Automated Inter-AS Traffic Engineering: An open source approach and operational considerations

Kostas Zorbadelos - Lead Network Architect May 2, 2023

CANAL+ Telecom https://www.peeringdb.com/asn/21351

Problem Statement & Design Goals

An open source automation tool for BGP Traffic Engineering

Considerations using on-demand BGP announcements

Problem Statement & Design Goals

Problem Description

- IP Network with multiple points of presense / geographically dispersed
- Service Provider network with customer transit services
- Having multiple Internet transit providers and peers (in IXes or PNIs)
- Varying costs in transit capacity, submarine capacity can also be involved
- In traditional ISP networks downstream traffic is dominant

Need to optimize incoming traffic streams and distribute them among available capacity. Also need to divert traffic on-demand for security reasons (eg DDoS attacks).

IP Network - Geolocations



CANAL+ Telecom (AS21351) geolocations

AS21351 - Border / Peering Points



Public Peering Excha	nge Points	Filter		
Exchange 12 IPv4	ASN IPv6	Speed	RS Peer	
Equinix Sydney 45.127.173.38	21351 2001:de8:8::2	10G 1351:1	0	
EL-DX 206.41.108.49	21351 2001:504:40:	40G 108::1:49	0	
France-IX Paris 37.49.237.94	21351 2001:718:54::	20G 1:94	0	
Guyanix 194.57.235.199	21351	100M	0	
MARTINIX 194.254.208.133	21351	100M	0	
MIXP 196.223.0.6	21351 2001:43/8:27	1G 0:d0d0::6	0	
NAPAtrica IX Durban 196.10.141.46	21351 2001:43/8:6d	40G 2::46	0	
NAPAtrica IX Johannesburg	21351	60G	0	
REUNIX	2001:438.6d 21351	2.889 1G	0	

Q

28

SY/AU

Facility 1 ASN	Country City
Equinix MI1 - Miami, NOTA	United States of America
21351	Miami
Equinix SY1/SY2 - Sydney	Australia
21351	Sydney
GNU Caverine Collery	French Guiana
21351	Cayenne
Interxion Paris 1 (Aubervilliers	France
Cedex)	Paris
21351	
Intervion Paris 2 (Aubervilliers)	France
21351	Paris
Interxion Paris 5 (St Denis)	France
21351	Paris
MCN Le Lamentin	Martinique
21351	Le Lamentin
Telehouse - Paris 2 (Voltaire -	France
Léon Frot)	Paris
21351	
Teraco Durban, South Africa	South Africa
21351	Durban
Teraco Johannesburg Campus.	South Africa
South Africa	Johannesburg

- Manual configuration on routers is cumbersome
- Inconsistent configuration, error-prone
- Routers involved could be many, fast reaction not possible
- Configuration could be performed by network operators or even a program without human involvement
- Ideally vendor neutral (multiple vendor equipment in many networks)

Design Goals (continued)

- Do the traffic engineering reliably, without errors
- Easy to operate
- Do it quickly, even real time, depending on current traffic conditions or security incidents
- A link failure (especially in submarine capacity) would need proper action to bypass the failure
- Optimize economics (in transit services, capacities or peerings)
- Build a tool on a solid foundation that can grow in features
- Provide automation / scripting capabilities (future "self driving network")

An open source automation tool for BGP Traffic Engineering

• BGP is the exterior routing protocol between ASes



- BGP is used extensively for traffic engineering with various tricks (hacks?)
- Tool's purpose is to automate BGP announcements to peers, to affect in-bound traffic flow

- Centralized configuration point
- Sources of truth, representing the *intended* state (peerings and announcements)
- Standardized BGP policy configuration generated by automation tools (OUT-bound policies)
- Tagging of prefixes (BGP communities) affects policy
- Only need to think (or generate) the proper tags in the routes to get the desired outcome
- Design flexibility, all traffic engineering tricks should be supported

- RFC 8092 BGP large communities [3] a "recent" development (2017)
- 12 octets, three 4-byte integers (example 21351:602:6799)
- Overcome policy design limitations with 32-bit ASNs
- RFC 8195 Use of BGP Large Communities [8], informational RFC giving excellent policy examples
- Informational and action communities
- An IETF "blessed" way to create policies!
- Our design was based on this

Informational vs action communities

- Following RFC 8195 paradigm, second number in the large community is a field that contains a *function* identifier
- Informational communities are labels for various attributes
- Action Communities are added as labels to request that a route be treated in a particular way within an AS

```
*Informational communities example*
<ASN>:3:<TYPE_OF_ROUTE>
Contains the type of a route (eg internal loopback,
internal b2b customer, transit customer route,
BGP announcement)
```

```
*Action community example*
<ASN>:40:<PEER_ASN>
*Do not* announce a route to a peer ASN
```

Large community pattern					
<localasn>:40:0</localasn>					
<localasn>:40:<peerasn></peerasn></localasn>					
<localasn>:41:<peerasn></peerasn></localasn>					
<localasn>:6[N]:<peerasn></peerasn></localasn>					
<localasn>:400:0</localasn>					
<localasn>:400:<locationcode></locationcode></localasn>					
<localasn>:60[N]:<locationcode></locationcode></localasn>					

Tool implementation - Sources of Truth



We utilize NETBOX [6] as an IPAM system. Each prefix announcement is tagged accordingly, using BGP large communities. NETBOX contains the intended state of all the announcements of our AS (prefixes and policy for them).



We utilize Peering Manager [5] to hold all the information regarding eBGP peerings with transit providers and peers. We document both PNIs and peerings via Internet Exchanges. From this the configuration management engine generates configuration for the peerings plus the standardized OUT-bound policies.

Configuration management engine



- A lot of open source configuration management frameworks
- Salt (sometimes referred to as SaltStack) an open-source software for event-driven IT automation, remote task execution, and configuration management
- NAPALM is a vendor neutral, cross-platform open source project that provides a unified API to network devices
- All Python based
- Development based on Salt/NAPALM using Jinja templates [7]

https://github.com/kzorba/bgp-te-tool [9]

- Tool code, documentation and demonstration
- Simulated network using docker containers and docker-compose
- goBGP containers simulate peers and transit providers
- Current implementation supports Juniper routers (JunOS jinja templates)
- Contributions (eg other vendor support) highly welcome!

Tool usage - Peerings

٥			Home	- Peering	Manager - Chromium				* - ¤ ×
A Home - Peering Manag ×	+								· - · ·
← → C ▲ Not secure http	://peering-man	ager.infra.msv	narammina 🔺 Chra	nor - R	ProFiver - N Ø Sign In	B AkinitaCD	lohr BLCV	< 🛠 💀 🛛 4	>> LI 🚢 i
Peering Manager			ogramming V circ	103 1 9	, rivitikas it i sigiriti	Admittedit	3003 1 07	Q Last Search = AS2	21351 - 🚢 -
Autonomous Systems	Logged in	as kzorba.							×
BGP Groups									_
Internet Exchanges	Peering	Data			Deployment		Polic	cy Options	
Provisioning •	Autonon Networks	to peer with	6	7	Configurations Templates to build router config	urations () Routi Policie	ing Policies es filtering advertised/received rou	tes 52
Policy Options	200 000				e acatta				_
Deployment +	Groups of	BGP sessions	1	7	e-mails Templates to build e-mails	C) Tags f	for traffic engineering	7
3rd Party •	Internet	Exchange Poin	ls.		Routers				
Other •	Infrastruc	tures allowing p	eering	9	Network devices running BGP	11	1		
	Direct Pe BGP sessi IXP Peeri BGP sessi	ering Sessions ons for transit, P ng Sessions ons setup over I	Nis, etc. 4 XPs 21	5 4					
	Change	log							
	User	Action	Туре		Object			Time	
fr-1vm-peeringadm01.infra.msv	kzorba Updated Internet Exchange P		Internet Exchange P	ering Session	n France-IX Paris	- AS51706 - IP 200	1:7f8:54:251	2022-10-05 15:36	
(V1.5.2)	kzorba	kzorba Updated Internet Exchange Peering Se		ering Session	n France-IX Paris	- AS51706 - IP 200	2022-10-05 15:36		
API · E Docs · O GitHub	kzorba	Updated	Internet Exchange P	ering Session	n France-IX Paris	- AS51706 - IP 200	1:7f8:54:251	2022-10-05 15:36	

BGP information in Peering Manager

Tool usage - BGP announcements in IPAM

0			154.€	57.0.0/17 -	NetBox - Chi	romiun						+ -	• ×
⊯ 154.67.0.0/17 - NetBo × +													~
← → ♂ ▲ Not secure https://netboxinfra.msv/ipam/	iggregates/19/								< *	N O	4 *	•	8 E
🖿 Office365+ 🖿 Networking 🖿 C+ Telecom	🖿 Programming 🚸 C	hronos - P	🕇 Prefixes	- N 🕲 Sigr	n In 🖿 Akinita	GR 🖿 Job	os 🆿 CV 🚺	Calculator	f 🔳 How to Con 🔀 Liste de con		🖿 Other	bookn	narks
Transaction - Organization -	Devices - IPAM - V	firtualization +	Circuits +	Power - S	ecrets - Other	-			Search Q & k	zorba +			Ì
Aggregates / AFRINIC / 154.67	.0.0/17								Search aggregates	۹			
154.67.0.0/17									+ Clone / Edit 0	Delete			
Created April 10, 2020 - Updated 4 r	nonths, 2 weeks ago												
Aggregate Change Log									Show available Hide ava	ilable			
Aggregate					Tags					-			
Family IPv4					alloca	tion 21351	1:3:1999 213	51:40:32812	6 preference:130 21351:400:663 21351:400:8	340			
RIR AFRIND	5					-				/			
Utilization 0%								_					
Date Added Sept. 30	, 2013												
Description AFRINI	C Allocation												
Child Prefixes													
Prefix	Status	VRF	Utilizatio	m	Tenant	Site	VLAN	Role	Description				
154,67.0.0/24	Available	Global	-		-	-	-	-	-				
154.67.1.0/28	Active	Global		0%	-	-	-	-	Cashas 017 0- 00886 800				
154,67.1.16/28	Available	Global	-		-	-	_	-					
54.67.1.32/28	Active	Global		0%	-		-	-	Cashes 017 C+ 00886 0PL				
154.67.1.48/28	Available	Global	-		-	-	-	-	-				
154.67.1.64/28	Available	Global	-		-	-	-	-	-				
154.67.1.128/26	Available	Global	_		-	-	-	-	-				
IT 154 67 1 100/08	Action	Global		0%	_	-	-	_	Participation (1971 Participation)				

Prefix large community tagging in netbox

Network configuration with Salt



Applying state in a subset of routers

Production Network rollout

- Gradual deployment in AS21351 (router by router and location by location)
- Operations currently handled by the engineering team, a very small circle
- Training to operations teams will follow
- Possibility to rollout the tool POP by POP and in each (geo) location at a time was very beneficial for controlled deployment
- Up until now, rollout was smooth and controlled
- Tool currently handles ~260 peerings in 6 diverse geo-locations and 6 IXes

- Design (mostly) and implementation not trivial
- Very limited human resources and day to day operations required attention
- Very demanding preparations tasks in the network to prepare the ground for the tool introduction
- 90% of the time was low level details and thinking
- Operations now require a paradigm shift (not easy)
- After all the work, the basis is there for future development as well

Considerations using on-demand BGP announcements

- Good MANRS: Route policy, contacts and intended announcements SHOULD be documented in IRR. Accurate route filtering necessary
- What about on-demand announcements?
- Sometimes, not practical to pre-provision every possible route/route6 object
- Updating the IRR on the time of need leads to late reaction (upstream filter updates?)
- RPKI ROAs another source of validation

- RPKI ROA maxLength can provide the necessary flexibility for TE
- Security implications (forged-origin prefix/subprefix hijack)
- draft-ietf-sidrops-rpkimaxlen RFC 9319/BCP 185: The Use of maxLength in the Resource Public Key Infrastructure (RPKI)
 [2]
- Greater harm for non announced address space
- Minimal ROAs recommended (whenever possible)

Experiences with IP transit providers

- Different levels of flexibility in transit provider networks
- From completely manual communication in case of need for a new announcement to various automated options in-between
- Respected providers care a lot about security and stability of the routing system (of course)
- Some providers give options to announce sub-prefixes and are responsive
- Automated generation of filters takes time and is performed a limited number of times within the day
- Flexibility required from customers

Monitoring of BGP Announcements

- Need a way to monitor our current BGP announcements
 - 1. as seen from the outside world (Internet)
 - 2. as sent from our own routers to peers
- For (1) various services exist (eg RIPE RIS live feed)
- For (2) best solution is BMP (BGP Monitoring Protocol [Adj-RIB-Out] [1], supported in pmacct [4] tool)



• Currently work in progress

- Best practices (in a document form) for operators and transit providers regarding on-demand announcements?
- Are the security implications with RPKI ROA maxLength a blocking point?
- Flexibility/operation agility vs security tradeoff
- How "real-time" should traffic engineering actions be and how often?
- AntiDDoS defences and big failures?

THANK YOU!

A special thanks to the authors/contibutors of the great open source tools available and the relevant communities.

QUESTIONS?

References

[1] T. Evens et al. RFC 8671: Support for Adj-RIB-Out in the BGP Monitoring Protocol (BMP). https://datatracker.ietf.org/doc/html/rfc8671. Nov.

2019.

[2] Y. Gilad et al. The Use of maxLength in the Resource Public Key Infrastructure (RPKI). https://datatracker.ietf.org/doc/rfc9319/.

References ii

[3] J. Heitz et al. *RFC 8092: BGP Large Communities Attribute.*

https://datatracker.ietf.org/doc/html/rfc8092.Feb. 2017.

- [4] Paolo Lucente. pmacct BMP daemon. http://www.pmacct.net/.
- [5] Guillaume Mazoyer and Contributors. Peering Manager. https://peering-manager.net.
- [6] **NETBOX.** https://github.com/netbox-community/netbox.
- [7] Network Automation with Salt.

https://docs.saltproject.io/en/3002/topics/network_
automation/index.html.

- [8] J. Snijders, J. Heasley, and M. Schmidt. RFC 8195: Use of BGP Large Communities. https://datatracker.ietf.org/doc/html/rfc8195. June 2017.
- Kostas Zorbadelos. A tool for automated BGP Traffic Engineering. https://github.com/kzorba/bgp-te-tool.