

# Large BGP Communities & Shutdown Communication.

David Freedman

david.freedman@uk.clara.net

Claranet

# Network Operators Use BGP Communities

- [RFC 1997](#) style communities have been available for the past 20 years
  - Encodes a 32-bit value displayed as: “16-bit ASN:16-bit value”
  - Designed to simplify Internet routing policies
  - Signals routing information between networks so that an action can be taken
- Broad support in BGP implementations
- Widely deployed and required by network operators for Internet routing

Community	Local-pref	Description
(default)	120	customer
65520:nnnn	50	only within country <nnnn> (see country list below)
65530:nnnn	50	only within region <nnnn> (see region list below)
2914:435	50	only beyond the connected country
2914:436	50	only beyond the connected region
2914:450	96	customer fallback
2914:460	98	peer backup
2914:470	100	peer
2914:480	110	customer backup
2914:490	120	customer default
2914:666		<a href="#">blackhole</a>

RFC 1997 Communities Examples

# Needed RFC 1997 Style Communities, but Larger

- We knew we'd run out of 16-bit ASNs eventually and came up with 32-bit ASNs
  - RIRs started allocating 32-bit ASNs by request in 2007, no distinction between 16-bit and 32-bit ASNs now
- However, you can't fit a 32-bit value into a 16-bit field
  - Can't use native 32-bit ASNs with RFC 1997 communities
- Needed an Internet routing communities solution for 32-bit ASNs for almost 10 years
  - Parity and fairness so everyone can use their globally unique ASN



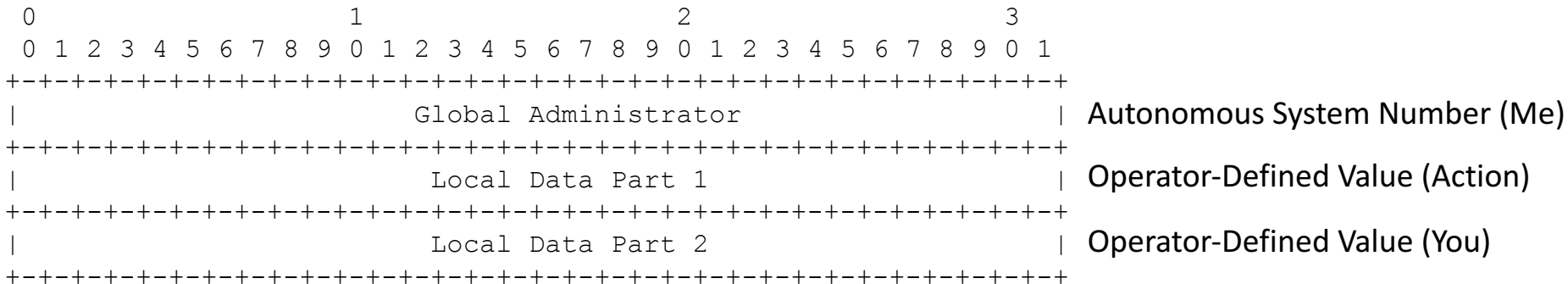
# The Solution: [RFC 8092](#)

## “BGP Large Communities Attribute”

- Idea progressed rapidly from inception in March 2016
- First I-D in September 2016 to RFC publication on February 16, 2017 in **just seven months**
- Final standard, plus a number of implementation and tools developed as well
- Network operators can test and deploy the new technology now



# Encoding and Usage



- **A unique namespace** for all 16-bit and 32-bit ASNs
  - No namespace collisions between ASNs
- Large communities are encoded as a 96-bit quantity and displayed as “32-bit ASN:32-bit value:32-bit value”
- Canonical representation is \$Me:\$Action:\$You

# Planning for Large Communities

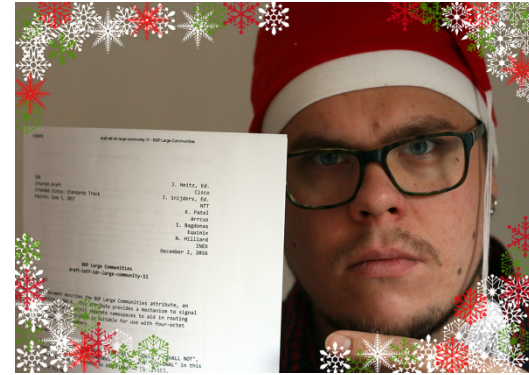
- The entire network ecosystem needs to support large communities in order to provision, deploy and troubleshoot them
- Ask your vendors and implementers for software support
- Update your tools and provisioning software
- Extend your routing policies, and openly publish this information
- Train your technical staff

Image sources: <https://www.sunet.se/blogg/all-i-want-for-christmas-is-large-bgp-communities/>

“All i want for christmas is ... Large BGP Communities” by Fredrik "Hugge" Korsbäck

20/04/2017

UKNOF37, Manchester



# Develop a Comprehensive Communities Policy

- Classic RFC 1997 communities will continue to be used together with large communities
  - There's no flag day to convert, large communities simply provide an additional way to signal information
- Your existing routing policy with classic communities is still valid
- Well-known communities such as “no-advertise”, “no-export”, “blackhole”, etc. are still used
- Extend your policy with large communities that allow network operators to signal the same information as they can with classic communities

# BGP Large Community Examples

RFC 1997 (Current)	BGP Large Communities	Action
65400: <i>peer-as</i>	2914:65400: <i>peer-as</i>	Do not Advertise to <i>peer-as</i> in North America (NTT)
43760: <i>peer-as</i>	43760:1: <i>peer-as</i>	Announce a prefix to a certain peer (INEX)
0:43760	43760:0: <i>peer-as</i>	Prevent announcement of a prefix to a certain peer (INEX)
65520: <i>nnn</i>	2914:65520: <i>nnn</i>	Lower Local Preference in Country <i>nnn</i> (NTT)
2914:410	2914:400:10	Route Received From a Peering Partner (NTT)
2914:420	2914:400:20	Route Received From a Customer (NTT)

- No namespace collisions or use of reserved ASNs
- Enables operators to use 32-bit ASNs in \$Me and \$You values



# Communities Policy Development

- [draft-ietf-grow-large-communities-usage](#) is a new [RFC 1998](#) style I-D in the IETF GROW Working Group
- Provides examples and inspiration for network operators to use large communities
- Also provides many examples on how to develop a communities policy
  - Informational communities
  - Action communities

#### 4.2.2. Location Based Selective AS\_PATH Prepending

AS 64497 might assign function identifier 7 to allow prepending of the AS\_PATH on propagation of routes to on any EBGP neighbor's interconnection in the geographical entity listed in the second Local Data field. This example follows the ISO 3166-1 numeric regions codes in the Local Data 2 field.

Large BGP Community	Meaning
64497:7:528	Prepend once to EBGP neighbors in the Netherlands
64497:7:392	Prepend once to EBGP neighbors in Japan
64497:7:840	Prepend once to EBGP neighbors in United States of America

Example documentation for AS 64497 offering Action Communities to trigger prepending of the AS\_PATH only when propagating the route to a certain geographical region.

Table 7: Action: Prepend in Region

# Informational Communities

- An informational label to mark a route with
  - Its origin: ISO 3166-1 numeric country ID and UIC M.49 geographic region
  - Relation or propagation: internal, customer, peer, transit
- Provides information for debugging or capacity planning
- The Global Administrator field is set to the ASN that labels the routes
- Most useful for downstream networks and the Global Administrator itself

# Information Communities Example

ISO 3166-1 Country ID		+	UN M.49 Region		+	Relation	
Large Community	Description		Large Community	Description		Large Community	Description
64497:1:528	Netherlands		64497:2:2	Africa		64497:3:1	Internal
64497:1:392	Japan		64497:2:9	Oceania		64497:3:2	Customer
64497:1:840	USA		64497:2:30	Eastern Asia		64497:3:3	Peering
			64497:2:150	Europe		64497:3:4	Transit

- For example, a communities value of “64497:1:528 64497:2:150 64497:3:2” would indicate that it was learned in the Netherlands, in Europe, from a customer

# CDN / Eyeball Example – You do a lot with 32 bits!

UK Postal Codes (~31 Bits)		or	GPS Coordinates	
Large Community	Postal Code		Large Community	Location
64497:9:849701135	E1W 1LB (London)		64497:10:1281024	Amsterdam
64497:9:1345374681	M90 1QX (Manchester)			(52.37783, 4.87995)

- Location encoding can be used to provide very accurate location information attached to more-specific routes announced to CDN caches
- UK postal codes can be encoded by stripping the whitespace and assuming they are base36 encoded, a decode results in a decimal.
- GPS coordinates can be encoded with [GeoHash](#)
  - For example 52.37783, 4.87995 (Amsterdam) encoded with 600 meter precision
  - Python: `import Geohash; Geohash.encode(52.37783, 4.87995, precision=6)`
  - Geohash result: `u173zp`
  - Convert `u173zp` from base36 to decimal = 1281024

# Action Communities

- An action label to request that a route be treated in a particular way within an AS
  - Propagation characteristics: export, selective export, no export
  - Local preference: influence ingress traffic within the AS
  - AS Path: influence traffic from outside the AS
- The Global Administrator field is set to the ASN which has defined the functionality of the community
  - Also is the AS that is expected to perform the action
- Most useful for transit providers taking action on behalf of a customer or the Global Administrator

# Action Communities Example

- Selective no export
  - ASN based selective no export
  - Location based selective no export
- Selective AS path prepending
  - ASN based selective AS path prepending
  - Location based selective AS path
- Local preference
  - Global local preference
  - Region based local preference

ASN Based No Export	
Large Community	Description
64497:4:64498	AS 64498
64497:4:64499	AS 64499
64497:4:65551	AS 65551

Location Based No Export	
Large Community	Description
64497:5:528	Netherlands
64497:5:392	Japan
64497:5:840	USA

# Getting Started With Large Communities

- 2018 is the year of large BGP communities
  - Preparation, testing, training and deployment can take weeks, months or even over a year
  - Start the work now, so you are ready when customers want to use large communities
- Lots of resources are available to help network operators learn about large communities
  - BGP speaker implementations
  - Analysis and ecosystem tools
  - Presentations (<http://largebgpcommunities.net/talks/>)
  - Documentation for each implementation
  - Configuration examples (<http://largebgpcommunities.net/examples/>)

# Large Communities Beacon Prefixes

- The following prefixes are announced with AS path 2914\_15562\$
  - 192.147.168.0/24 ([looking glass](#))
  - 2001:67c:208c::/48 ([looking glass](#))
  - BGP Large Community: 15562:1:1

## Cisco IOS Output (Without Large Communities Support)

```
route-views>show ip bgp 192.147.168.0
BGP routing table entry for 192.147.168.0/24, version 98399100
Paths: (39 available, best #30, table default)
  Not advertised to any peer
  Refresh Epoch 1
  701 2914 15562
    137.39.3.55 from 137.39.3.55 (137.39.3.55)
      Origin IGP, localpref 100, valid, external
      unknown transitive attribute: flag 0xE0 type 0x20 length 0xC
      value 0000 3CCA 0000 0001 0000 0001
      rx pathid: 0, tx pathid: 0
```

## BIRD Output (With Large Communities Support)

```
COLOCLUE1 11:06:17 from 94.142.247.3] (100/-) [AS15562i]
Type: BGP unicast univ
BGP.origin: IGP
BGP.as_path: 8283 2914 15562
BGP.next_hop: 94.142.247.3
BGP.med: 0
BGP.local_pref: 100
BGP.community: (2914,410) (2914,1206) (2914,2203) (8283,1)
BGP.large_community: (15562, 1, 1)
```



# BGP Speaker Implementation Status

Implementation	Software	Status	Details
Arista	<a href="#">EOS</a>	Planned	Feature Requested <a href="#">BUG169446</a>
Cisco	<a href="#">IOS XR</a>	✓ Done!	Beta (perhaps in 6.3.2 for real?)
cz.nic	<a href="#">BIRD</a>	✓ Done!	BIRD 1.6.3 ( <a href="#">commit</a> )
ExaBGP	<a href="#">ExaBGP</a>	✓ Done!	<a href="#">PR482</a>
FreeRangeRouting	<a href="#">frr</a>	✓ Done!	<a href="#">Issue 46</a> ( <a href="#">commit</a> )
Juniper	<a href="#">Junos OS</a>	Planned	Second Half 2017 (perhaps 17.3R1?)
MikroTik	<a href="#">RouterOS</a>	Won't Implement Until RFC	Feature Requested 2016090522001073
Nokia	<a href="#">SR OS</a>	Planned	Third Quarter 2017
nop.hu	<a href="#">freeRouter</a>	✓ Done!	
OpenBSD	<a href="#">OpenBGPD</a>	✓ Done!	OpenBSD 6.1 ( <a href="#">commit</a> )
OSRG	<a href="#">GoBGP</a>	✓ Done!	<a href="#">PR1094</a>
rtbrick	<a href="#">Fullstack</a>	✓ Done!	FullStack 17.1
Quagga	<a href="#">Quagga</a>	✓ Done!	Quagga 1.2.0 <a href="#">875</a>
Ubiquiti	<a href="#">EdgeOS</a>	Planned	<a href="#">Internal Enhancement Requested</a>
VyOS	<a href="#">VyOS</a>	Requested	Feature Requested <a href="#">T143</a>

Visit <http://largebgpcommunities.net/implementations/> for the Latest Status

# Tools and Ecosystem Implementation Status

Implementation	Software	Status	Details
DE-CIX	<a href="#">pbgpp</a>	✓ Done!	<a href="#">PR16</a>
FreeBSD	tcpdump	✓ Done!	<a href="#">PR213423</a>
Marco d'Itri	<a href="#">zebra-dump-parser</a>	✓ Done!	<a href="#">PR3</a>
OpenBSD	tcpdump	✓ Done!	OpenBSD 6.1 ( <a href="#">patch</a> )
pmacct.net	<a href="#">pmacct</a>	✓ Done!	<a href="#">PR61</a>
RIPE NCC	<a href="#">bgpdump</a>	✓ Done!	<a href="#">Issue 41</a> ( <a href="#">commit</a> )
tcpdump.org	<a href="#">tcpdump</a>	✓ Done!	<a href="#">PR543</a> ( <a href="#">commit</a> )
Yoshiyuki Yamauchi	<a href="#">mrtparse</a>	✓ Done!	<a href="#">PR13</a>
Wireshark	<a href="#">Dissector</a>	✓ Done!	18172 ( <a href="#">patch</a> )

Visit <http://largebgpcommunities.net/implementations/> for the Latest Status

# BGP

## Shutdown Communication

# Communication can be a challenge...



# Communication can be a challenge...

- [draft-nalawade-bgp-inform-02](#) – Died 2002 due to lack of adoption.
- [draft-nalawade-bgp-soft-notify-01](#) – Died 2005 due to lack of adoption.
- [draft-ietf-idr-advisory-00](#) – Adopted (IDR) in 2009. Died due to incorporation into [draft-frs-bgp-operational-message-00](#)
- [draft-ietf-idr-operational-message-00](#) – Adopted (IDR) in 2012. Died due to lack of progression.



# Get messaging back on the table

- *'The IETF has become a dumping ground for ideas. There are too many "researchers" in the IETF now. We don't implement every RFC anymore. The demand/complexity ratio is what counts now.'* – Anonymous large router vendor.
- Need something simple, effective, easy to implement...

IDR  
Internet-Draft  
Updates: 4486 (if approved)  
Intended status: Standards Track  
Expires: August 4, 2017

J. Snijders  
NTT  
J. Heitz  
Cisco  
J. Scudder  
Juniper  
January 31, 2017

**BGP Administrative Shutdown Communication  
draft-ietf-idr-shutdown-05**

Abstract

This document enhances the BGP Cease NOTIFICATION message "Administrative Shutdown" and "Administrative Reset" subcodes for operators to transmit a short freeform message to describe why a BGP session was shutdown or reset.

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 [RFC2119].



```

      0                1                2                3
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Error code 6 |   Subcode   |   Length   |   ...   | \
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+ /
\                                                                    \
/                               ... Shutdown Communication ...      /
\                                                                    \
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

**Subcode:** the Error Subcode value **MUST** be one of the following values: 2 ("Administrative Shutdown") or 4 ("Administrative Reset").

**Length:** this 8-bit field represents the length of the Shutdown Communication field in octets. The length value **MUST** range from 0 to 128 inclusive. When the length value is zero, no Shutdown Communication field follows.

**Shutdown Communication:** to support international characters, the Shutdown Communication field **MUST** be encoded using UTF-8. A receiving BGP speaker **MUST NOT** interpret invalid UTF-8 sequences. Note that when the Shutdown Communication contains multibyte characters, the number of characters will be less than the length value.

Mechanisms concerning the reporting of information contained in the Shutdown Communication are implementation specific but **SHOULD** include methods such as SYSLOG [RFC5424]

# Sending a shutdown communication

```
$ bgpctl neighbor 165.254.255.24 down \
```

```
    "[TICKET-1-1438367390] we are upgrading to  
openbsd 6.1, be back in 30 minutes"
```

```
request processed
```

# On the receiving side:

```
Jan  8 19:28:54 shutdown bgpd[50719]: neighbor  
165.254.255.26: received notification:  
Cease, administratively down
```

```
Jan  8 19:28:54 shutdown bgpd[50719]: neighbor  
165.254.255.26: received shutdown reason:  
" [TICKET-1-1438367390] we are upgrading to  
openbsd 6.1, be back in 30 minutes"
```

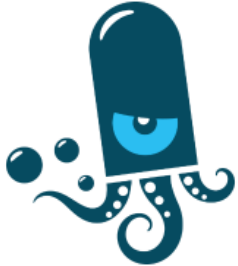
# Implementations so far...



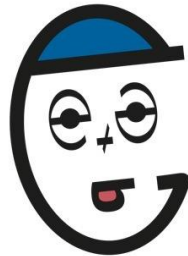
Openbsd / OpenBGPD



GoBGP



PMAcct



ExaBGP



Wireshark

**IETF Status:**  
(almost) Last call

Believed to be in the works:



And yes, UTF-8 / UNICODE works too...



# Questions?