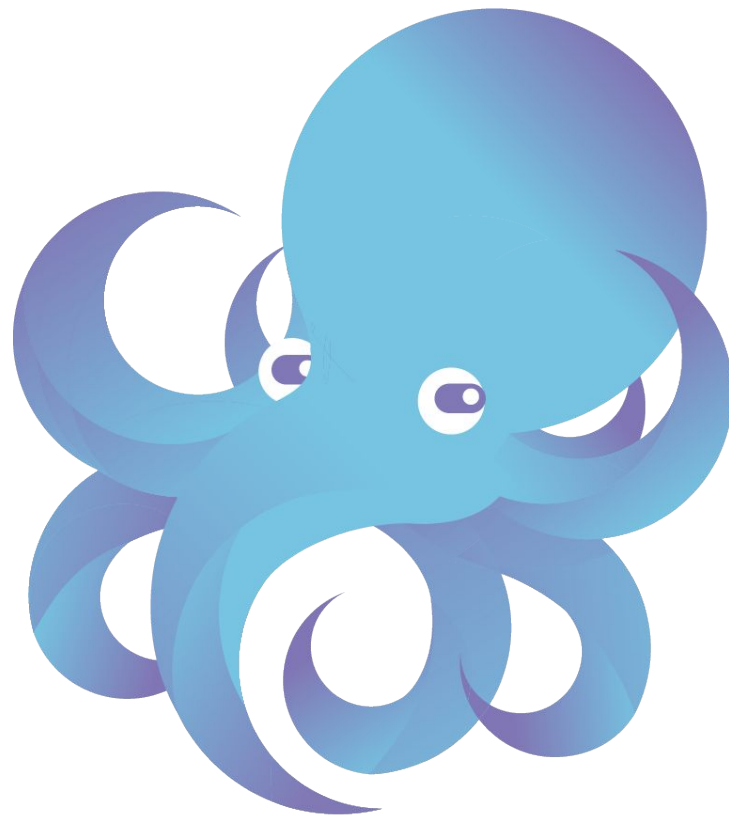




Cloudflare and RPKI at scale

Louis Poinsignon

Martin J. Levy



Introduction

Louis Poinsignon:

Network Engineer at Cloudflare in San Francisco

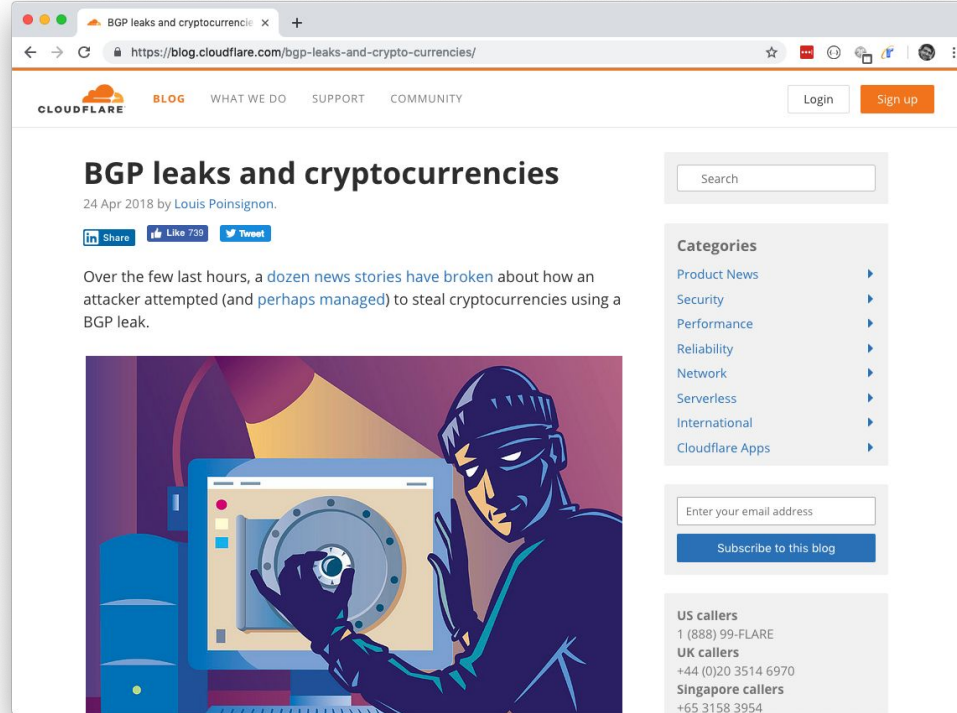
Open-source projects including flows and RPKI

Network data collection (BGP, flows, peering-portal)



Featured in RIPE's video!
<https://www.youtube.com/watch?v=Y9vbbxr-Gbl>

How did it start?



The Initial Story

Authoritative DNS route hijack in April 2018

DNS route announced via peering session (in Chicago)

This affected our network, hence our DNS Resolver

What should we do?

The Initial Story

At the time...

150+ PoPs

26,000 BGP sessions

IP space from five RIRs

Just the RIPE Validator ^[1]

How to distribute a prefix list efficiently?

The Initial Story

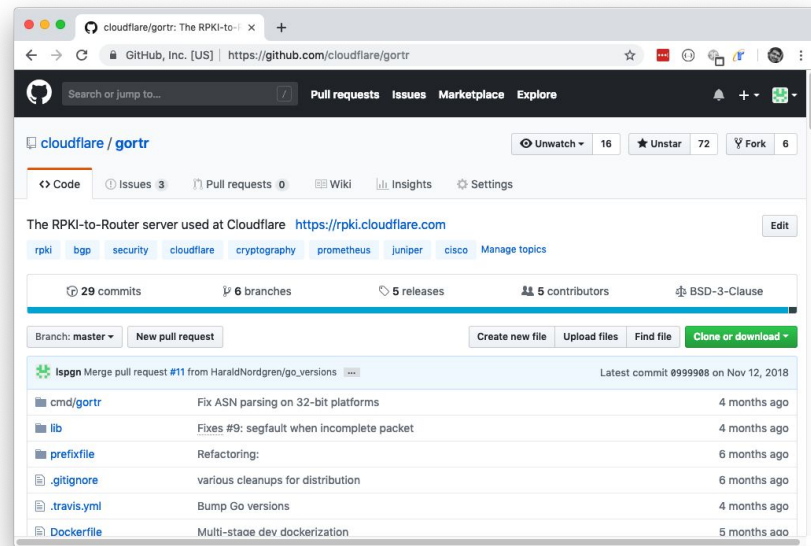
July: started deploying internally GoRTR.

August: open-source release.

<https://github.com/cloudflare/gortr>

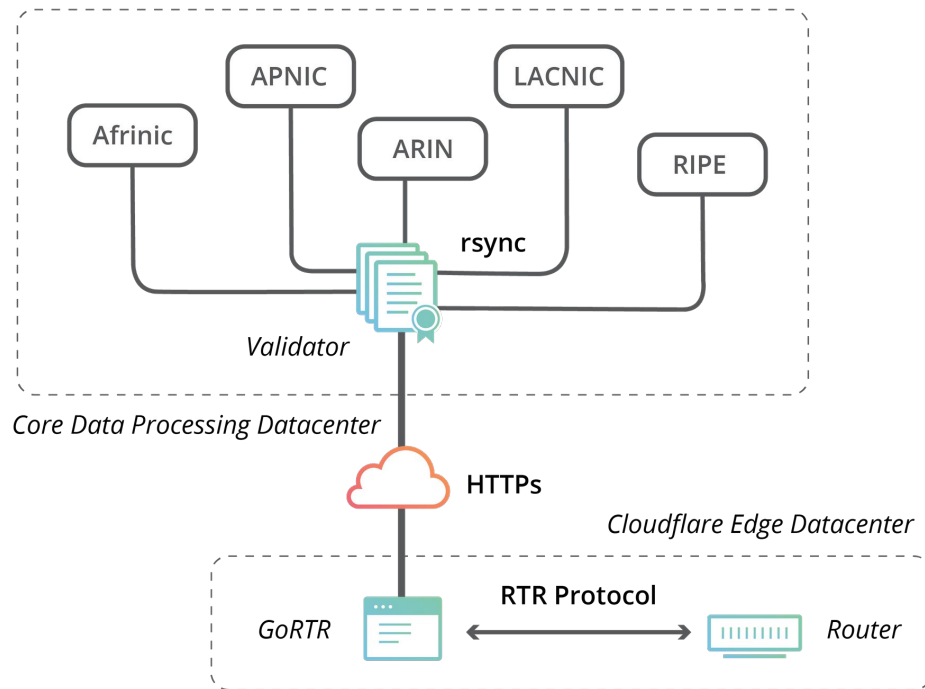
September → December:

- Turn up RTR sessions
- Signing prefixes`



<https://github.com/cloudflare/gortr>

Diagram



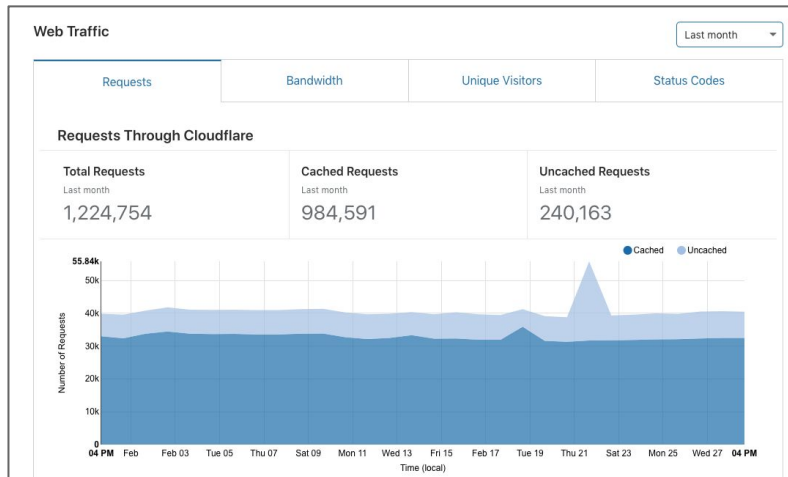
Behind the scene (until January 2019)

RIPE Validator providing list of prefixes

Running in a Mesos cluster

With a cronjob:

- Fetching the data
- Filtering (remove $> /24$ and $> /48$ and duplicates)
- Signing it
- Making it available to our edge



`https://rpki.cloudflare.com/rpki.json` was born.

Effects

The question everyone asked us:

How much traffic was affected?

Many invalids. Little traffic in practice

(we had a default or valid less specific)

Except in one place: Few gigabits per seconds displaced due to geographical more specific



<https://www.flickr.com/photos/thure/6287816628/>

Accounting

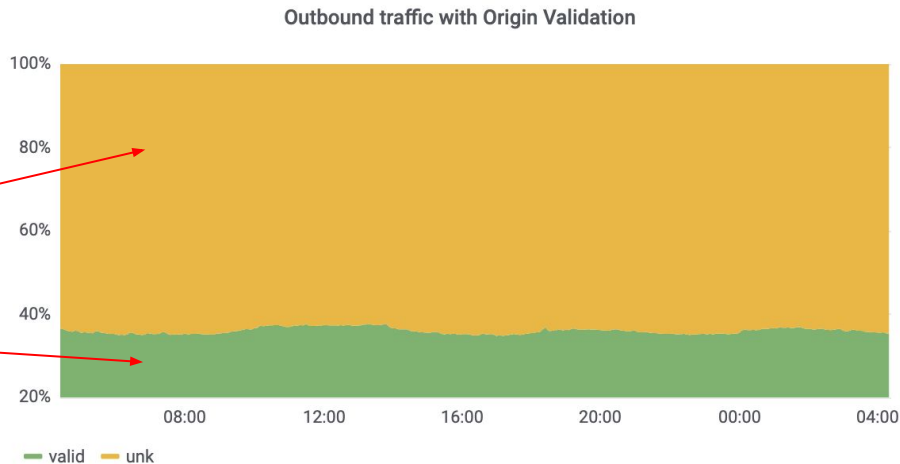
Using flows, we see **at least 30%** of our traffic being **valid**

Very **little/none invalid**

We use **GoFlow** for accounting. (Other tools compatible with flows: pmacct, Kentik, etc)

Traffic with a no ROA

Traffic with a valid ROA



Signing the routes

Signing the routes

Cloudflare has IP space from five RIRs

(no space from twnic/jpnic/cnnic)

Not a unified
experience!

RIR	Features	Ease of use	API
AFRINIC	★	★	★
APNIC	★★	★★	★
ARIN	★★	★★	★★
LACNIC	★	★★★★	★
RIPE	★★★	★★★★	★★★

Rankings

Features: RRDP, 2 factors, extra info, CA

Ease of use: steps to sign a ROA, multi user

API: functional, complete and accessible

Comparison - AFRINIC

Hard to set up: client TLS certificate to create (BPKI) in order to do RPKI.

Buggy.

No RRDP.

No API.

No auto-renew.

Hosted CA possible.





Extensive certificate informations.

Manage Your RPKI Resources

You can manage your RPKI resources from this page

RPKI Operations
List certificates
View ROAs
View Old ROAs
Issue ROA

Add ROA

* Name:	Please enter a unique ROA name. Spaces will be replaced by '_'. <input type="text"/>
Your AS Numbers:	Please select your ASN from this list or enter any other valid ASN in the field below. <input type="text" value="Select an ASN"/>
* AS Number:	ASN must be between 0 - 4294967295 in ASPLAIN format. "Reserved" and "Unallocated" ASNs will be rejected. <input type="text"/>
IPv4 address range:	Please select your prefix in the drop down list and click the '+' button, then you can specify the details <input type="text" value="Select an IPv4 prefix"/> 
IPv6 address range:	Please select your prefix in the drop down list and click the '+' button, then you can specify the details <input type="text" value="Select an IPv6 prefix"/> 
* Not Valid Before (YYYY-MM-DD):	<input type="text" value="YYYY-MM-DD"/> 
* Not Valid After (YYYY-MM-DD):	<input type="text" value="YYYY-MM-DD"/> 

Comparison - APNIC

Two factors or client certificate.

RRDP.

Auto-renew.

Allow BGP batch signing.

(slight bugs with large amount of prefixes).

Hosted CA possible.

Draft for API:

<https://www.apnic.net/manage-ip/apnic-services/services-roadmap/public-api-draft-for-members/>

Routes

Routes

Register your routes in MyAPNIC using the tool below. It will automatically create route objects authorized. RPKI ROAs will also be created at the same time, if the ROA option is enabled (changes to ROA status will not be updated until then).

Import routes

BGP announcements associated with your resources but not managed under this tool were found.

Review & Import from BGP

Dismiss

Create route

Prefix Route's prefix. E.g. 203.10.0.0/20

Origin AS Route's origin. E.g. AS123

Most specific announcement Route's most specific announcement. E.g. /22

ROA ☒ Enabled

Whois ☒ Enabled

☐ Define Whois route attributes

Options ☐ Notify additional contacts

Cancel

Next

Comparison - ARIN

Two factors. Separate signing key.

No RRDp.

No auto-renew.

Semi-functional API (add).

Dashboard not easy to find.

Hosted CA possible.

Slow rsync update (4 times a day).

Some certificate information.

Create a Route Origin Authorization

[Browser Signed](#) [Signed](#)

* denotes required field

***ROA Name:** ⓘ
Any name of your choosing.

***Origin AS:**
The AS Number you are authorizing.

***Start Date:**
The first date your ROA can be considered valid.

***End Date:**
The last date your ROA can be considered valid.

***Prefixes:** ⓘ
[+ Add Prefix](#)
The prefixes you authorize to originate from this AS.

***Private Key:** [Browse](#)
This key will not be uploaded to ARIN.

Comparison - LACNIC

No two factors. Single user.

No RRDP.

No API.

Auto-renew opt-in.

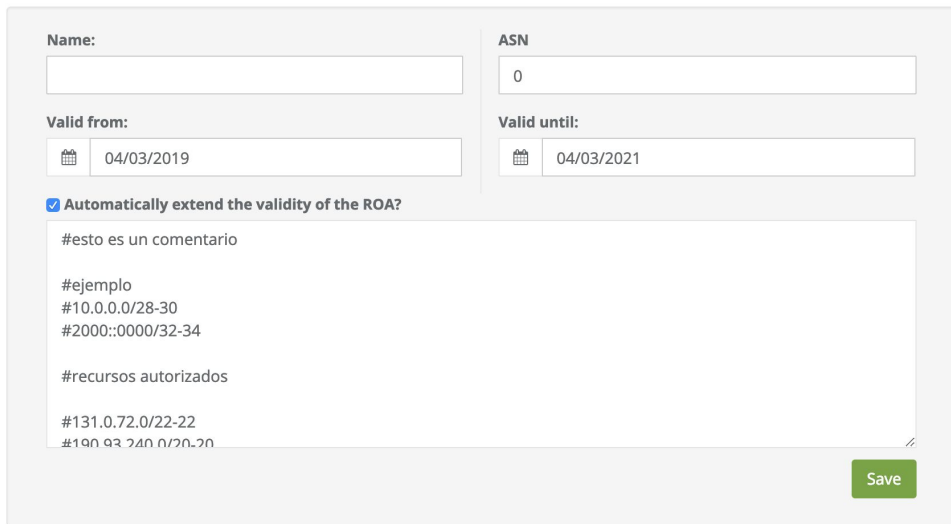
Allow BGP batch signing.

Based off RIPE.

No Hosted CA.

Some extra info (revoked, path).

Incorrect certificate encoding (BER). High turnover of certificate (few days).



The screenshot shows the LACNIC ROA (Route Origin Authorization) form. It includes fields for Name, ASN, Valid from, and Valid until. The Valid from and Valid until fields are set to 04/03/2019 and 04/03/2021 respectively. There is a checkbox for "Automatically extend the validity of the ROA?" which is checked. Below this is a large text area for comments and authorized resources. The text area contains the following content:

```
#esto es un comentario

#ejemplo
#10.0.0.0/28-30
#2000::0000/32-34

#recursos autorizados

#131.0.72.0/22-22
#190 93 240 0/20-20
```

A green "Save" button is located at the bottom right of the form.

Comparison - RIPE

Two factors.

RRDP.

Auto-renew.

Nice API.

Allow BGP batch signing.

No Hosted CA (theoretically).

No extra information. But history.

Incorrect certificate encoding (BER).

The screenshot shows the RIPE database's 'Route Origin Authorisations (ROAs)' management page. At the top, there are three tabs: 'BGP Announcements' (highlighted in teal), 'Route Origin Authorisations (ROAs)', and 'History'. To the right of these tabs is a search bar labeled 'Search...'. Below the tabs, there are two buttons: 'Discard Changes' and 'Delete ROAs'. To the right of these are two status filters: 'Causing Problems' (with a warning icon) and 'Not Causing Problems' (checked, with a checkmark icon). A '+ New ROA' button is on the far right. Below these elements is a table with four columns: 'AS number', 'Prefix', 'Most specific length allowed', and 'Affects'. At the bottom of the table, there are input fields for 'AS Number', 'Prefix', and 'Max length', followed by a save icon and a refresh icon.

Automation

We automated prefixes adding on **ARIN and RIPE** with a **Salt state**.

Two secrets to store (API key and signing key).

Cannot delete or list via API for ARIN: very prone to mistakes if user wants to reduce the amount of ROA files.

```
def _format_payload(roas, signature):
    template = """-----BEGIN ROA REQUEST-----
{roas}
-----END ROA REQUEST-----
-----BEGIN SIGNATURE-----
{signature}
-----END SIGNATURE-----
"""
    payload = template.format(
        roas=roas, signature="\n".join(textwrap.wrap(signature, width=64))
    )
    return payload

def _make_roa(name, asn, t, start_val, end_val, prefix, length, maxlength):
    template = (
        '1|{time}|{name}|{asn}|{start_val}|{end_val}|{prefix}|{length}|{maxlength}|'
    )
    time_str = calendar.timegm(t.timetuple())
    start_val_str = start_val.strftime(_TIME_FORMAT)
    end_val_str = end_val.strftime(_TIME_FORMAT)
    roa = template.format(
        time=time_str,
        name=name,
        asn=asn,
        start_val=start_val_str,
        end_val=end_val_str,
        prefix=prefix,
        length=length,
        maxlength=maxlength,
    )
    return roa

def _sign(pkey, roas):
    signature = pkey.sign(roas.encode('utf-8'), padding.PKCS1v15(), hashes.SHA256())
    return base64.b64encode(signature).decode('utf-8')
```

Validator

Why write a new validator?

November 2018: First release of NLnet Labs Routinator 3000 ^[1]

We were still using RIPE Validator

We wanted something more custom: with monitoring and RRDP

By building it in Golang:

- Many APIs and easy concurrency
- Community doing cryptography
- Cloudflare uses Golang a lot (cfssl, sidh, etc.)

Challenges

Juniper bugs: Routing Validation disabled

Difficulties: rsync, BER encoded instead of DER, conditions in cryptography

The TAL is an ordered sequence of:

- 1) a URI section,
- 2) a <CRLF> or <LF> line break,
- 3) a subjectPublicKeyInfo [[RFC5280](#)] in DER format [[X.509](#)], encoded in Base64 (see [Section 4 of \[RFC4648\]](#)). To avoid long lines, <CRLF> or <LF> line breaks MAY be inserted into the Base64-encoded string.

Cloudflare's RPKI Toolkit

Sets of libraries and tools written in Go

Including **OctoRPKI** 

<https://github.com/cloudflare/cfrpki>

Cloudflare's RPKI Toolkit

Libraries

- CER/ROA/MFT decoder
- PKI manager (exploring, validating)
- RRDP/rsync fetcher
- Validation of prefixes

Software

- Local validator (without RRDP/rsync)
- API tools for a distributed version without filesystem
- OctoRPKI
- Certificate Transparency tool



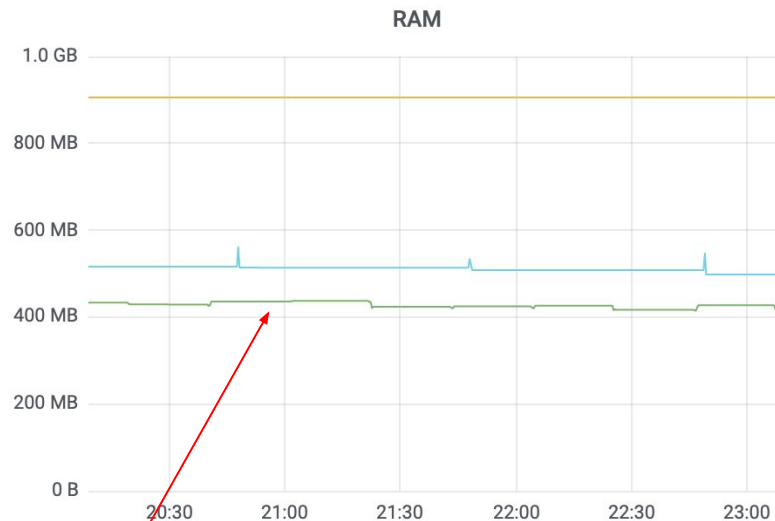
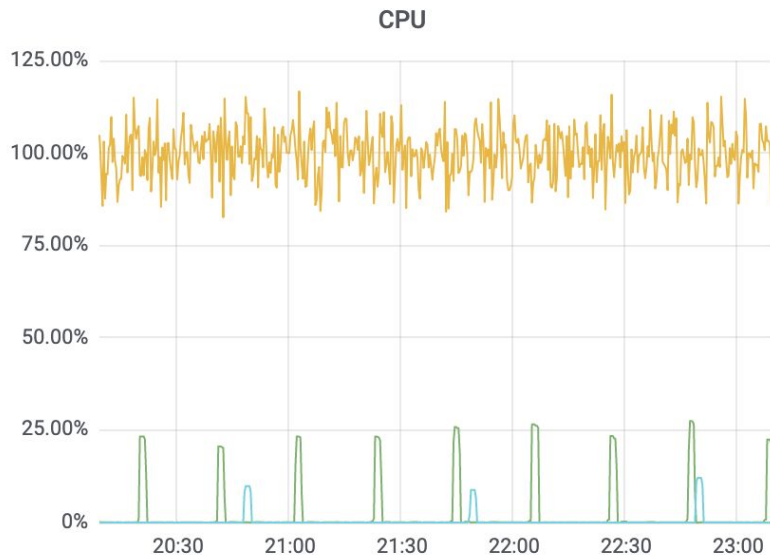
OctoRPKI - Features (1/2)

- Decodes TAL/CER/ROA/MFT
- Explore via Manifest or directory.
- RRDP support (and failover to rsync)
- Monitoring (Prometheus and JSON API which includes logs)
- Dockerizeable
- Handle stability (generate file when done)

OctoRPKI - Features (2/2)

- Full compatibility with GoRTR (including signing the JSON file)
- Server + caching options for generated file (CDN friendly)
- Configuration options
 - Disable/Enable components
 - Modes (server, one-off)
- ~5-15 minutes for a full cold-start sync

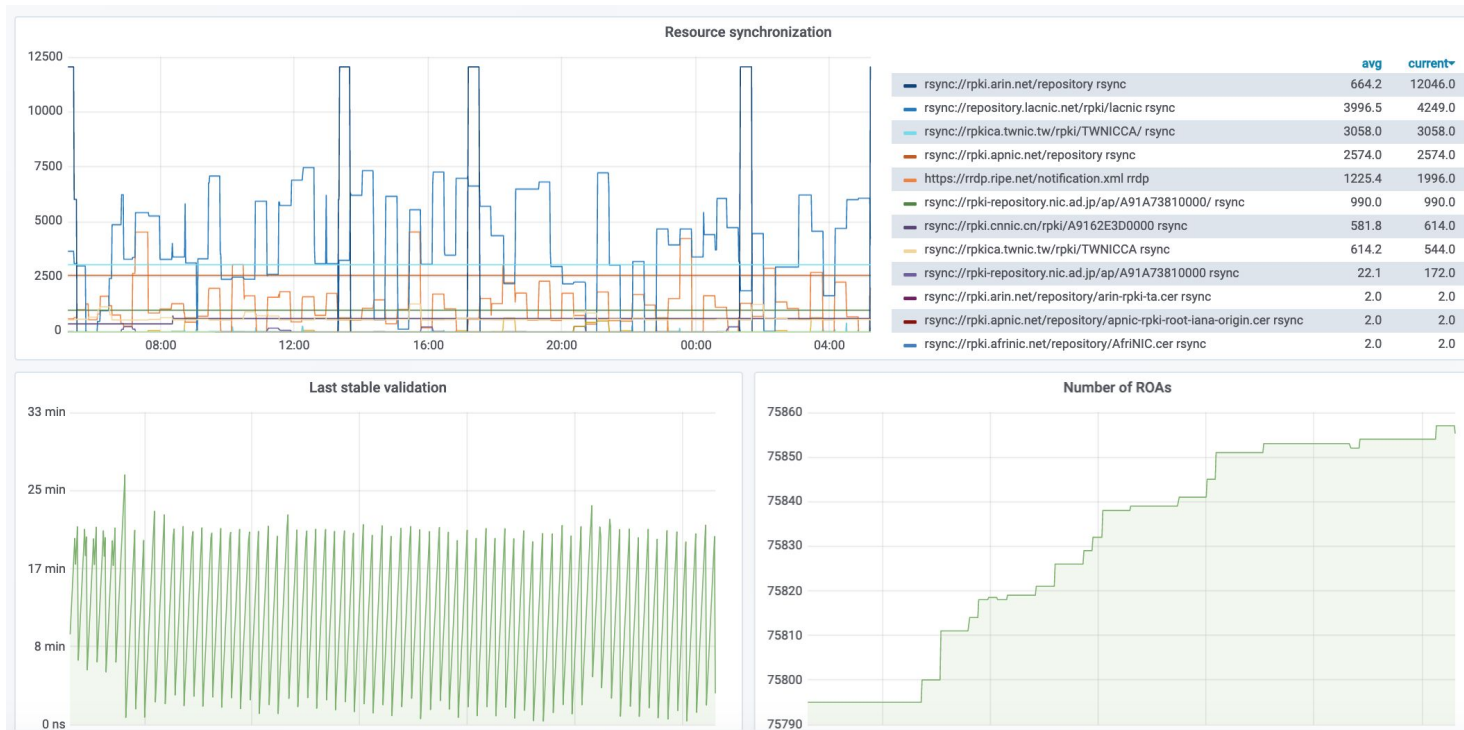
OctoRPKI - Compute footprint



	avg
rpki-benchmark-octorpki	428 MB
rpki-benchmark-ripe	906 MB
rpki-benchmark-routinator	511 MB

OctoRPKI v1.1.3
RIPE Validator v2.25
Routinator v3.3.0

Monitoring



```
1. docker
INFO[0002] Rsync sync rsync://rpki.arin.net/repository/arin-rpki-ta.cer
INFO[0002] Rsync sync rsync://rpki.ripe.net/ta/ripe-ncc-ta.cer
INFO[0004] Rsync sync rsync://repository.lacnic.net/rpki/lacnic/rta-lacnic-rpki.cer
INFO[0007] Rsync sync rsync://rpki.afrinic.net/repository/AfriNIC.cer
ERROR[0012] Error exploring file: open cache/rpki.afrinic.net/repository/04E8B0D80F4D11E0B657D8931367AE7D/62gPOPXWxxu0sQa4vQZYUBLaMb
Y.mft: no such file or directory
ERROR[0012] Error exploring file: open cache/rpki.apnic.net/repository/838DB214166511E2B38C286172FD1FF2/C5zKkN0Neo3ZmsZIX_g2EA3t6I.
mft: no such file or directory
ERROR[0012] Error exploring file: open cache/rpki.arin.net/repository/arin-rpki-ta/arin-rpki-ta.mft: no such file or directory
ERROR[0012] Error exploring file: open cache/repository.lacnic.net/rpki/lacnic/rta-lacnic-rpki.mft: no such file or directory
ERROR[0012] Error exploring file: open cache/rpki.ripe.net/repository/ripe-ncc-ta.mft: no such file or directory
INFO[0012] Still exploring. Revalidating now
INFO[0012] RRDP sync https://rrdp.apnic.net/notification.xml
INFO[0012] RRDP: Downloading root notification https://rrdp.apnic.net/notification.xml
INFO[0014] RRDP: https://rrdp.apnic.net/notification.xml Downloading snapshot at: https://rrdp.apnic.net/fa64523b-7381-4fda-9eb9-b1
233b30f503/83968/snapshot.xml
INFO[0064] RRDP sync https://rrdp.ripe.net/notification.xml
INFO[0064] RRDP: Downloading root notification https://rrdp.ripe.net/notification.xml
INFO[0064] RRDP: https://rrdp.ripe.net/notification.xml Downloading snapshot at: https://rrdp.ripe.net/8ab7553b-b124-4717-b20c-105a
da07476c/866/snapshot.xml
INFO[0177] Rsync sync rsync://repository.lacnic.net/rpki/lacnic
INFO[0241] Rsync sync rsync://rpki.arin.net/repository
INFO[0298] Rsync sync rsync://rpki.afrinic.net/repository
INFO[0309] Rsync sync rsync://rpki.apnic.net/repository
```

Validator

```
127.0.0.1:8080/infos
127.0.0.1:8080/output.json
127.0.0.1:8080/infos
{
  "uri": "rsync://repository.lacnic.net/rpki/lacnic/rta-lacnic-rpki.cer",
  "file-count": 1,
  "iteration": 1,
  "errors": 0,
  "duration": 2.9613296,
  "last-fetch": 1550342002
},
{
  "uri": "rsync://rpki.afrinic.net/repository/AfriNIC.cer",
  "file-count": 1,
  "iteration": 1,
  "errors": 0,
  "duration": 4.3905323,
  "last-fetch": 1550342006
},
{
  "uri": "https://rrdp.apnic.net/notification.xml",
  "file-count": 6734,
  "iteration": 1,
  "errors": 0,
  "duration": 52.0518939,
  "last-fetch": 15503420590,
  "rrdp-serial": 83968,
  "rrdp-sessionid": "fa64523b-7381-4fda-9eb9-b1233b30f503",
  "rrdp-last-file": "https://rrdp.apnic.net/fa64523b-7381-4fda-9eb9-b1233b30f503/83968/snapshot.xml"
},
{
  "uri": "https://rrdp.ripe.net/notification.xml",
  "file-count": 83968,
  "iteration": 1,
  "errors": 0,
  "duration": 52.0518939,
  "last-fetch": 15503420590,
  "rrdp-serial": 83968,
  "rrdp-sessionid": "fa64523b-7381-4fda-9eb9-b1233b30f503",
  "rrdp-last-file": "https://rrdp.ripe.net/8ab7553b-b124-4717-b20c-105a-da07476c/866/snapshot.xml"
}
```

API

```
$
```



```
INFO[1065] Stable state. Revalidating in 20m0s
```

ROA list

```
$ cat output.json | jq '.roas | length'
```

```
76058
```

OctoRPKI - Run it yourself

```
$ docker run -ti \  
  -p 8080:8080 \  
  -v $PWD/cache:/cache \  
  -v $PWD/tals/arin.tal:/tals/arin.tal \  
  cloudflare/octorpmi
```

Container image

Adding ARIN TAL

Use cache folder on host

Open port 8080 on host

GoRTR

```
$ docker run -ti \  
  -p 8082:8082 \  
  -v $PWD/example.pub:/example.pub \  
  cloudflare/gortr \  
    -verify.key /example.pub \  
    -cache https://YOUR_ROA_URL
```

OctoRPKI does not embed a RTR server. Modular and independence!

Fully compatible with **GoRTR**

<https://github.com/cloudflare/gortr>

Signs the prefix list to ensure a safe distribution of the file.

Can run natively on Juniper!

GoRTR

The only software to support **plaintext**, **SSH** and **TLS** as transports

Compatibility matrix

A simple comparison between software and devices. Implementations on versions may vary.

Device/software	Plaintext	TLS	SSH	Notes
RTRdump	Yes	Yes	Yes	
Juniper	Yes	No	No	
Cisco	Yes	No	Yes	Only SSH password
Alcatel	Yes	No	No	
Arista	No	No	No	
FRRouting	Yes	No	Yes	Only SSH password
Bird	Yes	No	Yes	Only SSH key
Quagga	Yes	No	No	

GoRTR without installing anything

SSH:

`rtr.rpki.cloudflare.com:8283 (user: rpki / pass: rpki)`

Plaintext:

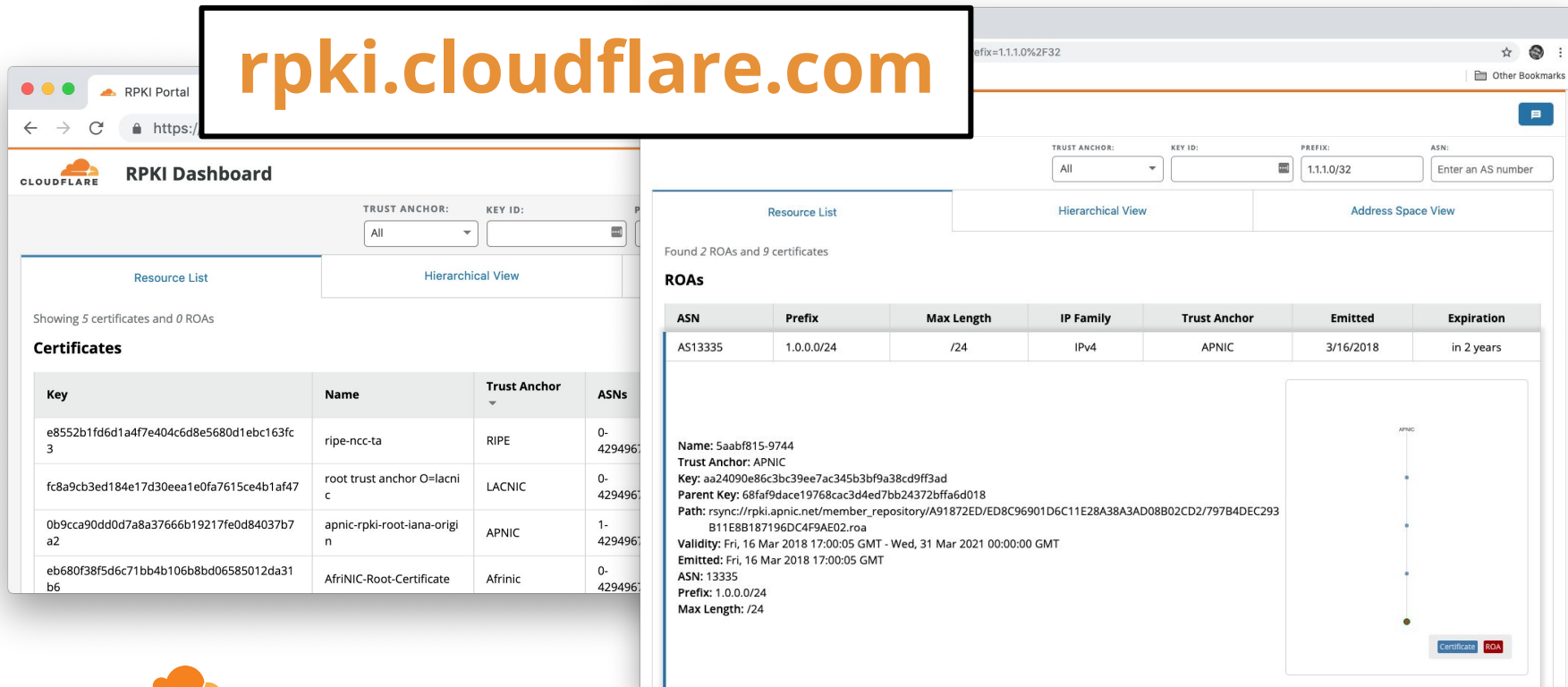
`rtr.rpki.cloudflare.com:8282`

Just configure into your router and go!

```
router bgp 65001
  rpki server 192.168.1.100
    transport tcp port 8282
  !
!
```

Cloudflare's Internal Version

rpki.cloudflare.com



The screenshot displays the Cloudflare RPKI Dashboard. The left sidebar shows the 'RPKI Dashboard' header and navigation tabs for 'Resource List' and 'Hierarchical View'. The main content area is divided into two sections: 'Certificates' and 'ROAs'.

Certificates

Key	Name	Trust Anchor	ASNs
e8552b1fd6d1a4f7e404c6d8e5680d1ebc163fc3	ripe-ncc-ta	RIPE	0-429496
fc8a9cb3ed184e17d30eea1e0fa7615ce4b1af47	root trust anchor O=Iacnic	LACNIC	0-429496
0b9cca90dd0d7a8a37666b19217fe0d84037b7a2	apnic-rpki-root-iana-origin	APNIC	1-429496
eb680f38f5d6c71bb4b106b8bd06585012da31b6	AfrinIC-Root-Certificate	AfrinIC	0-429496

ROAs

Found 2 ROAs and 9 certificates

ASN	Prefix	Max Length	IP Family	Trust Anchor	Emitted	Expiration
AS13335	1.0.0.0/24	/24	IPv4	APNIC	3/16/2018	in 2 years

Resource Details for AS13335:

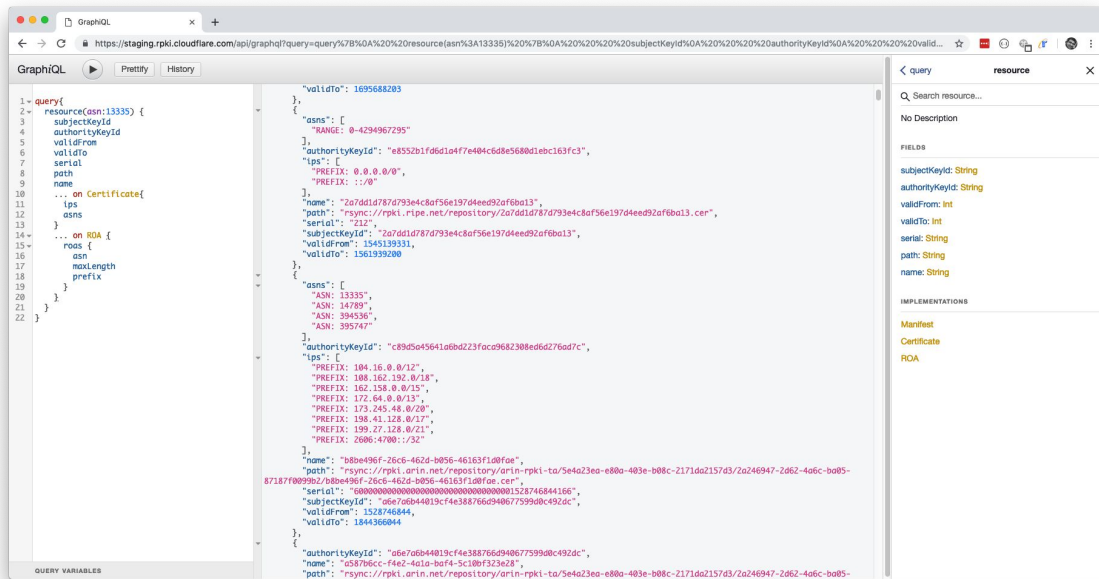
- Name: 5aabf815-9744
- Trust Anchor: APNIC
- Key: aa24090e86c3bc39ee7ac345b3bf9a38cd9ff3ad
- Parent Key: 68faf9dace19768cac3d4ed7bb24372bffa6d018
- Path: rsync://rpki.apnic.net/member_repository/A91872ED/ED8C96901D6C11E28A38A3AD08B02CD2/797B4DEC293B11E8B187196DC4F9AE02.roa
- Validity: Fri, 16 Mar 2018 17:00:05 GMT - Wed, 31 Mar 2021 00:00:00 GMT
- Emitted: Fri, 16 Mar 2018 17:00:05 GMT
- ASN: 13335
- Prefix: 1.0.0.0/24
- Max Length: /24

A diagram on the right shows the trust path from the APNIC Trust Anchor to the specific ROA, with a legend for 'Certificate' (blue) and 'ROA' (red).

Cloudflare's Internal Version

Provides <https://rpki.cloudflare.com/rpki.json>

Also a **GraphQL** API for the dashboard

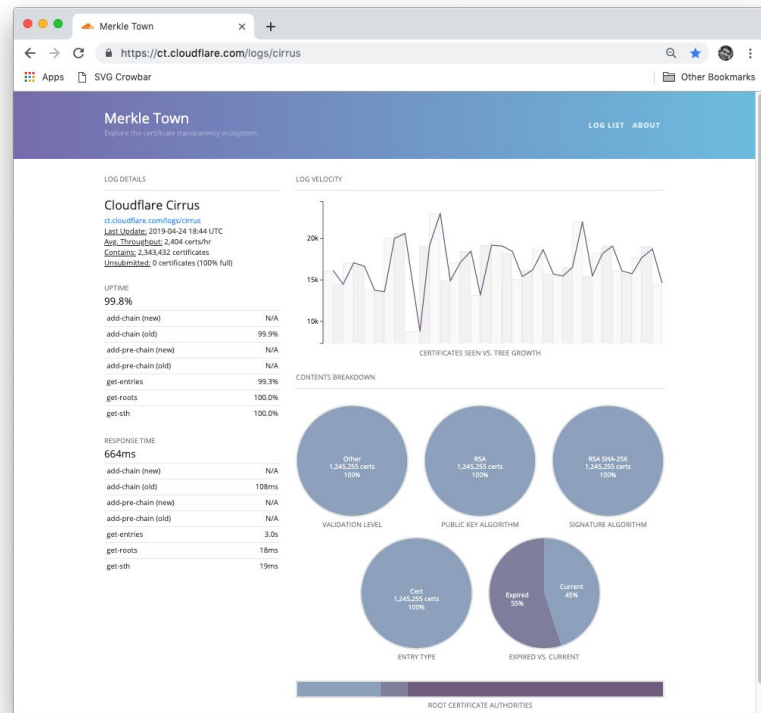


Certificate Transparency

Historical records of certificates

Contains a chain (root → ROA)

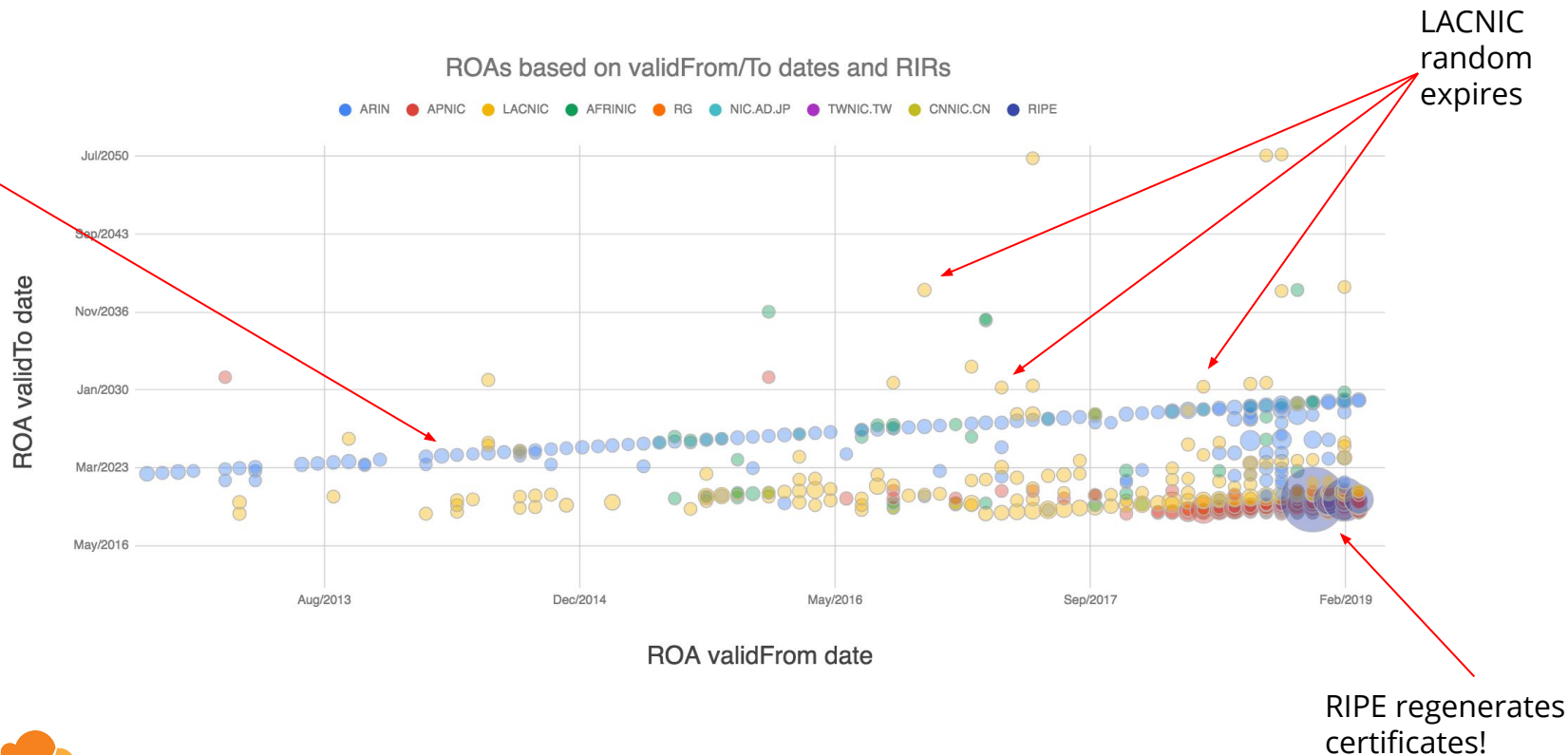
Sent by our validator



Other data

Other data - so how fresh are those ROAs?

ARIN uses
ten year
expire



Future projects

Future projects or ideas

Certificate encoder, ASPA.

More toolings and visualizations around RPKI (BGP collection):

- Integration in our portal <https://peering.cloudflare.com/> (*ask for your free access*)

The Probing Project

Could we probe the entire Internet ^[1] to see who is doing validation?

Who?

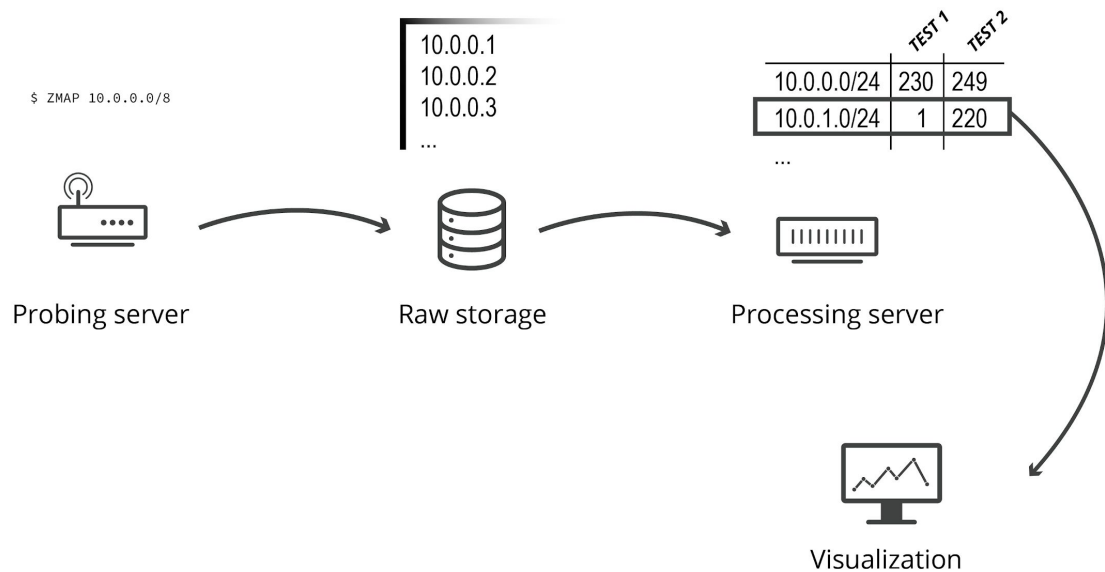
Involved in this project and special thanks:

- Job Snijders: NTT, NLNOG
- Jérôme Fleury, Vasco Asturiano: Cloudflare

Methodology

1. Run two tests with zmap (~2hrs/test):
 - one test with behind an RPKI valid prefix and ;
 - one test behind an RPKI invalid prefix
2. Sum the IPs that replied by range for each test
3. Visualize the ratio of replies between the two tests per prefix

Methodology



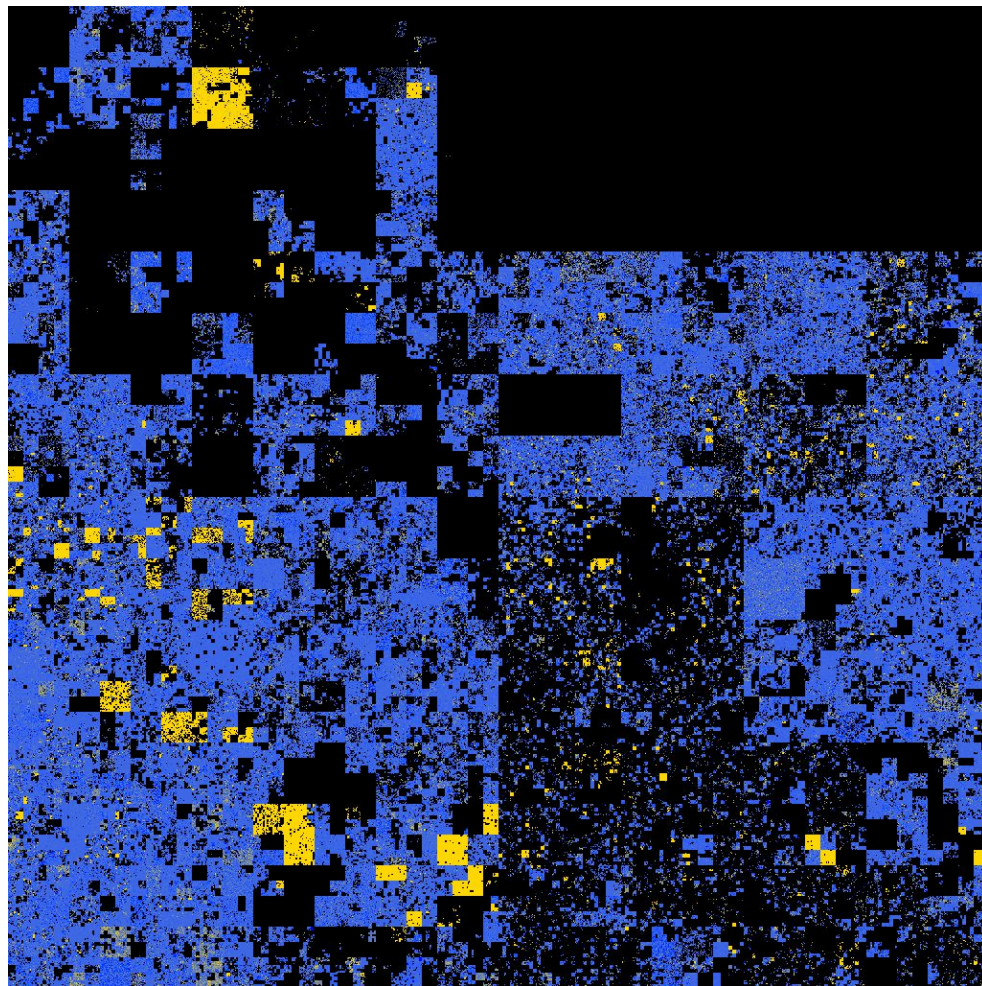
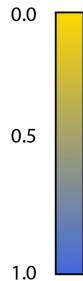
Hilbert

$1024 \times 1024 \rightarrow 1 \text{ pixel per } /20$

Ratio

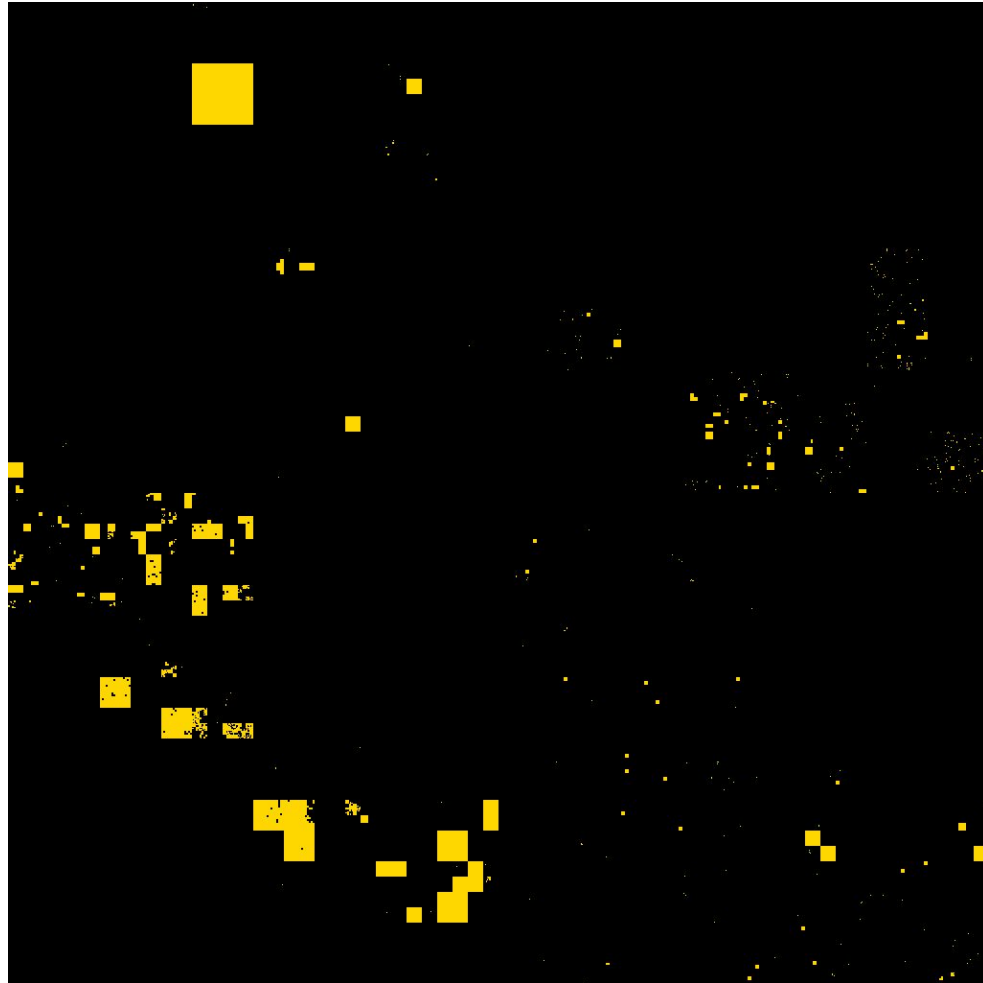
$$\frac{\text{Received with source as invalid prefix}}{\text{Received with source as valid prefix}}$$

Good!



Hilbert

AT&T
Announced IP addresses



Learnings

AT&T has ~400k /24.

Easiest to visualize because big blocks.

Could we **automatically** identify networks?

Going from "0 to 1" to "0 and 1" (0 good and 1 bad)

Machine learning

Simplify the identification of “validating prefixes”.

Requires **training** instead of defining specific ratios.

*Little particles of artificial
intelligence*

Your network



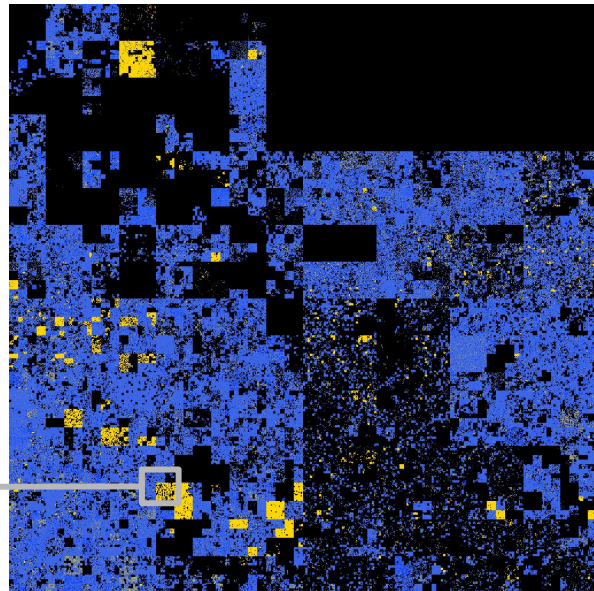
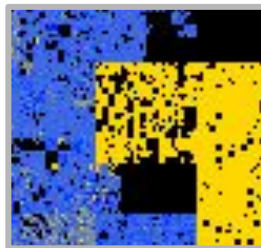
<https://www.pexels.com/photo/man-pouring-condiment-on-plate-of-food-2494658/>

Training

Take an easy sample to manually identify.

~6,500 /20 prefixes

Assisted by script: (with ASN mapping)

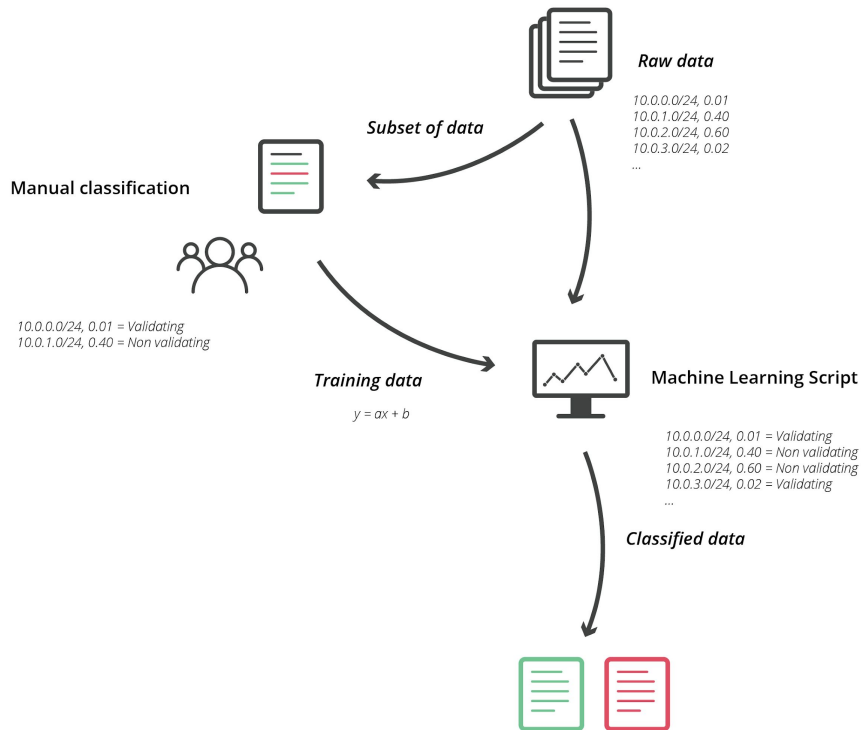


```
$ ./create_training.py \  
  --input stripped-data.csv \  
  --asn-valids 7018 \  
  --asn-invalids 701 \  
  --subnets 96.0.0.0/8 97.0.0.0/8 98.0.0.0/8 \  
  --output training_data.csv
```

Classification

The SVM model will take a few minutes to run on 433k /20's.

Model improvements: adds more variables (proximity to other low-ratio prefixes...).



Results

Detected validating:

/20 detected: 28,199

/24 detected: 320k

less than AT&T total prefixes assumed validating:

not everything is responding to ICMP

Work in progress!

Recent Leaks And Conclusions

Summary of Amazon Route Hijack

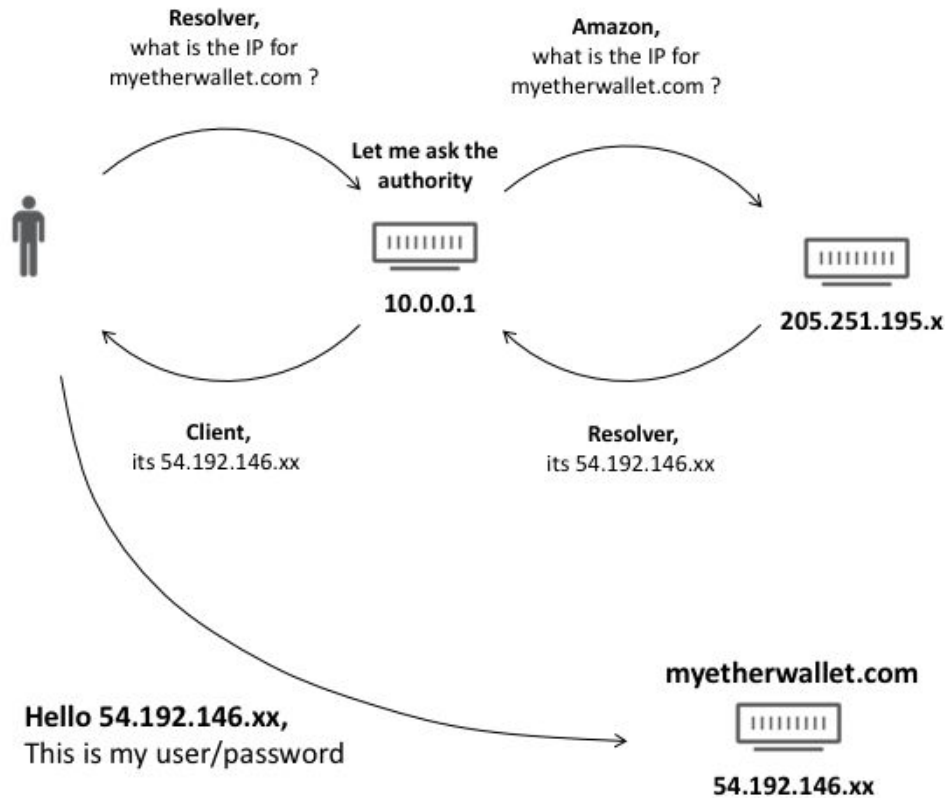
An attacker announces Amazon Authority DNS prefixes.

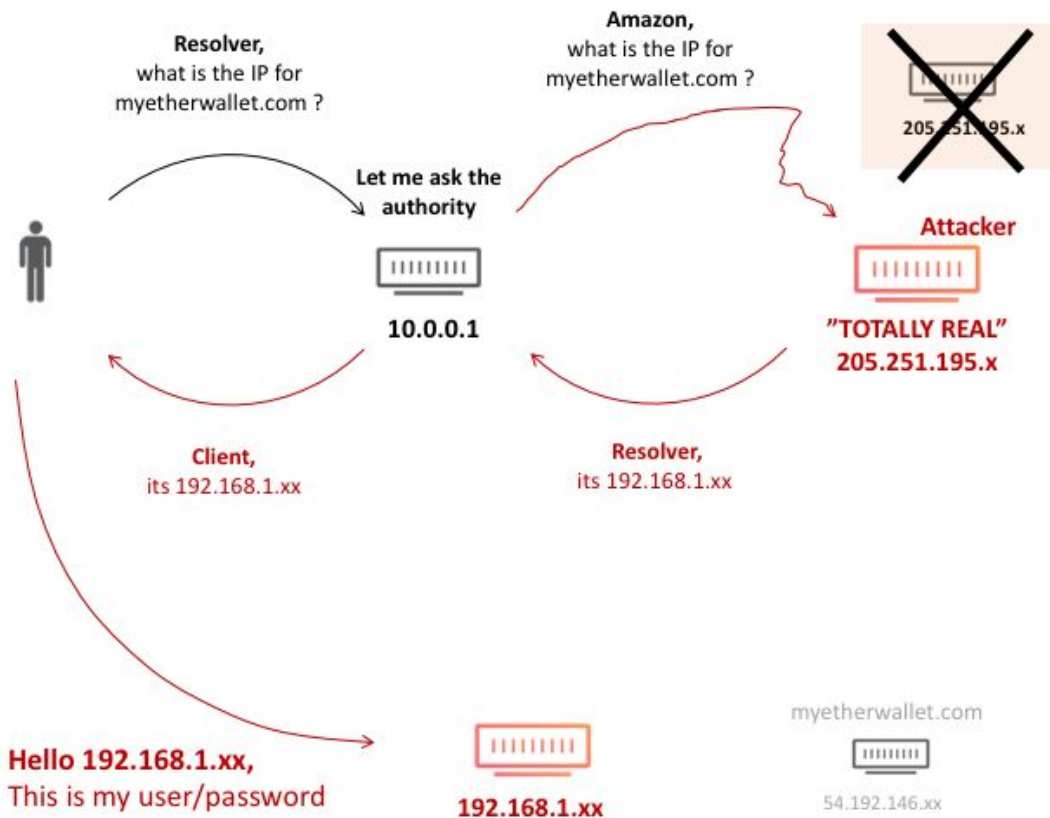
Cloudflare and Google accept them in specific locations.

Cloudflare and Google DNS resolvers use this route when clients request the website, the attacker's server is returned.

The server has a phishing website for the client.

Attacker gather credentials and steals Bitcoins.





Summary of Amazon Route Hijack

Amazon did not have signed routes

Cloudflare did not do RPKI validation + route filtering

If RPKI was deployed:

Route would have been rejected because wrong origin

Summary of Verizon Route Leak

A company has two Internet accesses: Verizon and another ISP.

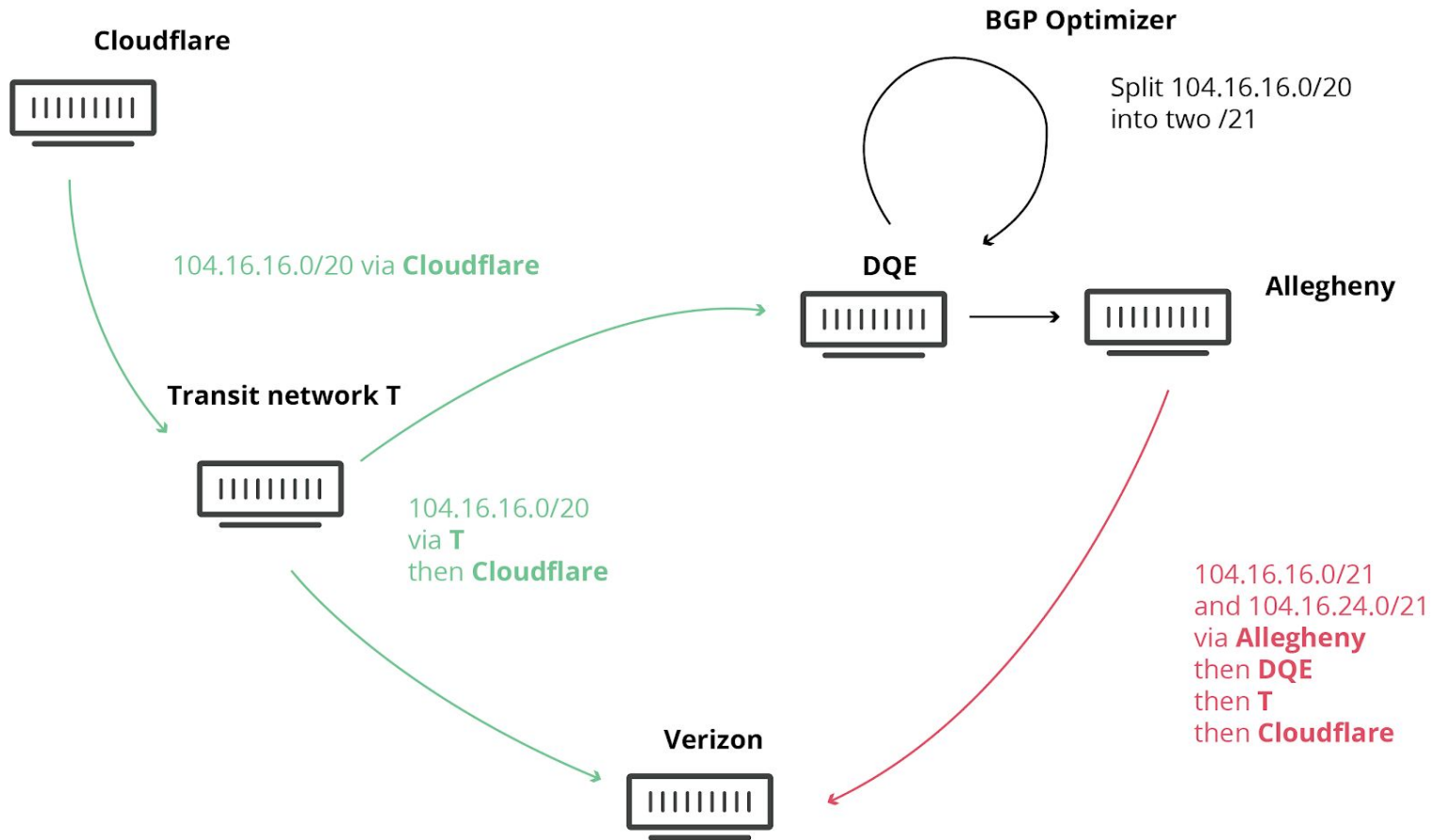
The ISP has a BGP optimizer which feeds more-specific routes.

Unfortunately, the ISP sends the routes to the company which end up being sent to Verizon.

Verizon did not filter them and re-announces them to its peers and clients.

Cloudflare loses traffic.





Summary of Verizon Route leak

Cloudflare had signed routes.

Verizon did not filter. Many networks accepted the leak.

Cloudflare filtering routes did not matter here.

If basic filtering was deployed:

Peering sessions would have been removed when going above prefix threshold.

AS-Path filtering could have avoided accepting routes.

If RPKI was deployed:

Routes would have been rejected because wrong length.

What we learned

RPKI will not be the solution to everything. But in our stories...

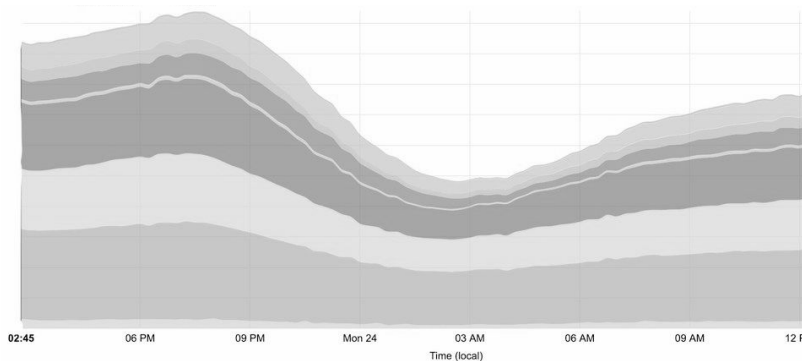
Filtering solves Amazon being hijacked

Signing helps your network not being leaked

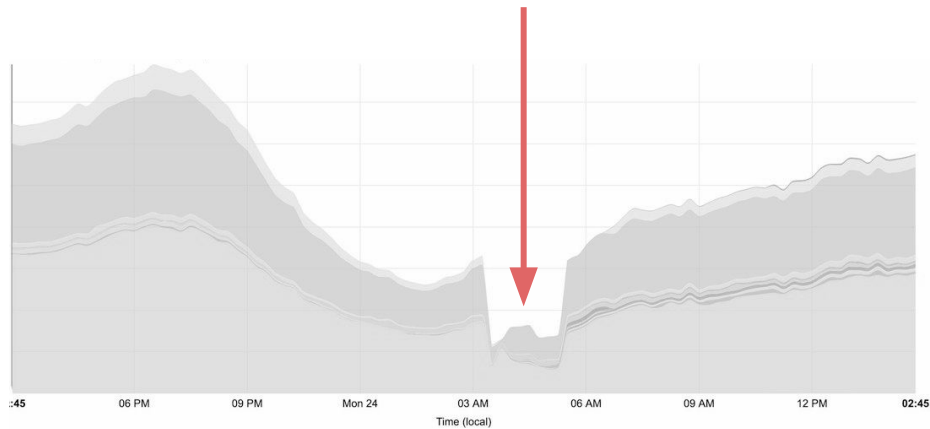
Deploy RPKI now

Because tomorrow is already too late

With filtering



Without filtering



Thank you

Questions?

louis@cloudflare.com
[@lpoinsig](https://twitter.com/lpoinsig) (twitter)

