

Internet Engineering Task Force
Internet-Draft
Intended status: Best Current Practice
Expires: July 18, 2014

J. Durand
CISCO Systems, Inc.
I. Pepelnjak
NIL
G. Doering
SpaceNet
January 14, 2014

BGP operations and security
draft-ietf-opsec-bgp-security-02.txt

Abstract

BGP (Border Gateway Protocol) is the protocol almost exclusively used in the Internet to exchange routing information between network domains. Due to this central nature, it is important to understand the security measures that can and should be deployed to prevent accidental or intentional routing disturbances.

This document describes measures to protect the BGP sessions itself (like TTL, MD5, control plane filtering) and to better control the flow of routing information, using prefix filtering and automatization of prefix filters, max-prefix filtering, AS path filtering, route flap dampening and BGP community scrubbing.

Foreword

A placeholder to list general observations about this document.

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](#) [1].

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

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

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."

This Internet-Draft will expire on July 18, 2014.

Copyright Notice

Copyright (c) 2014 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	3
2.	Definitions and Acronyms	3
3.	Protection of the BGP router	4
4.	Protection of BGP sessions	4
4.1.	Protection of TCP sessions used by BGP	4
4.2.	BGP TTL security (GTSM)	5
5.	Prefix filtering	5
5.1.	Definition of prefix filters	5
5.1.1.	Special purpose prefixes	5
5.1.2.	Prefixes not allocated	6
5.1.3.	Prefixes too specific	10
5.1.4.	Filtering prefixes belonging to the local AS and downstreams	10
5.1.5.	IXP LAN prefixes	10
5.1.6.	The default route	12
5.2.	Prefix filtering recommendations in full routing networks	12
5.2.1.	Filters with internet peers	12
5.2.2.	Filters with customers	14
5.2.3.	Filters with upstream providers	15
5.3.	Prefix filtering recommendations for leaf networks . . .	15
5.3.1.	Inbound filtering	15
5.3.2.	Outbound filtering	16
6.	BGP route flap dampening	16
7.	Maximum prefixes on a peering	16
8.	AS-path filtering	17
9.	Next-Hop Filtering	18

10.	BGP community scrubbing	19
11.	Possible future work	19
12.	Change logs	19
12.1.	Diffs between draft-jdurand-bgp-security-01 and draft-jdurand-bgp-security-00	19
12.2.	Diffs between draft-jdurand-bgp-security-02 and draft-jdurand-bgp-security-01	20
12.3.	Diffs between draft-ietf-opsec-bgp-security-00 and draft-jdurand-bgp-security-02	21
12.4.	Diffs between draft-ietf-opsec-bgp-security-01 and draft-ietf-opsec-bgp-security-00	22
12.5.	Diffs between draft-ietf-opsec-bgp-security-02 and draft-ietf-opsec-bgp-security-01	22
13.	Acknowledgements	23
14.	IANA Considerations	23
15.	Security Considerations	24
16.	References	24
16.1.	Normative References	24
16.2.	Informative References	25
	Authors' Addresses	27

[1.](#) Introduction

BGP [\[6\]](#) is the protocol used in the internet to exchange routing information between network domains. This protocol does not directly include mechanisms that control that routes exchanged conform to the various rules defined by the Internet community. This document intends to both summarize common existing rules and help network administrators apply coherent BGP policies.

[2.](#) Definitions and Accronyms

- o ACL: Access Control List
- o IRR: Internet Routing Registry
- o IXP: Internet eXchange Point
- o LIR: Local Internet Registry
- o pMTUd: Path MTU Discovery
- o RIR: regional Internet Registry
- o Tier 1 transit provider: an IP transit provider which can reach any network on the internet without purchasing transit services
- o uRPF: Unicast Reverse Path Forwarding

3. Protection of the BGP router

The BGP router needs to be protected from stray packets. This protection should be achieved by an access control list (ACL) which would discard all packets directed to TCP port 179 on the local device and sourced from an address not known or permitted to become a BGP neighbor. If supported, an ACL specific to the control-plane of the router should be used (receive-ACL, control-plane policing, etc.), to avoid configuration of data-plane filters for packets transiting through the router (and therefore not reaching the control plane). If the hardware can not do that, interface ACLs can be used to block packets to the local router.

Some routers automatically program such an ACL upon BGP configuration. On other devices this ACL should be configured and maintained manually or using scripts.

The filtering of packets destined to the local router is a wider topic than "just for BGP" (if you bring down a router by overloading one of the other protocols from remote, BGP is harmed as well). For a more detailed recommendation, see [RFC6192](#) [22].

4. Protection of BGP sessions

Current issues of TCP-based protocols (therefore including BGP) have been documented in [29]. The following sub-sections recall the major points raised in this RFC and gives best practices for BGP operation.

4.1. Protection of TCP sessions used by BGP

Attacks on TCP sessions used by BGP (ex: sending spoofed TCP RST packets) could bring down the TCP session. Following a successful ARP spoofing attack (or other similar Man-in-the-Middle attack), the attacker might even be able to inject packets into the TCP stream (routing attacks).

TCP sessions used by BGP can be secured with a variety of mechanisms. MD5 protection of TCP session header [14] is the most common one as it was the first mechanism widely implemented on routers. IPsec or TCP Authentication Option (TCP-AO, [10]) offer stronger protection and should now be preferred when available.

The drawback of TCP session protection is additional configuration and management overhead for authentication information (ex: MD5 password) maintenance. Protection of TCP sessions used by BGP is thus recommended when peerings are established over shared networks where spoofing can be done (like IXPs).

You SHOULD block spoofed packets (packets with a source IP address belonging to your IP address space) at all edges of your network, making the protection of TCP sessions used by BGP unnecessary on iBGP or eBGP sessions run over point-to-point links.

[4.2.](#) BGP TTL security (GTSM)

BGP sessions can be made harder to spoof with the Generalized TTL Security Mechanisms (aka TTL security) [[9](#)]. Instead of sending TCP packets with TTL value = 1, the routers send the TCP packets with TTL value = 255 and the receiver checks that the TTL value equals 255. Since it's impossible to send an IP packet with TTL = 255 to a non-directly-connected IP host, BGP TTL security effectively prevents all spoofing attacks coming from third parties not directly connected to the same subnet as the BGP-speaking routers. Network administrators SHOULD implement TTL security on directly connected BGP peerings.

Note: Like MD5 protection, TTL security has to be configured on both ends of a BGP session.

[5.](#) Prefix filtering

The main aspect of securing BGP resides in controlling the prefixes that are received/advertised on the BGP peerings. Prefixes exchanged between BGP peers are controlled with inbound and outbound filters that can match on IP prefixes (prefix filters, [Section 5](#)), AS paths (as-path filters, [Section 8](#)) or any other attributes of a BGP prefix (for example, BGP communities, [Section 10](#)).

[5.1.](#) Definition of prefix filters

This section list the most commonly used prefix filters. Following sections will clarify where these filters should be applied.

[5.1.1.](#) Special purpose prefixes

[5.1.1.1.](#) IPv4 special purpose prefixes

IPv4 registry [[36](#)] maintains the list of IPv4 special purpose prefixes and their routing scope. Reader will refer to this registry in order to configure prefix filters. Only prefixes with value "False" in column "Global" MUST be discarded on Internet BGP peerings.

5.1.1.2. IPv6 special purpose prefixes

IPv6 registry [37] maintains the list of IPv6 special purpose prefixes and their routing scope. Reader will refer to this registry in order to configure prefix filters. Only prefixes with value "False" in column "Global" MUST be discarded on Internet BGP peerings.

At the time of the writing of this document, the list of IPv6 prefixes that MUST not cross network boundaries can be simplified as IANA allocates at the time being prefixes to RIR's only in 2000::/3 prefix [35]. All other prefixes (ULA's, link-local, multicast... are outside of that prefix) and therefore the simplified list becomes:

- o 2001:DB8::/32 and more specifics - documentation [16]
- o Prefixes more specifics than 2002::/16 - 6to4 [3]
- o 3FFE::/16 and more specifics - was initially used for the 6Bone (worldwide IPv6 test network) and returned to IANA
- o All prefixes that are outside 2000::/3 prefix

5.1.2. Prefixes not allocated

IANA allocates prefixes to RIRs which in turn allocate prefixes to LIRs. It is wise not to accept in the routing table prefixes that are not allocated. This could mean allocation made by IANA and/or allocations done by RIRs. This section details the options for building a list of allocated prefixes at every level. It is important to understand that filtering prefixes not allocated requires constant updates as prefixes are continually allocated. Therefore automation of such prefix filters is key for the success of this approach. One SHOULD probably NOT consider solutions described in this section if they are not capable of maintaining updated prefix filters: the damage would probably be worse than the intended security policy.

5.1.2.1. IANA allocated prefix filters

IANA has allocated all the IPv4 available space. Therefore there is no reason why one would keep checking prefixes are in the IANA allocated IPv4 address space [38]. No specific filters need to be put in place by administrators who want to make sure that IPv4 prefixes they receive in BGP updates have been allocated by IANA.

For IPv6, given the size of the address space, it can be seen as wise accepting only prefixes derived from those allocated by IANA.

Administrators can dynamically build this list from the IANA allocated IPv6 space [39]. As IANA keeps allocating prefixes to RIRs, the aforementioned list should be checked regularly against changes and if they occur, prefix filters should be computed and pushed on network devices. The list could also be pulled directly by routers when they implement such mechanisms. As there is delay between the time a RIR receives a new prefix and the moment it starts allocating portions of it to its LIRs, there is no need doing this step quickly and frequently. Based on past experience, authors recommend that the process in place makes sure there is no more than one month between the time the IANA IPv6 allocated prefix list changes and the moment all IPv6 prefix filters are updated.

If process in place (manual or automatic) cannot guarantee that the list is updated regularly then it's better not to configure any filters based on allocated networks. The IPv4 experience has shown that many network operators implemented filters for prefixes not allocated by IANA but did not update them on a regular basis. This created problems for latest allocations and required a extra work for RIRs that had to "de-bogonize" the newly allocated prefixes.

5.1.2.2. RIR allocated prefix filters

A more precise check can be performed as one would like to make sure that prefixes they receive are being originated or transited by autonomous systems entitled to do so. It has been observed in the past that one could easily advertise someone else's prefix (or more specific prefixes) and create black holes or security threats. To partially mitigate this risk, administrators would need to make sure BGP advertisements correspond to information located in the existing registries. At this stage 2 options can be considered (short and long term options). They are described in the following subsections.

5.1.2.3. Prefix filters creation from Internet Routing Registries (IRR)

An Internet Routing Registry (IRR) is a database containing internet routing information, described using Routing Policy Specification Language objects [17]. Network administrators are given privileges to describe routing policies of their own networks in the IRR and information is published, usually publicly. A majority of Regional Internet Registries do also operate an IRR and can control that registered routes conform to prefixes allocated or directly assigned.

It is possible to use the IRR information to build, for a given neighbor autonomous system, a list of prefixes originated or transited which one may accept. This can be done relatively easily using scripts and existing tools capable of retrieving this

information in the registries. This approach is exactly the same for both IPv4 and IPv6.

The macro-algorithm for the script is described as follows. For the peer that is considered, the distant network administrator has provided the autonomous system and may be able to provide an AS-SET object (aka AS-MACRO). An AS-SET is an object which contains AS numbers or other AS-SETs. An operator may create an AS-SET defining all the AS numbers of its customers. A tier 1 transit provider might create an AS-SET describing the AS-SET of connected operators, which in turn describe the AS numbers of their customers. Using recursion, it is possible to retrieve from an AS-SET the complete list of AS numbers that the peer is likely to announce. For each of these AS numbers, it is also easy to check in the corresponding IRR for all associated prefixes. With these two mechanisms a script can build for a given peer the list of allowed prefixes and the AS number from which they should be originated. One could decide not use the origin information and only build monolithic prefix filters from fetched data.

As prefixes, AS numbers and AS-SETs may not all be under the same RIR authority, a difficulty resides choosing for each object the appropriate IRR to poll. Some IRRs have been created and are not restricted to a given region or authoritative RIR. They allow RIRs to publish information contained in their IRR in a common place. They also make it possible for any subscriber (probably under contract) to publish information too. When doing requests inside such an IRR, it is possible to specify the source of information in order to have the most reliable data. One could check a popular IRR containing many sources (such as RADB [\[40\]](#), the Routing Assets Database) and only select as sources some desired RIRs and trusted major ISPs.

As objects in IRRs may frequently vary over time, it is important that prefix filters computed using this mechanism are refreshed regularly. A daily basis could even be considered as some routing changes must be done sometimes in a certain emergency and registries may be updated at the very last moment. It has to be noted that this approach significantly increases the complexity of the router configurations as it can quickly add tens of thousands configuration lines for some important peers.

Last but not least, authors recommend that network administrators publish and maintain their resources properly in IRR database maintained by their RIR, when available.

5.1.2.4. SIDR - Secure Inter Domain Routing

An infrastructure called SIDR (Secure Inter-Domain Routing) [23] has been designed to secure internet advertisements. At the time this document is written, many documents have been published and a framework with a complete set of protocols is proposed so that advertisements can be checked against signed routing objects in RIR routing registries. There are basically two services that SIDR offers:

- o Origin validation [11] seeks at making sure that attributes associated with a routes are correct (the major point being the validation of the AS number originating this route). Origin validation is now operational (Internet registries, protocols, implementations on some routers...) and in theory it can be implemented knowing that the proportion of signed resources is still low at the time this document is written.
- o Path validation provided by BGPsec [42] seeks at making sure that no ones announce fake/wrong BGP paths that would attract traffic for a given destination [43]. BGPsec is still an on-going work item at the time this document is written and therefore cannot be implemented.

Implementing SIDR mechanisms is expected to solve many of BGP routing security problems in the long term but it may take time for deployments to be made and objects to become signed. It also has to be pointed that SIDR infrastructure is complementing (not replacing) the security best practices listed in this document. Authors therefore recommend to implement any SIDR proposed mechanism (example: route origin validation) on top of the other existing mechanisms even if they could sometimes appear targeting the same goal.

If route origin validation is implemented, authors recommend to refer to rules described in [45]. In short, each external route received on a router SHOULD be checked against the RPKI data set:

- o If a corresponding ROA is found and is valid then the prefix SHOULD be accepted.
- o If the ROA is found and is INVALID then the prefix SHOULD be discarded.
- o If an ROA is not found then the prefix SHOULD be accepted but corresponding route SHOULD be given a low preference.

Authors also recommend that network operators sign their routing objects so their routes can be validated by other networks running origin validation.

[5.1.3.](#) Prefixes too specific

Most ISPs will not accept advertisements beyond a certain level of specificity (and in return do not announce prefixes they consider as too specific). That acceptable specificity is decided for each peering between the 2 BGP peers. Some ISP communities have tried to document acceptable specificity. This document does not make any judgement on what the best approach is, it just recalls that there are existing practices on the internet and recommends the reader to refer to what those are. As an example the RIPE community has documented that IPv4 prefixes longer than /24 and IPv6 prefixes longer than /48 are generally not announced/accepted in the internet [[31](#)] [[32](#)].

[5.1.4.](#) Filtering prefixes belonging to the local AS and downstreams

A network SHOULD filter its own prefixes on peerings with all its peers (inbound direction). This prevents local traffic (from a local source to a local destination) from leaking over an external peering in case someone else is announcing the prefix over the Internet. This also protects the infrastructure which may directly suffer in case backbone's prefix is suddenly preferred over the Internet.

To an extent, such filters can also be configured on a network for the prefixes of its downstreams in order to protect them too. Such filters must be defined with caution as they can break existing redundancy mechanisms. For example in case an operator has a multihomed customer, it should keep accepting the customer prefix from its peers and upstreams. This will make it possible for the customer to keep accessing its operator network (and other customers) via the internet in case the BGP peering between the customer and the operator is down.

[5.1.5.](#) IXP LAN prefixes

[5.1.5.1.](#) Network security

When a network is present on an IXP and peers with other IXP members over a common subnet (IXP LAN prefix), it MUST NOT accept more specific prefixes for the IXP LAN prefix from any of its external BGP peers. Accepting these routes may create a black hole for connectivity to the IXP LAN.

If the IXP LAN prefix is accepted as an "exact match", care needs to be taken to avoid other routers in the network sending IXP traffic towards the externally-learned IXP LAN prefix (recursive route lookup pointing into the wrong direction). This can be achieved by preferring IGP routes before eBGP, or by using "BGP next-hop-self" on all routes learned on that IXP.

If the IXP LAN prefix is accepted at all, it **MUST** only be accepted from the ASes that the IXP authorizes to announce it - which will usually be automatically achieved by filtering announcements by IRR DB.

5.1.5.2. pMTUD and the loose uRPF problem

In order to have pMTUD working in the presence of loose uRPF, it is necessary that all the networks that may source traffic that could flow through the IXP (ie. IXP members and their downstreams) have a route for the IXP LAN prefix. This is necessary as "packet too big" ICMP messages sent by IXP members' routers may be sourced using an address of the IXP LAN prefix. In the presence of loose uRPF, this ICMP packet is dropped if there is no route for the IXP LAN prefix or a less specific route covering IXP LAN prefix.

In that case, any IXP member **SHOULD** make sure it has a route for the IXP LAN prefix or a less specific prefix on all its routers and that it announces the IXP LAN prefix or less specific (up to a default route) to its downstreams. The announcements done for this purpose **SHOULD** pass IRR-generated filters described in [Section 5.1.2.3](#) as well as "prefixes too specific" filters described in [Section 5.1.3](#). The easiest way to implement this is that the IXP itself takes care of the origination of its prefix and advertises it to all IXP members through a BGP peering. Most likely the BGP route servers would be used for this. The IXP would most likely send its entire prefix which would be equal or less specific than the IXP LAN prefix.

5.1.5.3. Example

Let's take as an example an IXP in the RIPE region for IPv4. It would be allocated a /22 by RIPE NCC (X.Y.0.0/22 in our example) and use a /23 of this /22 for the IXP LAN (let say X.Y.0.0/23). This IXP LAN prefix is the one used by IXP members to configure eBGP peerings. The IXP could also be allocated an AS number (AS64496 in our example).

Any IXP member **MUST** make sure it filters prefixes more specific than X.Y.0.0/23 from all its eBGP peers. If it received X.Y.0.0/24 or X.Y.1.0/24 this could seriously impact its routing.

The IXP SHOULD originate X.Y.0.0/22 and advertise it to its members through an eBGP peering (most likely from its BGP route servers, configured with AS64496).

The IXP members SHOULD accept the IXP prefix only if it passes the IRR generated filters (see [Section 5.1.2.3](#))

IXP members SHOULD then advertise X.Y.0.0/22 prefix to their downstreams. This announce would pass IRR based filters as it is originated by the IXP.

[5.1.6.](#) The default route

[5.1.6.1.](#) IPv4

The 0.0.0.0/0 prefix is likely not intended to be accepted nor advertised other than in specific customer / provider configurations, general filtering outside of these is RECOMMENDED.

[5.1.6.2.](#) IPv6

The ::/0 prefix is likely not intended to be accepted nor advertised other than in specific customer / provider configurations, general filtering outside of these is RECOMMENDED.

[5.2.](#) Prefix filtering recommendations in full routing networks

For networks that have the full internet BGP table, some policies should be applied on each BGP peer for received and advertised routes. It is recommended that each autonomous system configures rules for advertised and received routes at all its borders as this will protect the network and its peer even in case of misconfiguration. The most commonly used filtering policy is proposed in this section and uses prefix filters defined in previous section [Section 5.1](#).

[5.2.1.](#) Filters with internet peers

[5.2.1.1.](#) Inbound filtering

There are basically 2 options, the loose one where no check will be done against RIR allocations and the strict one where it will be verified that announcements strictly conform to what is declared in routing registries.

5.2.1.1.1. Inbound filtering loose option

In this case, the following prefixes received from a BGP peer will be filtered:

- o Prefixes not globally routable ([Section 5.1.1](#))
- o Prefixes not allocated by IANA (IPv6 only) ([Section 5.1.2.1](#))
- o Routes too specific ([Section 5.1.3](#))
- o Prefixes belonging to the local AS ([Section 5.1.4](#))
- o IXP LAN prefixes ([Section 5.1.5](#))
- o The default route ([Section 5.1.6](#))

5.2.1.1.2. Inbound filtering strict option

In this case, filters are applied to make sure advertisements strictly conform to what is declared in routing registries ([Section 5.1.2.2](#)). Warn is given as registries are not always accurate (prefixes missing, wrong information...) This varies accross the registries and regions of the Internet. Before applying a strict policy the reader SHOULD check the impact on the filter and make sure solution is not worse than the problem.

Also in case of script failure each administrator may decide if all routes are accepted or rejected depending on routing policy. While accepting the routes during that time frame could break the BGP routing security, rejecting them might re-route too much traffic on transit peers, and could cause more harm than what a loose policy would have done.

In addition to this, one could apply the following filters beforehand in case the routing registry used as source of information by the script is not fully trusted:

- o Prefixes not globally routable ([Section 5.1.1](#))
- o Routes too specific ([Section 5.1.3](#))
- o Prefixes belonging to the local AS ([Section 5.1.4](#))
- o IXP LAN prefixes ([Section 5.1.5](#))
- o The default route ([Section 5.1.6](#))

5.2.1.2. Outbound filtering

Configuration should be put in place to make sure that only appropriate prefixes are sent. These can be, for example, prefixes belonging to both the network in question and its downstreams. This can be achieved by using a combination of BGP communities, AS-paths or both. It can also be desirable that following filters are positioned before to avoid unwanted route announcement due to bad configuration:

- o Prefixes not globally routable ([Section 5.1.1](#))
- o Routes too specific ([Section 5.1.3](#))
- o IXP LAN prefixes ([Section 5.1.5](#))
- o The default route ([Section 5.1.6](#))

In case it is possible to list the prefixes to be advertised, then just configuring the list of allowed prefixes and denying the rest is sufficient.

5.2.2. Filters with customers

5.2.2.1. Inbound filtering

The inbound policy with end customers is pretty straightforward: only customers prefixes MUST be accepted, all others MUST be discarded. The list of accepted prefixes can be manually specified, after having verified that they are valid. This validation can be done with the appropriate IP address management authorities.

The same rules apply in case the customer is also a network connecting other customers (for example a tier 1 transit provider connecting service providers). An exception can be envisaged in case it is known that the customer network applies strict inbound/outbound prefix filtering, and the number of prefixes announced by that network is too large to list them in the router configuration. In that case filters as in [Section 5.2.1.1](#) can be applied.

5.2.2.2. Outbound filtering

The outbound policy with customers may vary according to the routes customer wants to receive. In the simplest possible scenario, the customer may only want to receive only the default route, which can be done easily by applying a filter with the default route only.

In case the customer wants to receive the full routing (in case it is multihomed or if wants to have a view of the internet table), the following filters can be simply applied on the BGP peering:

- o Prefixes not globally routable ([Section 5.1.1](#))
- o Routes too specific ([Section 5.1.3](#))
- o The default route ([Section 5.1.6](#))

There can be a difference for the default route that can be announced to the customer in addition to the full BGP table. This can be done simply by removing the filter for the default route. As the default route may not be present in the routing table, one may decide to originate it only for peerings where it has to be advertised.

[5.2.3.](#) Filters with upstream providers

[5.2.3.1.](#) Inbound filtering

In case the full routing table is desired from the upstream, the prefix filtering to apply is the same as the one for peers [Section 5.2.1.1](#) with the exception of the default route. The default route can be desired from an upstream provider in addition to the full BGP table. In case the upstream provider is supposed to announce only the default route, a simple filter will be applied to accept only the default prefix and nothing else.

[5.2.3.2.](#) Outbound filtering

The filters to be applied would most likely not differ much from the ones applied for internet peers ([Section 5.2.1.2](#)). But different policies could be applied in case it is desired that a particular upstream does not provide transit to all the prefixes.

[5.3.](#) Prefix filtering recommendations for leaf networks

[5.3.1.](#) Inbound filtering

The leaf network will position the filters corresponding to the routes it is requesting from its upstream. In case a default route is requested, a simple inbound filter can be applied to accept only the default route ([Section 5.1.6](#)). In case the leaf network is not capable of listing the prefixes because the amount is too large (for example if it requires the full internet routing table) then it should configure filters to avoid receiving bad announcements from its upstream:

- o Prefixes not routable ([Section 5.1.1](#))
- o Routes too specific ([Section 5.1.3](#))
- o Prefixes belonging to local AS ([Section 5.1.4](#))
- o The default route ([Section 5.1.6](#)) depending if the route is requested or not

[5.3.2.](#) Outbound filtering

A leaf network will most likely have a very straightforward policy: it will only announce its local routes. It can also configure the following prefixes filters described in [Section 5.2.1.2](#) to avoid announcing invalid routes to its upstream provider.

[6.](#) BGP route flap dampening

The BGP route flap dampening mechanism makes it possible to give penalties to routes each time they change in the BGP routing table. Initially this mechanism was created to protect the entire internet from multiple events impacting a single network. Studies have shown that implementations of BGP route flap dampening could cause more harm than they solve problems and therefore RIPE community has in the past recommended not using BGP route flap dampening [[30](#)]. Works have then been conducted to propose new route flap dampening thresholds in order to make the solution "usable" [[41](#)] and RIPE has reviewed its recommendations in [[33](#)]. New thresholds have been proposed to make BGP route flap dampening usable. Authors of this document propose to follow RIPE recommendations and only use BGP route flap dampening with adjusted configured thresholds.

[7.](#) Maximum prefixes on a peering

It is recommended to configure a limit on the number of routes to be accepted from a peer. Following rules are generally recommended:

- o From peers, it is recommended to have a limit lower than the number of routes in the internet. This will shut down the BGP peering if the peer suddenly advertises the full table. One can also configure different limits for each peer, according to the number of routes they are supposed to advertise plus some headroom to permit growth.
- o From upstreams which provide full routing, it is recommended to have a limit higher than the number of routes in the internet. A limit is still useful in order to protect the network (and in particular the routers' memory) if too many routes are sent by the

upstream. The limit should be chosen according to the number of routes that can actually be handled by routers.

It is important to regularly review the limits that are configured as the internet can quickly change over time. Some vendors propose mechanisms to have two thresholds: while the higher number specified will shutdown the peering, the first threshold will only trigger a log and can be used to passively adjust limits based on observations made on the network.

8. AS-path filtering

This section is listing rules that apply to BGP AS-paths (for both 16 and 32 bits Autonomous System Numbers):

- o You SHOULD accept from customers only AS(4)-Paths containing ASNs belonging to (or authorized to transit through) the customer. If you can not build and generate filtering expressions to implement this, consider accepting only path lengths relevant to the type of customer you have (as in, if they are a leaf or have customers of their own), try to discourage excessive prepending in such paths.
- o You SHOULD NOT advertise prefixes with non-empty AS-path unless you intend to be transit for these prefixes.
- o You SHOULD NOT advertise prefixes with upstream AS numbers in the AS-path to your peering AS unless you intend to be transit for these prefixes.
- o You SHOULD NOT accept prefixes with private AS numbers in the AS-path except from customers. Exception: an upstream offering some particular service like black-hole origination based on a private AS number. Customers should be informed by their upstream in order to put in place ad-hoc policy to use such services.
- o You SHOULD NOT advertise prefixes with private AS numbers in the AS-path unless you are a customer using BGP without your own AS number. In that case you SHOULD use private AS numbers to advertise your prefixes to your upstream. This private AS number is usually provided by the upstream.
- o You SHOULD NOT accept prefixes when the first AS number in the AS-path is not the one of the peer unless you the peering is done toward a BGP route-server [\[12\]](#) (connection on an IXP) with transparent AS path handling. In that case this verification needs to be de-activated as the first AS number will be the one of an IXP member whereas the peer AS number will be the one of the BGP route-server.

- o You SHOULD NOT override BGP's default behavior accepting your own AS number in the AS-path. In case an exception to this is required, impacts should be studied carefully as this can create severe impact on routing.

AS-path filtering should be further analyzed when ASN renumbering is done. Such operation is common and mechanisms exist to allow smooth ASN migration [44]. The usual migration technique, local to a router, consists in modifying the AS-path so it is presented to a peer as if no renumbering was done. This makes it possible to change ASN of a router without reconfiguring all eBGP peers at the same time (as this operation would require synchronization with all peers attached to that router). During this renumbering operation, rules described above may be adjusted.

9. Next-Hop Filtering

If peering on a shared network, like an IXP, BGP can advertise prefixes with a 3rd-party next-hop, thus directing packets not to the peer announcing the prefix but somewhere else.

This is a desirable property for BGP route-server setups [12], where the route-server will relay routing information, but has neither capacity nor desire to receive the actual data packets. So the BGP route-server will announce prefixes with a next-hop setting pointing to the router that originally announced the prefix to the route-server.

In direct peerings between ISPs, this is undesirable, as one of the peers could trick the other one to send packets into a black hole (unreachable next-hop) or to an unsuspecting 3rd party who would then have to carry the traffic. Especially for black-holing, the root cause of the problem is hard to see without inspecting BGP prefixes at the receiving router at the IXP.

Therefore, an inbound route policy SHOULD be applied on IXP peerings in order to set the next-hop for accepted prefixes to the BGP peer IP address (belonging to the IXP LAN) that sent the prefix (which is what "next-hop-self" would enforce on the sending side).

This policy MUST NOT be used on route-server peerings, or on peerings where you intentionally permit the other side to send 3rd-party next-hops.

This policy also MUST be adjusted if Remote Triggered Black Holing best practice (aka RTBH [26]) is implemented. In that case one would apply a well-known BGP next-hop for routes it wants to filter (if an internet threat is observed from/to this route for example). This

well known next-hop will be statically routed to a null interface. In combination with unicast RPF check, this will discard traffic from and toward this prefix. Peers can exchange information about black-holes using for example particular BGP communities. One could propagate black-holes information to its peers using agreed BGP community: when receiving a route with that community one could change the next-hop in order to create the black hole.

10. BGP community scrubbing

Optionally we can consider the following rules on BGP AS-paths:

- o Scrub inbound communities with your AS number in the high-order bits - allow only those communities that customers/peers can use as a signaling mechanism
- o Do not remove other communities: your customers might need them to communicate with upstream providers. In particular do not (generally) remove the no-export community as it is usually announced by your peer for a certain purpose.

11. Possible future work

Following propositions were made and could be added to the document:

- o Appendix with IRRTOOLSET examples
- o Improve IRR section and clarify who should do what and recommendations on object management
- o Change "filters" in "import/export policies"
- o Give rationale against community scrubbing
- o Integrate comments of Donald Smith regarding TTL security and MD5/TCP-AO

12. Change logs

12.1. Diffs between [draft-jdurand-bgp-security-01](#) and [draft-jdurand-bgp-security-00](#)

Following changes have been made since previous document [draft-jdurand-bgp-security-00](#):

- o "This documents" typo corrected in the former abstract
- o Add normative reference for [RFC5082](#) in former [section 3.2](#)

- o "Non routable" changed in title of former [section 4.1.1](#)
- o Correction of typo for IPv4 loopback prefix in former [section 4.1.1.1](#)
- o Added shared transition space 100.64.0.0/10 in former [section 4.1.1.1](#)
- o Clarification that 2002::/16 6to4 prefix can cross network boundaries in former [section 4.1.1.2](#)
- o Rationale of 2000::/3 explained in former [section 4.1.1.2](#)
- o Added 3FFE::/16 prefix forgotten initially in the simplified list of prefixes that MUST not be routed by definition in former [section 4.1.1.2](#)
- o Warn that filters for prefixes not allocated by IANA MUST only be done if regular refresh is guaranteed, with some words about the IPv4 experience, in former [section 4.1.2.1](#)
- o Replace RIR database with IRR. A definition of IRR is added in former [section 4.1.2.2](#)
- o Remove any reference to anti-spoofing in former [section 4.1.4](#)
- o Clarification for IXP LAN prefix and pMTUd problem in former [section 4.1.5](#)
- o "Autonomous filters" typo (instead of Autonomous systems) corrected in the former [section 4.2](#)
- o Removal of an example for manual address validation in former [section 4.2.2.1](#)
- o [RFC5735](#) obsoletes [RFC3300](#)
- o Ingress/Egress replaced by Inbound/Outbound in all the document

12.2. Diffs between [draft-jdurand-bgp-security-02](#) and [draft-jdurand-bgp-security-01](#)

Following changes have been made since previous document [draft-jdurand-bgp-security-01](#):

- o 2 documentation prefixes were forgotten due to errata in [RFC5735](#). But all prefixes were removed from that document which now point

to other references for sake of not creating a new "registry" that would become outdated sooner or later

- o Change MD5 section with global TCP security session and introducing TCP-AO in former [section 3.1](#). Added reference to [BCP38](#)
- o Added new [section 3](#) about BGP router protection with forwarding plane ACL
- o Change text about prefix acceptable specificity in former [section 4.1.3](#) to explain this doc does not try to make recommendations
- o Refer as much as possible to existing registries to avoid creating a new one in former [section 4.1.1.1](#) and 4.1.1.2
- o Abstract reworded
- o 6to4 exception described (only more specifics MUST be filtered)
- o More specific -> more specifics
- o should -> MUST for the prefixes an ISP needs to filter from its customers in former [section 4.2.2.1](#)
- o Added "plus some headroom to permit growth" in former [section 7](#)
- o Added new section on Next-Hop filtering

12.3. Diffs between [draft-ietf-opsec-bgp-security-00](#) and [draft-jdurand-bgp-security-02](#)

Following changes have been made since previous document [draft-jdurand-bgp-security-02](#):

- o Added a subsection for RTBH in next-hop section with reference to [RFC6666](#)
- o Changed last sentence of introduction
- o Many edits throughout the document
- o Added definition of tier 1 transit provider
- o Removed definition of a BGP peering
- o Removed description of routing policies for IPv6 prefixes in IANA special registry as this now contains a routing scope field

- o Added reference to [RFC6598](#) and changed the IPv4 prefixes to be filtered by definition section
- o IXP added in acronym/definition section and only term used throughout the doc now

12.4. Diffs between [draft-ietf-opsec-bgp-security-01](#) and [draft-ietf-opsec-bgp-security-00](#)

Following changes have been made since previous document [draft-ietf-opsec-bgp-security-00](#):

- o Obsolete [RFC2385](#) moved from normative to informative reference
- o Clarification of preference of TCP-AO over MD5 in former [section 4.1](#)
- o Mentioning KARP efforts in TCP session protection section in former [section 4](#) and adding 3 RFC as informative references: 6518, 6862 and 6952
- o Removing reference to SIDR working-group
- o Better dissociating origin validation and path validation to clarify what's potentially available for deployment
- o Adding that SIDR mechanisms should be implemented in addition to the other ones mentioned throughout this document
- o Added a paragraph in former [section 8](#) about ASN renumbering
- o Change of security considerations section
- o Added the newly created IANA IPv4 Special Purpose Address Registry instead of references to RFCs listing these addresses

12.5. Diffs between [draft-ietf-opsec-bgp-security-02](#) and [draft-ietf-opsec-bgp-security-01](#)

Following changes have been made since previous document [draft-ietf-opsec-bgp-security-01](#):

- o Added a reference to [draft-ietf-sidr-origin-ops](#)
- o Added a reference to [RFC6811](#) and [RFC6907](#)
- o Changes "Most of RIR's" to "A majority of RIR's" on IRR availability

- o Various edits
- o Added NIST BGP security recommendations document
- o Added that it's possible to get info from ISPs from RADB
- o Correction of the url for IPv4 special use prefixes repository
- o Clarification of the fact only prefixes with Global Scope set to False MUST be discarded
- o IANA list could be pulled directly by routers (not just pushed on routers).
- o Warning added when prefixes are checked against IRR
- o Recommend network operators to sign their routing objects
- o Recommend network operators to publish their routing objects in IRR of their IRR when available
- o Dissociate rules for local AS and downstreams in former [section 5.1.4](#)

[13.](#) Acknowledgements

The authors would like to thank the following people for their comments and support: Marc Blanchet, Ron Bonica, Randy Bush, David Freedman, Wesley George, Daniel Ginsburg, David Groves, Mike Hugues, Joel Jaeggli, Tim Kleefass, Warren Kumari, Jacques Latour, Jerome Nicolle, Hagen Paul Pfeifer, Thomas Pinaud, Carlos Pignataro, Jean Rebiffe, Donald Smith, Kotikalapudi Sriram, Matjaz Straus, Tony Tauber, Gunter Van de Velde, Sebastian Wiesinger, Matsuzaki Yoshinobu.

Authors would like to thank once again Gunter Van de Velde for presenting the draft at several IETF meetings in various working groups, indeed helping dissemination of this document and gathering of precious feedback.

[14.](#) IANA Considerations

This memo includes no request to IANA.

15. Security Considerations

This document is entirely about BGP operational security. It depicts best practices that one should adopt to secure its BGP infrastructure: protecting BGP router and BGP sessions, adopting consistent BGP prefix and AS-path filters and configure other options to secure the BGP network.

On the other hand this document doesn't aim at depicting existing BGP implementations and their potential vulnerabilities and ways they handle errors. It will not detail how protection could be enforced against attack techniques using crafted packets.

16. References

16.1. Normative References

- [1] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997, <<http://xml.resource.org/public/rfc/html/rfc2119.html>>.
- [2] Rose, M., "Writing I-Ds and RFCs using XML", [RFC 2629](#), June 1999.
- [3] Carpenter, B. and K. Moore, "Connection of IPv6 Domains via IPv4 Clouds", [RFC 3056](#), February 2001.
- [4] Huitema, C. and B. Carpenter, "Deprecating Site Local Addresses", [RFC 3879](#), September 2004.
- [5] Hinden, R. and B. Haberman, "Unique Local IPv6 Unicast Addresses", [RFC 4193](#), October 2005.
- [6] Rekhter, Y., Li, T., and S. Hares, "A Border Gateway Protocol 4 (BGP-4)", [RFC 4271](#), January 2006.
- [7] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", [RFC 4291](#), February 2006.
- [8] Huitema, C., "Teredo: Tunneling IPv6 over UDP through Network Address Translations (NATs)", [RFC 4380](#), February 2006.
- [9] Gill, V., Heasley, J., Meyer, D., Savola, P., and C. Pignataro, "The Generalized TTL Security Mechanism (GTSM)", [RFC 5082](#), October 2007.

- [10] Touch, J., Mankin, A., and R. Bonica, "The TCP Authentication Option", [RFC 5925](#), June 2010.
- [11] Mohapatra, P., Scudder, J., Ward, D., Bush, R., and R. Austein, "BGP Prefix Origin Validation", [RFC 6811](#), January 2013.
- [12] "Internet Exchange Route Server", <<http://tools.ietf.org/id/draft-ietf-idr-ix-bgp-route-server-00.txt>>.

16.2. Informative References

- [13] Crocker, D., Ed. and P. Overell, "Augmented BNF for Syntax Specifications: ABNF", [RFC 2234](#), November 1997.
- [14] Heffernan, A., "Protection of BGP Sessions via the TCP MD5 Signature Option", [RFC 2385](#), August 1998.
- [15] Ferguson, P. and D. Senie, "Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing", [BCP 38](#), [RFC 2827](#), May 2000.
- [16] Huston, G., Lord, A., and P. Smith, "IPv6 Address Prefix Reserved for Documentation", [RFC 3849](#), July 2004.
- [17] Blunk, L., Damas, J., Parent, F., and A. Robachevsky, "Routing Policy Specification Language next generation (RPSLng)", [RFC 4012](#), March 2005.
- [18] Crocker, D., Ed. and P. Overell, "Augmented BNF for Syntax Specifications: ABNF", [RFC 4234](#), October 2005.
- [19] Blanchet, M., "Special-Use IPv6 Addresses", [RFC 5156](#), April 2008.
- [20] Cotton, M. and L. Vegoda, "Special Use IPv4 Addresses", [RFC 5735](#), January 2010.
- [21] Arkko, J., Cotton, M., and L. Vegoda, "IPv4 Address Blocks Reserved for Documentation", [RFC 5737](#), January 2010.
- [22] Dugal, D., Pignataro, C., and R. Dunn, "Protecting the Router Control Plane", [RFC 6192](#), March 2011.
- [23] Lepinski, M. and S. Kent, "An Infrastructure to Support Secure Internet Routing", [RFC 6480](#), February 2012.

- [24] Lebovitz, G. and M. Bhatia, "Keying and Authentication for Routing Protocols (KARP) Design Guidelines", [RFC 6518](#), February 2012.
- [25] Weil, J., Kuarsingh, V., Donley, C., Liljenstolpe, C., and M. Azinger, "IANA-Reserved IPv4 Prefix for Shared Address Space", [BCP 153](#), [RFC 6598](#), April 2012.
- [26] Hilliard, N. and D. Freedman, "A Discard Prefix for IPv6", [RFC 6666](#), August 2012.
- [27] Lebovitz, G., Bhatia, M., and B. Weis, "Keying and Authentication for Routing Protocols (KARP) Overview, Threats, and Requirements", [RFC 6862](#), March 2013.
- [28] Manderson, T., Sriram, K., and R. White, "Use Cases and Interpretations of Resource Public Key Infrastructure (RPKI) Objects for Issuers and Relying Parties", [RFC 6907](#), March 2013.
- [29] Jethanandani, M., Patel, K., and L. Zheng, "Analysis of BGP, LDP, PCEP, and MSDP Issues According to the Keying and Authentication for Routing Protocols (KARP) Design Guide", [RFC 6952](#), May 2013.
- [30] Smith, P. and C. Panigl, "RIPE-378 - RIPE Routing Working Group Recommendations On Route-flap Damping", May 2006.
- [31] Smith, P., Evans, R., and M. Hughes, "RIPE-399 - RIPE Routing Working Group Recommendations on Route Aggregation", December 2006.
- [32] Smith, P. and R. Evans, "RIPE-532 - RIPE Routing Working Group Recommendations on IPv6 Route Aggregation", November 2011.
- [33] Smith, P., Bush, R., Kuhne, M., Pelsser, C., Maennel, O., Patel, K., Mohapatra, P., and R. Evans, "RIPE-580 - RIPE Routing Working Group Recommendations On Route-flap Damping", January 2013.
- [34] Doering, G., "IPv6 BGP Filter Recommendations", November 2009, <<http://www.space.net/~gert/RIPE/ipv6-filters.html>>.
- [35] "IANA IPv6 Address Space", <<http://www.iana.org/assignments/ipv6-address-space/ipv6-address-space.xml>>.

- [36] "IANA IPv4 Special Purpose Address Registry",
<<http://www.iana.org/assignments/iana-ipv4-special-registry/iana-ipv4-special-registry.xhtml>>.
- [37] "IANA IPv6 Special Purpose Address Registry",
<<http://www.iana.org/assignments/iana-ipv6-special-registry/iana-ipv6-special-registry.xml>>.
- [38] "IANA IPv4 Address Space Registry", <<http://www.iana.org/assignments/ipv4-address-space/ipv4-address-space.xml>>.
- [39] "IANA IPv6 Address Space Registry", <<http://www.iana.org/assignments/ipv6-unicast-address-assignments/ipv6-unicast-address-assignments.xml>>.
- [40] "Routing Assets Database", <<http://www.radb.net>>.
- [41] "Making Route Flap Damping Usable", <<http://tools.ietf.org/id/draft-ietf-idr-rfd-usable-04.txt>>.
- [42] "Security Requirements for BGP Path Validation",
<<http://datatracker.ietf.org/doc/draft-ietf-sidr-bgpsec-reqs/>>.
- [43] "Threat Model for BGP Path Security",
<<http://datatracker.ietf.org/doc/draft-ietf-sidr-bgpsec-threats/>>.
- [44] "Autonomous System (AS) Migration Features and Their Effects on the BGP AS_PATH Attribute",
<<http://datatracker.ietf.org/doc/draft-ga-idr-as-migration/>>.
- [45] "RPKI-Based Origin Validation Operation",
<<http://tools.ietf.org/html/draft-ietf-sidr-origin-ops-23/>>.
- [46] , , and , "Border Gateway Protocol Security",
<<http://csrc.nist.gov/publications/nistpubs/800-54/SP800-54.pdf>>.

Authors' Addresses

Jerome Durand
CISCO Systems, Inc.
11 rue Camille Desmoulins
Issy-les-Moulineaux 92782 CEDEX
FR

Email: jerduran@cisco.com

Ivan Pepelnjak
NIL Data Communications
Tivolska 48
Ljubljana 1000
Slovenia

Email: ip@ipspace.net

Gert Doering
SpaceNet AG
Joseph-Dollinger-Bogen 14
Muenchen D-80807
Germany

Email: gert@space.net

