GROW Working Group Internet-Draft Intended status: Informational Expires: November 30, 2012 R. Raszuk NTT MCL J. Heitz Ericsson A. Lo Arista L. Zhang UCLA X. Xu Huawei May 29, 2012

Simple Virtual Aggregation (S-VA) draft-ietf-grow-simple-va-07.txt

Abstract

The continued growth in the Default Free Routing Table (DFRT) stresses the global routing system in a number of ways. One of the most costly stresses is FIB size: ISPs often must upgrade router hardware simply because the FIB has run out of space, and router vendors must design routers that have adequate FIB.

FIB suppression is an approach to relieving stress on the FIB by NOT loading selected RIB entries into the FIB. Simple Virtual Aggregation (S-VA) is a simple form of Virtual Aggregation (VA) that allows any and all edge routers to shrink their RIB and FIB requirements substantially and therefore increase their useful lifetime.

S-VA does not increase FIB requirements for core routers. S-VA is extremely easy to configure considerably more so than the various tricks done today to extend the life of edge routers. S-VA can be deployed autonomously by an ISP (cooperation between ISPs is not required), and can co-exist with legacy routers in the ISP.

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

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 <u>http://datatracker.ietf.org/drafts/current/</u>.

Internet-Drafts are draft documents valid for a maximum of six months

Raszuk, et al.

Expires November 30, 2012

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 November 30, 2012.

Copyright Notice

Copyright (c) 2012 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>http://trustee.ietf.org/license-info</u>) 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

$\underline{1}$. Introduction	•	•	•		•	•	•	•		•	•		•	<u>4</u>
<u>1.1</u> . Scope of this Document														<u>5</u>
<u>1.2</u> . Requirements notation														<u>5</u>
<u>1.3</u> . Terminology														<u>5</u>
$\underline{2}$. Operation of S-VA														<u>6</u>
<u>3</u> . Deployment considerations														7
4. IANA Considerations														<u>9</u>
5. Security Considerations .														<u>9</u>
<u>6</u> . Acknowledgements														<u>9</u>
<u>7</u> . References														<u>10</u>
<u>7.1</u> . Normative References .														<u>10</u>
7.2. Informative References														<u>10</u>
Authors' Addresses														<u>10</u>

<u>1</u>. Introduction

ISPs today manage constant DFRT growth in a number of ways. One way, of course, is for ISPs to upgrade their router hardware before DFRT growth outstrips the size of the FIB. This is too expensive for many ISPs. They would prefer to extend the lifetime of routers whose FIBs can no longer hold the full DFRT.

A common approach taken by lower-tier ISPs is to default route to their providers. Routes to customers and peer ISPs are maintained, but everything else defaults to the provider. This approach has several disadvantages. First, packets to Internet destinations may take longer-than-necessary AS paths.

This problem can be mitigated through careful configuration of partial defaults, but this can require substantial configuration overhead. A second problem with defaulting to providers is that the ISP is no longer able to provide the full DFRT to its customers. Finally, provider defaults prevents the ISP from being able to detect martian packets. As a result, the ISP transmits packets that could otherwise have been dropped over its expensive provider links.

An alternative is for the ISP to maintain full routes in its core routers, but to filter routes from edge routers that do not require a full DFRT. These edge routers can then default route to the core or exit routers. This is often possible with edge routers that interface to customer networks. The problem with this approach is that it cannot be used for all edge routers. For instance, it cannot be used for routers that connect to transits. It should also not be used for routers that connect to customers which wish to receive the full DFRT.

This draft describes a very simple technique, called Simple Virtual Aggregation (S-VA), that allows any and all edge routers to have substantially reduced FIB requirements even while still advertising and receiving the full DFRT over BGP. The basic idea is as follows. Core routers in the ISP maintain the full DFRT in the FIB and RIB. Edge routers maintain the full DFRT in the BGP protocol RIB, but suppress certain routes from being installed in RIB and FIB tables. Edge routers install a default route to core routers, to ABRs which are installed on the POP to core boundary or to the ASBR routers.

S-VA requires no changes to BGP and no changes to any choice of forwarding mechanisms in routers. Configuration is extremely simple: S-VA must be enabled on the edge router which needs to save its RIB and FIB space. In the same time operator must inject into his intradomain routing a new prefix further called virtual aggregate (VAprefix) which will be used as the aggregate forwarding reference by

the edge routers performing S-VA. Everything else is automatic. ISPs can deploy FIB suppression autonomously and with no coordination with neighbor ASes.

<u>1.1</u>. Scope of this Document

The scope of this document is limited to Intra-domain S-VA operation. In other words, the case where a single ISP autonomously operates S-VA internally without any coordination with neighboring ISPs.

Note that this document assumes that the S-VA "domain" (i.e. the unit of autonomy) is the AS (that is, different ASes run S-VA independently and without coordination). For the remainder of this document, the terms ISP, AS, and domain are used interchangeably.

This document applies equally to IPv4 and IPv6 both unicast and multicast address families.

S-VA may operate with a mix of upgraded routers and legacy routers. There are no topological restrictions placed on the mix of routers. S-VA functionality is local to the router on which it is enabled and routing correctness is guaranteed.

Note that S-VA is a greatly simplified variant of "full VA" [<u>I-D.ietf-grow-va</u>]. With full VA, all routers (core or otherwise) can have reduced FIBs. However, full VA requires substantial new configuration and operational complexity compared to S-VA. Full VA also requires the use of MPLS LSPs between all routers. Note that S-VA was formerly specified in [<u>I-D.ietf-grow-va</u>]. It has been moved to this separate draft to simplify its understanding.

<u>1.2</u>. Requirements notation

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 [RFC2119].

<u>1.3</u>. Terminology

- RIB/FIB-Installing Router (FIR): An router that does not suppress any routes, and advertises itself as a default route for 0/0. Typically a core router, POP to core boundary router or an ASBR would be configured as an FIR.
- RIB/FIB-Suppressing Router (FSR): An S-VA router that installs a route to 0/0, and may suppress other routes. Typically an edge router would be configured as an FSR.

- Install and Suppress: The terms "install" and "suppress" are used to describe whether a protocol local RIB entry has been loaded or not loaded into the global RIB and FIB. In other words, the phrase "install a route" means "install a route into the global RIB and FIB", and the phrase "suppress a route" means "do not install a route from BGP into the global RIB and FIB".
- Legacy Router: A router that does not run S-VA, and has no knowledge of S-VA.
- Global Routing Information Base (RIB): The term global RIB is used to indicate the router's main routing information base. That RIB is normally used to populate FIB tables of the router. It needs to be highlighted that unless FIB compression is used global RIB and FIB tables are in sync.
- Local/Protocol Routing Information Base (loc-RIB): The term local RIB is used to indicate the protocol's table where product of SPF or BGP best path selection is kept before being installed in global RIB. For example, in some protocol implementations BGP loc-RIB can be further divided into Adj-RIBs-In, the Loc-RIB, and the Adj-RIBs-Out.

2. Operation of S-VA

There are three types of routers in S-VA, FIB-Installing routers (FIR), FIB-Suppressing routers (FSR), and optionally legacy routers. While any router can be an FIR or an FSR (there are no topology constraints), the most simple form of deployment is for AS border or POP border routers to be configured as FIRs, and for customer facing edge routers respectively in the AS or in the POP to be configured as FSRs.

FIRs must originate a default BGP route to NLRI 0/0 [<u>RFC4271</u>]. The ORIGIN is set to INCOMPLETE (value 2) and the BGP NEXT_HOP is set to match the other BGP routes which are also advertised by said FIR. The ATOMIC_AGGREGATE and AGGREGATOR attributes are not included. The FIR MUST attach a NO_EXPORT Community Attribute [<u>RFC1997</u>] to the default route.

FIRs should not FIB-suppress any routes. They may, however, still use some form of local FIB compression algorithm if deemed necessary.

FSRs must detect the VA prefix 0/0 and install it both in loc-RIB, RIB and FIB. Following that FSR may suppress any more specific routes which carry the same next hop as the VA prefix. To guarantee semantical correctness FSR by default should also be able to detect installation of not matching next hop route and reinstall all the more specifics which were previously eligible for suppression to maintain semantical forwarding correctness.

Generally, any more specific route which carries the same next hop as the VA-prefix 0/0 is eligible for suppression. However, provided that there was at least one less specific prefix (e.g., 1.0.0.0/8) and the next-hop of such prefix was different from that of the VA 0/0, those more specific prefixes (e.g., 1.1.1.0/24) which are otherwise subject to suppression would not be eligible for suppression anymore.

Similarly when IBGP multipath is enabled and when multiple VA prefixes are detected which are multipath candidates under given network condition only those more specific prefixes are subject to suppression which have the identical set of next hops as multipath set of VA prefixes.

We illustrate the expected behavior on the figure below. This figure shows an autonomous system with a FIR FIR1 and an FSR FSR1. FSR1 is an ASBR and is connected to two remote ASBRs, EP1 and EP2.

+		-+
Autonomous System		++
1		EP1
	/	-+
To\ ++	++ /	++
Other \ FIR1	FSR1 /	
Routers /		
+	++ \	++
	\	-+ EP2
		++
+		-+

Suppose that FSR1 has been enabled to perform S-VA. Originally it receives all routes from FIR1 (doing next hop self) as well as directly connected EBGP peers EP1 and EP2. FIR1 now will advertise a VA prefix 0/0 with next hop set to himself. That will trigger detection of such prefix on FSR1 and suppression all routes which have the same next hop as VA prefix and which otherwise would be installed in RIB and FIB. However it needs to be observed that FSR1 will not suppress any EBGP routes received from his peers EP1 and EP2 due to next hop being different from the one assigned to VA-prefix.

<u>3</u>. Deployment considerations

The simplest deployment model of S-VA is its use within the POP. In such model the POP to core boundary routers (usually RRs in the data path) would act as FIRs and would inject VA-prefix 0/0 to all of its clients within the POP. In such model of operation an observation

Internet-Draft

can be made that such ABRs do have full routing knowledge and client to ABR distance is negligible as compared with client to intra-domain exit distance.

Therefore under the above intra POP S-VA deployment model clients can be configured that even in the event of lack of ABR to ABR advertisement symmetry there is still no need to monitor if more specific unsuppressed route would cover suppressed one. Thus in this particular deployment model there is no need to detect and reinstall the previously suppressed ones.

Another deployment consideration should be given to networks which may utilize route reflection. In the event of enabling IBGP multipath a special care must be taken that both outbound prefixes as well as VA-prefixes would pass via said route reflectors to their clients.

In order to address the above aspects the following solutions could be considered:

- Use of intra-POP S-VA
- Full mesh Small or medium side networks where S-VA can be deployed are normally fully meshed and do not use route reflection. It also needs to pointed out that some large networks are also fully meshed today.
- Use of add-paths Use of add-paths new BGP encoding will allow to distribute more then one overall best path from RR to each client.
- Alternate advertisement of VA-prefix S-VA prefix does not need to be advertised in BGP. The BGP suppression will happen as long as we configure the S-VA with next hop(s) and implementation verifies that such VA-prefix is installed in the RIB and FIB.

In some deployment scenarios BGP routes could be used to resolve other BGP routes - commonly process called double or multi-level BGP recursion. If such recursion involves specific route resolution policy a special care must be taken to either automatically or manually exclude such routes matching given policy from suppression.

Route resolution over default route is a special case. Most network operating systems can be configured by the operator to enable route resolution over default route(s). In simple-va all default routes are intra-domain routes and their objective it to shift full lookup from edge router to more powerful pop to core boundary router or exit ASBR. In those cases simple-va should be configured in concert with global configuration regarding resolution via default route. In the event of actually using default for next hop resolution the worse case scenario is that the packets may be forwarded one more hop then dropped if more specific destination route is not found there.

Operators are advised to keep effect in mind when choosing a policy for use of default route for next hop resolution.

Selected BGP routes in the RIB may be redistributed to other protocols. If they no longer exist in the RIB, they will not be redistributed. This is especially important when the conditional redistribution is taking place based on the length of the prefix, community value etc .. In those cases where redistribution policy is in place simple-va code should refrain from suppressing prefixes matching such policy.

In the case where operator injects a default at the pop to core boundary into the pop or alternatively when intra-domain default route is injected into autonomous system by set of ASBRs peering with their upstreams a special care needs to be take to make sure that any aggregate subnet is advertised only from the BGP speakers which inject the default route and therefor attract traffic to non existing destinations. This will allow to completely mitigate potential forwarding issue while not specific to simple-va, but applicable to the general use of default routes.

<u>4</u>. IANA Considerations

There are no IANA considerations.

5. Security Considerations

The authors are not aware of any new security considerations due to S-VA.

<u>6</u>. Acknowledgements

The concept for Virtual Aggregation comes from Paul Francis. In this document authors only simplified some aspects of its behavior to allow simpler adoption by some operators.

Authors would like to thank Clarence Filsfils and Nick Hilliard for their valuable input.

7. References

S-VA

7.1. Normative References

- [RFC1997] Chandrasekeran, R., Traina, P., and T. Li, "BGP Communities Attribute", <u>RFC 1997</u>, August 1996.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, March 1997.
- [RFC4271] Rekhter, Y., Li, T., and S. Hares, "A Border Gateway Protocol 4 (BGP-4)", <u>RFC 4271</u>, January 2006.

7.2. Informative References

[I-D.ietf-grow-va]

Francis, P., Xu, X., Ballani, H., Jen, D., Raszuk, R., and L. Zhang, "FIB Suppression with Virtual Aggregation", <u>draft-ietf-grow-va-06</u> (work in progress), December 2011.

Authors' Addresses

Robert Raszuk NTT MCL 101 S Ellsworth Avenue Suite 350 San Mateo, CA 94401 US

Email: robert@raszuk.net

Jakob Heitz Ericsson 300 Holger Way San Jose, CA 95135 USA

Phone: Email: jakob.heitz@ericsson.com

Internet-Draft

Alton Lo Arista Networks 5470 Great America Parkway Santa Clara, CA 95054 USA

Phone: Email: altonlo@aristanetworks.com

Lixia Zhang UCLA 3713 Boelter Hall Los Angeles, CA 90095 US

Phone: Email: lixia@cs.ucla.edu

Xiaohu Xu Huawei Technologies No.3 Xinxi Rd., Shang-Di Information Industry Base, Hai-Dian District Beijing, Beijing 100085 P.R.China

Phone: +86 10 82836073 Email: xuxh@huawei.com