BGP Fundamentals

Aaron Atac - Subspace



NANOG

About Me







About Me





Where I've been and where I am now





BGP (Border Gateway Protocol)

- Originally sketched over lunch by Kirk Lougheed and Yakov Rekhter, with the help of Len Bosack at IETF 12, January 1989.
- 32 years folks.
- AKA TNP (Three Napkin Protocol) Design by pragmatism
- Original napkins apparently got too much ketchup on them, they then wrote these three pieces of paper:

there first hop galaway metric metric trought count of AS Schirection (AS #	4 bytes 2 bytes 1 byte 1 byte 2 byte 2 byte	eut. "n times
hitriation: De oprode dute	2 byles variable	

	convection.	oot of	ayuc .	- dala is	lest	block	vereived
	(TCP elle	ose after	- perclast	sen +)			

- . open ronfinmed
- invalid block they type (date is I byte block type invalid version (date is (byte version)
- lougheed Q cisco. co XAFOV @IBH. CON

326-1941 (11-7) PST 5-3896 (8-5) FST



BGP at 18: Lessons In Protocol Design - Dr. Yakov Rekhter

BGP (Border Gateway Protocol)

block 2 bytes (second not) B.6. P block length block length block type mulory hold down timer (winutes) 2 bytas protorok version is rorreatly 1 undate - 2 types notification - 4 knopshie - 8 my 45 # 2 byte open: link type 1 byte down - Z (not used in uplate dissolving field) internal - 4 H-link - 8 outh type cale 1 byle 0 - none authentication voridale undate network # N bytes first hop gateway 4 bytes metric 2 bytes count of AS 1 ovfe direction 1 byte 3 vepeat "rout" 2 byle 15 # times notification: an oprode 2 bytes date variable

State Wingram link type error in opece - my view of correct link type (16yte) as initial glate in DISCONDECT unknown outh type coole OPEN - SEND - 40 dote disrouxoct authentication forture (no date) open update error - data is block in error vecu send confirm 5 parting toop the upstate Two phose error is optile defim data is subrode (Rbyte) followed by wait update block in grayhor (1 no twork only) raine involut notwork field oubrodes - 1 roctives 2. invalid first hop gw invalid direction rocke estob lishal invalid to routing loop two-phase error convection out of ayue - date is list block vacatured (TCP close after packet sout) oppu rontinued 6. involid block they type (date is 1 byte block type) 415-326-1941 (11-7) PST lougheed & cisco.com involid version number (date is I byte version) XAKOV@IBH.COM (914) 845-3896 (8-5)



BGP at 18: Lessons In Protocol Design - Dr. Yakov Rekhter

prid v

oot of syne (loral close

FOT

vousate dage

gang down

TCP trinal out

Gend a gend open confirm

The network at that time

- 13 Nodes
- 170 Networks
- By July 1988 all were interconnected at 1.5 Mbps (T-1)
- IBM RT PCs connected by token ring IEEE 802.5

Claffy, Kimberly C.; Braun, Hans-Werner; Polyzos, George C. (August 1994). "Tracking long-term growth of the NSFNET". *Communications of the ACM*.





What and why?

- EGP and Policy Based Routing in the New NSFNET Backbone
 <u>RFC1092</u> February 1989
 - EGP as reachability protocol
 - "It should be noted that the use of EGP is only viewed as an interim measure until better inter autonomous system protocols are defined and widely deployed for gateways used by regional networks."
 - "The EGP model assumes an engineered spanning tree topology, however, the NSFNET (due to the presence of backdoor routes) does not fit into this model."



What and why?

- The NSFNET Routing Architecture <u>RFC1093</u> February 1989
 - "In the longer run the hope is to replace the EGP interface with a new inter Autonomous System protocol. Such a new protocol should also allow to move the filtering of network numbers or Autonomous Network number groups to the regional gateways in order for the regional gateways to decide as to what routing information they wish to receive."



- A Border Gateway Protocol (BGP) <u>RFC1105</u> June 1989
- Incremental updates instead of periodic
- Use TCP/179 unicast as a reliable transport
- Have Autonomous Systems (AS) make up AS_PATHs
 - Used to provide information and prevent loops



- A Border Gateway Protocol (BGP) <u>RFC1163</u> June 1990
- Path Attributes
 - Mandatory vs optional, transitive vs non-transitive attributes
 - ORIGIN, AS_PATH, and NEXT_HOP are the only mandatory attributes.
- "The notion of Up/Down/Horizontal relations present in RFC1105 has been removed from the protocol."
- Marker field



- A Border Gateway Protocol 3 (BGP-3) <u>RFC1267</u> -October 1991
- Optimize and simplify the exchange of information about previously reachable routes
- Connection Collision Detection



- A Border Gateway Protocol 4 (BGP-4) <u>RFC1771</u> -March 1995
- Support CIDR!
 - Encode reachability information as variable length prefixes instead of fixed length based on classful networks
- LOCAL_PREF attribute



The Basics of Paths

- The mandatory attributes that have to be propagated
 - ORIGIN: IGP or EGP or Unknown
 - AS_PATH: Shortest path wins
 - NEXT_HOP
- Others
 - LOCAL_PREF: Highest wins
 - MULTI_EXIT_DISC (MED): Lowest wins
 - Communities: Value that you can filter on.
 - Some are "well-known" meaning they have been defined and reserved.
 - Examples: NO_ADVERTISE, NO_EXPORT, NO_PEER, BLACKHOLE, GRACEFUL_SHUTDOWN.



What were some of the original problems?

- Scalability / iBGP full-mesh requirement
 - Fix Autonomous System Confederations for BGP <u>RFC1965</u> -June 1996
 - Partition AS into sub-ASes
 - Only mesh within sub-AS
 - BGP Route Reflection An alternative to full mesh IBGP <u>RFC1966</u> -June 1996
 - Hub-and-spoke



What were some of the original enhancements?

- Multiprotocol Extensions for BGP-4 <u>RFC2283</u> -February 1998
 - BGP-4 was only carrying routing information for IPv4
- Protection of BGP Sessions via the TCP MD5 Signature Option <u>RFC2385</u> - August 1998
 - Marker field introduced in BGP-2
 - Limited value, under-specified, under-implemented



What were some of the original enhancements?

- Capabilities Advertisement with BGP-4 <u>RFC2842</u> May 2000
 - Version number wasn't going to be efficient.
 - Why we still are on BGP-4 in February of 2021
 - BGP speaker telling the other side its supported features
- Route Refresh Capability for BGP-4 <u>RFC2918</u> September 2000
- Carrying Label Information in BGP-4 <u>RFC3107</u> May 2001



Newer Enhancements

- Graceful Restart Mechanism for BGP <u>RFC4724</u> January 2007
- BGP Support for Four-octet AS Number Space <u>RFC4893</u> May 2007
- The TCP Authentication Option <u>RFC5925</u> June 2010
- An Infrastructure to Support Secure Internet Routing <u>RFC6480</u> -February 2012
- A Profile for Route Origin Authorizations (ROAs) <u>RFC6482</u> -February 2012
- BGP Prefix Origin Validation <u>RFC6811</u> January 2013



Newer Enhancements

- BGP Monitoring Protocol (BMP) <u>RFC7854</u> June 2016
- Advertisement of Multiple Paths in BGP <u>RFC7911</u> July 2016
- BGP Administrative Shutdown Communication <u>RFC8203</u> July 2017
- Default External BGP (EBGP) Route Propagation Behavior without Policies <u>RFC8212</u> - July 2017
- Graceful BGP Session Shutdown <u>RFC8326</u> March 2018
- Support for Adj-RIB-Out in the BGP Monitoring Protocol (BMP) <u>RFC8671</u> - November 2019
- Extended BGP Administrative Shutdown Communication <u>RFC9003</u> -January 2021



Filtering and Validity of Paths

- Received Prefixes
 - You don't have to accept everything you receive.
 - You can modify what you receive.
 - Example: Change the ORIGIN of a path matching a certain community.
 - Different networks will give you different views, and these views may not be all that they seem.
 - Views may contain short paths, but they may be geographically distant.
 - Views may contain effects of others misconfigurations.
 - Views may contain lossy paths.
 - Views may contain hijacked and/or leaked paths.



Filtering and Validity of Paths

- What to do with these received prefixes and views of all these different paths?
 - Filter malicious and illegitimate (aka bogon) routes
 - Apply business logic and traffic engineering desired
 - Don't leak them (unintentionally announcing prefixes to a different network that will cause undesirable traffic shifts)
- IRR/route/route6/AS-SET, <u>PeeringDB</u>, <u>bgpq4</u>, and prefix lists
- BGP Prefix Origin Validation <u>RFC6811</u> January 2013
- NLNOG BGP Filtering Guide



The Land of ISPs

- IX aka Internet Exchange
 - Internet Exchange BGP Route
 Server <u>RFC7947</u>
- PNI/SFI:
 - Private Network Interconnect
 - Settlement Free Interconnect
- IP Transit aka Transit aka DIA aka Direct Internet Access
- The more you peer directly with networks the less you will need to use your transit connection.



Unknown

♦ N A N O G[™]

The Land of ISPs

Asia Pacific Europe/Middle East/Central Asia/Africa North America Latin American and Caribbean RFC1918 IP Addresses Unknown





opte.org

Routing Policy

- Announcing and receiving prefixes sounds so simple...yet it's not when you look under the hood at what the decision-making humans do to influence paths.
- Traffic engineering applies both to how you announce and how you receive.



External Peering and Traffic Engineering

- Most specific always wins
- ORIGIN
- MED
- Communities
- AS_PATH prepending
 - BGP assumes shortest AS_PATH is best, not always the case
 - AS_PATH length doesn't always compare with geographic distance
 - Congestion
 - Loss
 - Return path
- Maybe I don't actually want to peer with you...at least at this location



Volatility of the DFZ

- Routes flap
 - Route Flap Damping exists, but is often too aggressive and does more harm than good. Vendor default settings are generally too reactive to small offenders.
 - Even if you're bringing down infrastructure it doesn't mean you must withdraw an aggregate. Announce aggregates externally, route to more specific internally.
 - If you really really must, use the <u>RIPE-580</u> guidelines, not the vendor defaults!
- Routes leak
 - This is unfortunately common and can subject eyeballs to degraded experiences.
- Routes get hijacked
 - Unfortunately common as well, fat-fingering happens, other times it's intentional, serial independent actors or govts/nation-states.



Volatility of the DFZ

- DFZ = Default Free Zone aka BGP Global Routing Table
- AS6447 RouteViews Wed Feb 24 02:20:00 2021 AET
 - IPv4: 903983
 - IPv6: 110993
- <u>Desperately Seeking Default</u> Geoff Huston APNIC
 - Poor orchestration and misguided security principles of some organizations and their routing policy can as well lead to the bgp global routing table looking a lot different depending on the feed you're looking at.



Public Information Gathering

- There's some pretty cool public collectors, feeds, and probes out there:
 - <u>University of Oregon Route Views Project</u>
 - <u>RIPE RIS Live, RIPEstat, RIPE Atlas</u> (shout out <u>Global</u> <u>Traceroute</u>)
 - <u>CAIDA BGPStream</u>, <u>ASRank</u>
 - <u>NLNOG RING</u>
 - The Oracle Internet Intelligence Routing 3D Visualization





en gitaszekaleginstálit = 6 szisz-it: Thakt = Austracy i anguagi, the*-

TEROPSERVICES;

CHARGET-CHARGET.UNIC HT HESTAGCOOK, IN FOTH &



Thank You

e e estatural constant = 5 syst. (): Chart = austract (anguare, the ')

CHUPSERVICES;

