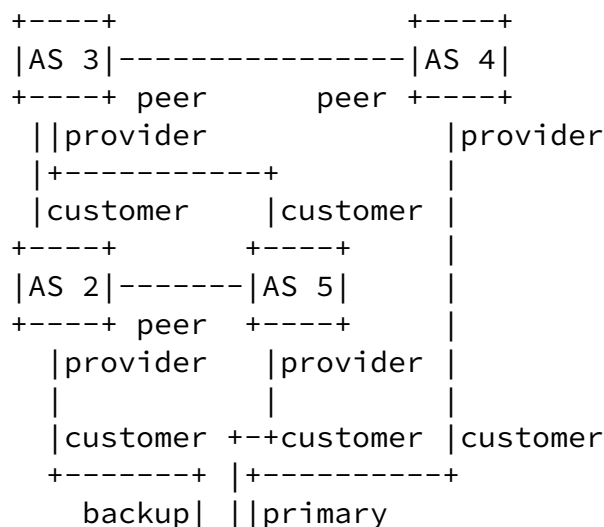
The intended configuration has all incoming traffic for addresses in the range 192.9.200.0/25 via the link from AS5, and all incoming traffic for addresses in the range 192.9.200.128/25 from AS2.

In this case if the link between AS3 and AS4 is reset, AS3 will learn both routes from AS2, and AS4 will learn both routes from AS5. As these customer routes are preferred over peer routes, when the link between AS3 and AS4 is restored, neither AS will alter its routing behavior with respect to AS1's routes. This situation is now wedged, in that there is no eBGP peering that can be reset that will flip BGP back to the intended state. This is an instance of a BGP Wedgie.

The restoration path here is that AS1 has to withdraw the backup advertisements on both paths and operate for an interval without backup, and then readvertise the backup prefix advertisements. The length of the interval cannot be readily determined in advance, as it has to be sufficiently long so as to allow AS2 and AS5 to learn of an alternate path to AS1. At this stage the backup routes can be readvertised.

4. Multi-Party BGP Wedgies

This situation can be more complex when three or more parties provide upstream transit services to an AS. An example is indicated in Figure 3.



```
+-----+
|AS 1|
+-----+
```

Figure 3

In this example the intended state is that AS2 and AS5 are both backup providers, and AS4 is the primary provider. When the link between AS1 and AS4 breaks and is subsequently restored, AS3 will continue to direct traffic to AS1 via AS2 or AS5. In this case a single reset of the link between AS2 and AS1 will not restore the original intended BGP state, as the BGP-selected best route to AS1 will switch to AS5, and AS2 and AS3 will learn a path to AS1 via AS5.

What AS1 is observing is incoming traffic on the backup link from AS2. Resetting this connection will not restore traffic back to the primary path, but instead will switch incoming traffic over to AS5. The action required to correct the situation is to simultaneously reset both the link to AS2, and also the link to AS5. This is not necessarily an intuitive solution, as at any point on time only one of these links will be carrying backup traffic, yet both BGP sessions need to be brought down at the same time in order to commence restoration of the intended primary and backup state.

5. BGP and Determinism

BGP does not behave deterministically in all cases, and, as a consequence, there is intended and unintended non-determinism in BGP. For example, the default final tie break in some implementations of BGP is to prefer the longest-lived route. To achieve determinism in this last step it would be necessary to use a comparison operator that has a predictable outcome, such as a comparison of router identifiers. This class of non-deterministic behavior is termed here

"intended" non-determinism, in that the policy interactions are, to some extent, predictable by network administrators.

BGP is also able to generate outcomes that can be described as "unintended non-determinism" that can result from unexpected policy interactions. These outcomes do not represent misconfiguration in the standard sense, since all policies may look completely rational

locally, but their interaction across multiple routing entities can cause unintended outcomes, and BGP may reach a state that includes such unintended outcomes in a non-deterministic manner.

Unintended non-determinism in BGP would not be as critical an issue if all stable routings were guaranteed to be consistent with the policy writer's intent. However, this is not always the case. The above examples indicate that the operation of BGP allows multiple stable states to exist from a single configuration state, where some of these states are not consistent with the policy writer's intent. These particular examples can be described as a form of "route pinning", where the route is pinned to a non-preferred path.

The challenge for the network administrator is to ensure that an intended state is maintained. Under certain circumstances this can only be achieved by deliberate service disruption, involving the withdrawal of routes being used to forward traffic, and re-advertising routes in a certain sequence in order to induce an intended BGP state. However, the knowledge that is required by any single network operator administrator in order to understand the reason why BGP has stabilized to an unintended state requires BGP policy configuration knowledge of remote networks. In effect there is insufficient local information for any single network administrator to correctly identify the root cause of the unintended BGP state, nor is there sufficient information to allow any single network administrator to undertake a sequence of steps to rectify the situation back to the intended routing state.

It is reasonable to anticipate that as the density of interconnection increases, and also that the capability for policy-based preference setting of learned and re-advertised routes will become more expressive. It is therefore reasonable to anticipate that the incidence of unintended BGP states will increase, and the ability to understand the necessary sequence of route withdrawals and re-advertisements will become more challenging to determine in advance.

Whether this could lead to BGP routing system reaching a point where each network consistently cannot direct traffic in a deterministic manner is at this stage a matter of speculation. BGP Wedgies are an illustration that a sufficiently complex interconnection topology,

coupled with a sufficiently expressive set of policy constructs, can lead to a number of stable BGP states, rather than a single intended state. As the topology complexity increases it is not possible to deterministically predict which state the BGP routing system may converge to. Paradoxically, the demands of inter-domain traffic engineering appear to require both greater levels of expressive capability in policy-based routing directives, operating across denser interconnectivity topologies in a deterministic manner. This may not be a sustainable outcome in BGP-based routing systems.

6. Security Considerations

BGP is a relaying protocol, where route information is received, processed and forwarded. BGP contains no specific mechanisms to prevent the unauthorized modification of the information by a forwarding agent, allowing routing information to be modified, deleted or false information to be inserted without the knowledge of the originator of the routing information or any of the recipients.

The memo proposes no modifications to the BGP protocol, nor does it propose any changes to the manner of deployment of BGP, and therefore introduces no new factors in terms of the security and integrity of inter-domain routing.

The memo illustrates that in attempting to create policy-based outcomes relating to path selection for incoming traffic it is possible to generate BGP configurations where there are multiple stable outcomes, rather than a single outcome. Furthermore, of these instances of multiple outcomes, there are cases where the BGP selection of a particular outcome is not a deterministic selection.

7. References

7.1 Normative References

[RFC1771] Rekhter, Y. and T. Li, "A Border Gateway Protocol 4 (BGP-4)", [RFC 1771](#), March 1995.

7.2 Informative References

[RFC1997] Chandrasekeran, R., Traina, P. and T. Li, "BGP Communities Attribute", [RFC 1997](#), August 1996.

Internet-Draft

BGP Wedgies

March 2004

Authors' Addresses

Tim Griffin
University of Cambridge

EMail: Timothy.Griffin@cl.cam.ac.uk

Geoff Huston
Asia Pacific Network Information Centre

EMail: gih@apnic.net

Internet-Draft

BGP Wedgies

March 2004

Intellectual Property Statement

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.

Disclaimer of Validity

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 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.

Copyright Statement

Copyright (C) The Internet Society (2004). 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.

Acknowledgment

Funding for the RFC Editor function is currently provided by the Internet Society.