

# Why is Securing the Internet's Routing System so Damn Difficult?

Geoff Huston  
APNIC  
February 2019

Usual BGP Disaster-Porn Clips

# Pakistan hijacks YouTube

Research // Feb 24, 2008 // Dyn Guest Blogs

Yawn

Late in the (UTC) day on 24 February 2008, Pakistan Telecom (AS 17557) began advertising a small part of YouTube's (AS 36561) assigned network. This story is almost as old as BGP. Old hands will recognize this as, fundamentally, the same problem as the infamous AS 7007 from 1997, a more recent ConEd mistake of early 2006 and even TTNet's Christmas Eve gift 2004.

FFS! No  
better 10  
years later

THE ACCIDENTAL LEAK —

# Google goes down after major BGP mishap routes traffic through China

Google says it doesn't believe leak was malicious despite suspicious appearances.

DAN GOODIN - 11/13/2018, 6:25 PM

Google lost control of several million of its IP addresses for more than an hour on Monday in an event that intermittently made its search and other services unavailable to many users and also caused problems for Spotify and other Google cloud customers. While Google said it had no reason to believe the mishap was a malicious hijacking attempt, the leak appeared suspicious to many, in part because it misdirected traffic to China Telecom, the Chinese government-owned provider that was recently caught [improperly routing traffic](#) belonging to a raft of Western carriers though mainland China.

The leak started at 21:13 UTC when [MainOne Cable Company](#), a small ISP in Lagos, Nigeria, suddenly updated tables in the Internet's global routing system to improperly declare that its [autonomous system 37282](#) was the proper path to reach [212 IP prefixes belonging to Google](#). Within minutes, China Telecom improperly accepted the route and announced it worldwide. The move by China Telecom, aka AS4809, in turn caused Russia-based [Transtelecom](#), aka AS20485, and other large service providers to also follow the route.



#### FURTHER READING

[Strange snafu misroutes domestic US Internet traffic through China Telecom](#)

According to [BGPmon on Twitter](#), the redirections came in five distinct waves over a 74-minute period. The redirected IP ranges transmitted some of Google's most sensitive communications, including the company's [corporate WAN infrastructure](#) and the [Google VPN](#). [This graphic](#) from regional Internet registry RIPE NCC shows how the domino effect played out over a two-hour span. The image below shows an abbreviated version of those events.

more

THE ACCIDENTAL LEAK —

# Google goes down due to BGP mishap routes

Google says it doesn't believe leak

DAN GOODIN - 11/13/2018, 6:25 PM

Google lost control of several million of IP addresses that intermittently made its search engine unavailable. The reason to believe the mishap was a malicious attack is that many, in part because it misdirected traffic to a local provider that was recently caught **improperly routing traffic** through mainland China.

The leak started at 21:13 UTC when **Mainland Company**, a small ISP in Lagos, Nigeria, sent an update to the Internet's global routing system to improperly declare that its **autonomous system 37282** was the proper path to reach **prefixes belonging to Google**. Within minutes, the move by **China Telecom**, aka AS20485, and other large

According to **BGPmon on Twitter**, the redirection period. The redirected IP ranges transmitted some of Google's most sensitive communications, including the company's **corporate WAN infrastructure** and the **Google VPN**. **This graphic** from regional Internet registry RIPE NCC shows how the domino effect played out over a two-hour span. The image below shows an abbreviated version of those events.

## Today's BGP leak in Brazil

Posted by Andree Toonk - October 21, 2017 - News and Updates - No Comments

Earlier today several people noticed network reachability problems for networks such as Twitter, Google and others. The root cause turned out to be another BGP mishap.

```

#BGP 40000 172.217.10.100
#Start: 2017-10-21T13:26:55+0200
#End: 31min@0.0s
# 1 AS5777 16.235.224.152 0.0% 0 10 19 71.9 37.6 56.7 71.9 9.3 56.0 19.8 0 1 16.2 22.4
# 2 AS62292 217.117.148.217 0.0% 0 10 19 94.8 47.8 26.4 40.2 0.2 54.2 7.7 7 16.7 51.2
# 3 AS717 844-nc-net 1593.189.137.1753 0.0% 0 10 19 94.8 47.8 26.4 40.2 0.2 54.2 7.7 7 16.7 51.2
# 4 AS6839 100g11-1-core1.via1.he.net (184.185.213.249) 10.0% 1 10 19 160.0 144.7 161.4 6.2 152.0 19.4 6 1 19.6 52.6
# 5 AS6839 100g11-2-core1.pac2.he.net (184.185.213.173) 0.0% 0 10 19 160.0 144.7 161.4 6.2 152.0 19.4 6 1 19.6 52.6
# 6 AS6839 100g10-2-core1.asn1.he.net (184.185.213.173) 0.0% 0 10 19 160.0 144.7 161.4 6.2 152.0 19.4 6 1 19.6 52.6
# 7 AS6839 100g10-1-core1.asn1.he.net (184.185.213.201) 0.0% 0 10 19 160.0 144.7 161.4 6.2 152.0 19.4 6 1 19.6 52.6
# 8 AS6839 100g4-1-core1.via1.he.net (184.185.213.201) 0.0% 0 10 19 174.5 168.1 172.2 181.7 7.0 172.1 8.7 6 1 21.0 59.1
# 9 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 10 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 11 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 12 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 13 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 14 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 15 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 16 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 17 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 18 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 19 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 20 AS5777 777 100.0 10.0 0 19 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
# 21 AS52328 200.16.70.280 80.0% 8 2 19 495.2 495.2 495.1 500.3 2.9 498.2 4.1 2 1 4.1 4.1
# 22 AS717 45.6.53.73 80.0% 8 1 19 930.0 930.0 930.0 930.0 0.0 930.0 0.0 0 0.0 0.0
# 23 AS5777 777 100.0 4.0 0 4 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0

```



Some Google services seem to have been hijacked for roughly 15 minutes. Seen anything? @atoonk @bgpmon @bgpstream MTR: [xor.meo.ws/P0SYOU7j-4Ftj...](https://xor.meo.ws/P0SYOU7j-4Ftj...)

26 10:32 PM - Oct 21, 2017

23 people are talking about this

Between 11:09 and 11:27 UTC traffic for many large CDN was rerouted through Brazil. Below an example for the Internet's most famous prefix 8.8.8.0/24 (Google DNS) At 2017-10-21 11:09:59 UTC, AS33362, an US based ISP saw the path towards Google's 8.8.8.0/24 like this:

```
33362 6939 16735 263361 15169
```

This shows the US based network AS33362, would have sent traffic to Google via 6939 (HE) to 16735 (Algar Telecom, Brazil), to 263361 infovale telecom which would have tried to delivered it to Google. The successful delivery of packets would have been unlikely, typically due to congestion which would have been the result of the increase in attracted traffic or an ACL blocking the unexpected traffic.



Why do we keep seeing these headlines?

### Popular Destinations rerouted to Russia

Posted by Andree Toonk - December 12, 2017 - Hijack - No Comments

# CommsWire

Essential daily reading for the communications industry executive

An iTWire publication www.itwire.com Editor: Stan Beer Friday 16 November 2018

## TELSTRA ROUTING ERROR TAKES DOWN INTERNET



include region span.

What makes this incident suspicious is the prefixes that were affected are all high profile destinations, as well as several more specific prefixes that aren't normally seen on the Internet. This means that this isn't a simple leak, but someone is intentionally inserting these more specific prefixes, possibly with the intent to attract traffic.

This graphic from over a two-hour

as Twitter,

il. Below an

gle's

39 (HE) to delivered it to ACL

# Degrees of Difficulty

Why are some issues so challenging to solve, while others seem to be effortless?

Why was the IPv4 Internet an unintended runaway success in the 90's, yet IPv6 has been a protracted exercise in industry-wide indecision?

# Internet Successes

- IPv4 (and datagram packet switching)
- Network Address Translators (perversely!)
- TCP evolution and adaptation
- DNS
- Content Distribution Systems
- Streaming



# Success Factors

- Piecemeal deployment without the requirement for central orchestration
- Competitive advantages to early adopters
- Economies of scale as adopter numbers increase
- Alignment of common benefit with individual benefit

# Internet ~~Non-Successes~~

Failures!

- SPAM
- DDOS defence
- BCP 38 deployment
- Secure end systems
- Secure networks
- Internet of Things
- IPv6 adoption (so far!)

# Failure Factors

- Need for orchestrated actions (flag days)
- Technologies that require universal or near universal adoption
- Where there are common benefits but not necessarily individual benefits
- Where there is no clear early adopter advantage

# What makes a problem "hard"?

It might be **technically challenging**: While we understand what we might want that does not mean we know how to construct a solution

It might be **economically hard**: The costs of a solution are not directly borne by the potential beneficiaries of deploying the solution

It might be **motivated by risk mitigation**: We are notorious for undervaluing future risk!

# Why is Securing Routing so Hard?

- Because no single entity is in charge
- Because we can't audit BGP, as we have no standard reference route set to compare with
- Because we can't arbitrate between conflicting BGP information (because there is no standard reference point)
- Because there are no credentials that allow a BGP update to be compared against the original route injection (because BGP is a hop-by-hop protocol)
- Because BGP is based on opaque local decisions

**Why should we worry?**

Because it's just too easy to  
be bad in routing!

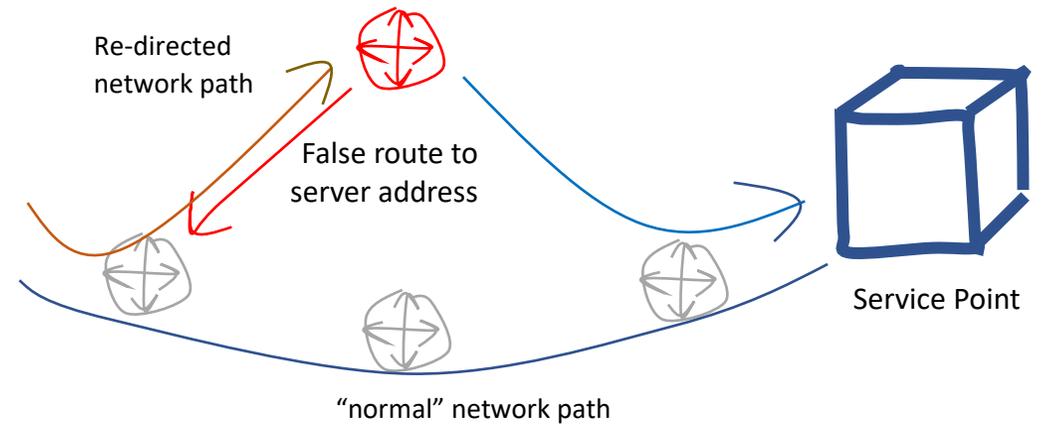


User traffic gets diverted enabling a Man-in-the-Middle attack on a service

What's the risk?



User





What's the  
risk?

## DOS Attack

Divert the traffic to a sinkhole

- Deny users access to the site
- Crude, but effective!



What's the  
risk?

## DNS Attacks

Divert DNS traffic to fake DNS servers and provide fake answers

- Very few domains are DNSSEC-signed and not enough resolvers perform DNSSEC validation
- So the faked answer can pass unchallenged



What's the  
risk?

## Server Attacks

Divert TCP traffic to fake servers and provide fake answers

- Collect user credentials while shadowing the actual site

# An attack vector on HTTPS...

- Let's say you can find an online trusted CA
  - that uses the DNS as proof-of-possession of a DNS name in order to mint a domain name certificate
  - And the DNS name is not DNSSEC protected
- You can mint a fake domain name certificate by:
  - Mount a routing attack on the DNS infrastructure with a fake DNS responder
  - Answer everything correctly except for \*.victim ACME DNS challenge from the CA
  - And for the \*.victim challenge queries respond with your own answer
  - Which means you can answer the CA's DNS challenge
- Now you have a trusted fake domain name certificate
- You are now able to pull off a MITM attack on a TLS 'protected' service

# I want a Pony Routing Security wish list

1. Identify whether an address is “bogus” or not
2. Assure that the address holder has given their permission for an address to be announced into the routing system
3. Identify which AS(s) have been given this permission
4. Identify if the AS Path is consistent with the ‘correct’ operation of BGP
5. Identify if the AS Path is consistent with the routing policies of the each of the Ases
6. Identify when routing information is being ‘incorrectly’ withheld

The saga so far...

# Internet Route Registries

- First used in the early 1990's as the Route Arbiter Database (RADB) as part of the NSFNET program
- Describes route origination and inter-AS routing policies
- An explicit declaration of intent in routing
- Route Registries can be used to filter BGP announcements, filtering out route advertisements that are not described in the route registry
  - Primary value in preventing neighbor route leaks
  - Can be used to prevent hijacks

# Route Registry Issues

- Poor Authority Model (or the complete lack of one in many cases!)
  - How can a user know that a RR entry is genuine and current?
  - How can a user know that a RR entry is maintained by an entity who is the authoritative “owner” of an IP address or ASN?
  - How can a user tell the difference between a current RR entry and a lapsed historical RR entry?
- Too many Route Registries
  - If two different RRs contain conflicting information, what are users meant to do?
- Incomplete Data
  - If a route is not described in a Route Registry is it just the registry that is missing data or is the route itself invalid?
- Scaling issues
  - No realistic way to apply IRR filters to upstreams
- RPSL got too geeky!
  - The Route Policy Language used by Route Registries got overly expressive and complex

# What's missing with RRs?

- If we want to improve the usefulness of route registries we probably need a **robust authority model**

## **How about Digital Signatures?**

- The signatures can provide currency and authenticity
- The authority model can allow RR entries to be seen as explicit authorities or permissions from address holders to network operators and from network operators to other networks

# X.509 Public Key Certificates for IP addresses and AS Numbers

- An X.509 Public key certificate that includes a set of IP addresses and AS numbers
- If a certificate can be validated against a trust anchor then it indicates that:
  - The IP addresses and/or AS numbers have been validly allocated
  - The holder of the subject key pair is the current holder of the IP addresses and/or AS numbers
  - Attestations validly signed using this key can be considered as genuine authorities that cannot be repudiated
- This is the foundation of the current work in routing security

# Route Origination Authority

- An address holder can convey a 'permission' for an AS to originate a BGP route for the address by signing a permission authority (ROA) using a signing key associated with a valid public key address certificate
- This authority:
  - can be validated by any interested party
  - is dated, so currency is known
  - cannot be repudiated

# If we all used ROAs then:

- ✓ 1. Identify whether an address is “bogus” or not
- ✓ 2. Assure that the address holder has given their permission for an address to be announced into the routing system
- ✓ 3. Identify which AS(s) have been given this permission
4. Identify if the AS Path is consistent with the ‘correct’ operation of BGP
5. Identify if the AS Path is consistent with the routing policies of the each of the Ases
6. Identify when routing information is being ‘incorrectly’ withheld

Is 3 out of 6 good enough to  
get a pony?

**NO!**

- The hijack can reproduce the origin and if the ROA is sloppy then it can use a more specific
- Even if the ROA is tight the conflicting routes can still support a desired attack profile

# From ROAs to a fully secure BGP

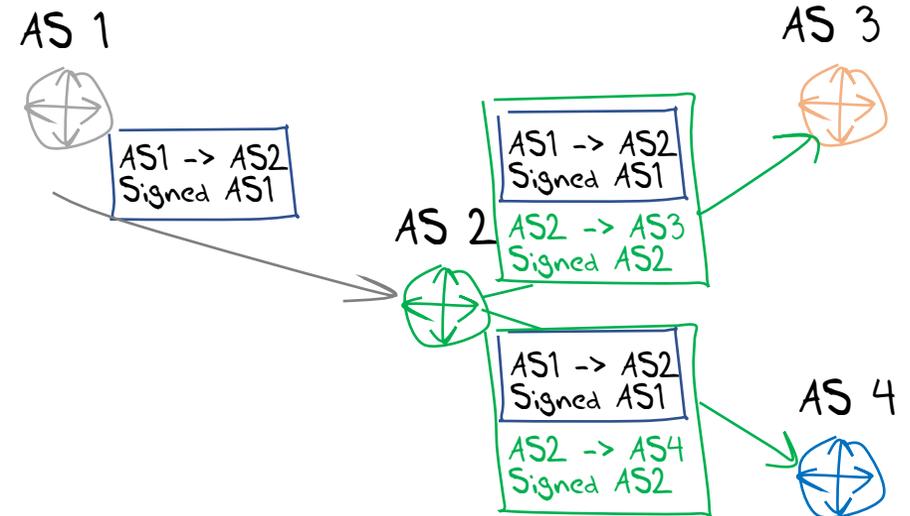
ROAs are good, but probably not enough to stop a determined routing attacker

- The attacker simply needs to replicate the BGP origination in the AS path to be accepted as “good”

So we really need to secure the BGP AS Path as well

We can do this with RPKI certs!

- Every eBGP speaker has a key that is certified by the AS
- When an update is passed to a neighbor AS, the router signs across the existing AS Path signature and the neighbor AS
- A BGPSEC speaker validates a received update by checking that
  - there is a current ROA to describe the address and origin AS
  - The received AS Path can be validated as a sequence of sign-over-sign operations by the AS keys



# But ASPath protection is hard...

- BGPSEC cannot cope with partial adoption
  - It cannot jump across non-participating networks
- It has a high crypto overhead for session restarts
- It does not define how to promulgate the collection of certificates required to validate the digital signatures
- It does not necessarily identify and prevent route leaks
- Which means that BGPSEC is not looking like its going to be deployed everywhere
  - Which means that there is little value in deploying it anywhere

# What's going wrong?

The economics of this situation work against it

- There are inadequate commercial drivers to undertake extensive informed route monitoring that would enable hijack suppression at source
- Probably because integrity of common infrastructure is everyone's problem which in turn quickly becomes nobody's problem
- And we have no 'forcing' authority to compel network operators

# Where to from here?

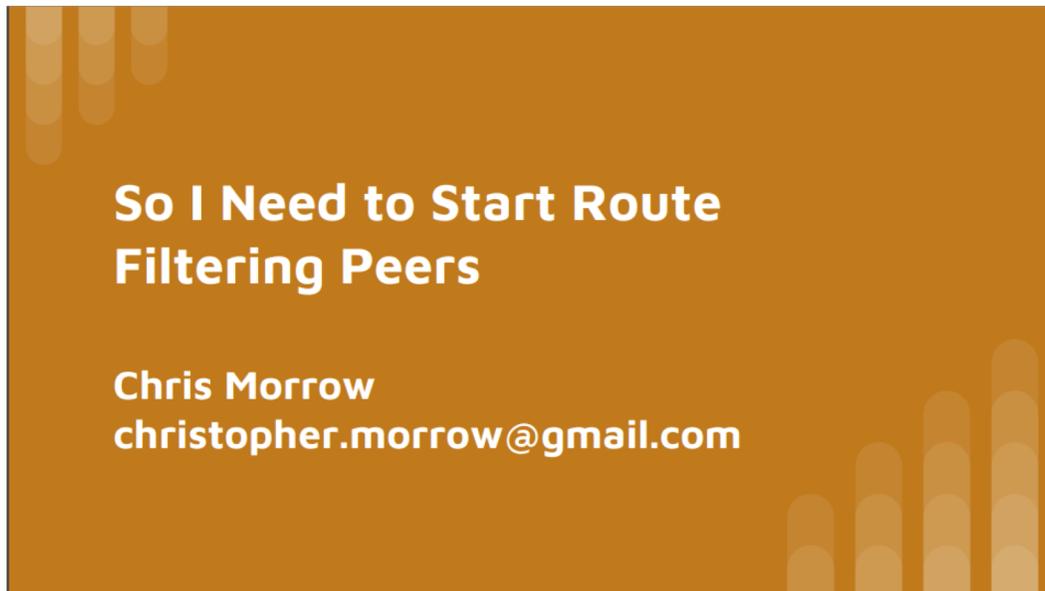
- We are pretty convinced about the value of RPKI certificates and digital signatures
  - Because we really have nothing better to offer in their place
- But the AS Path protection elements of BGPSEC are a critical problem!
- In the IETF we are working on approaches that address the issues with BGPSEC and AS Path protection
  - But that effort could take years
  - And there is no guarantee of success!

# Where are we heading?

- The problem is not going to go away
- So we probably need to look at other ways to secure the propagation of routing information:
  - What if we decoupled origination, topology and policy validation?
  - Will open market disciplines lead us to a secure Internet environment or are we necessarily looking at regulatory imposts to force universal adoption?
  - What could we gain by using deliberate efforts at asymmetric partial adoption?
    - What's more important in routing security: client routes or server routes?
    - i.e. should we concentrate on IXPs and CDN routes as points of active route policing?

# In the meantime...

Over at the Google ranch:



**So I Need to Start Route  
Filtering Peers**

**Chris Morrow**  
christopher.morrow@gmail.com



## Conclusion

Goal is to start marking routes based on filter inclusion / exclusion by 01/2019

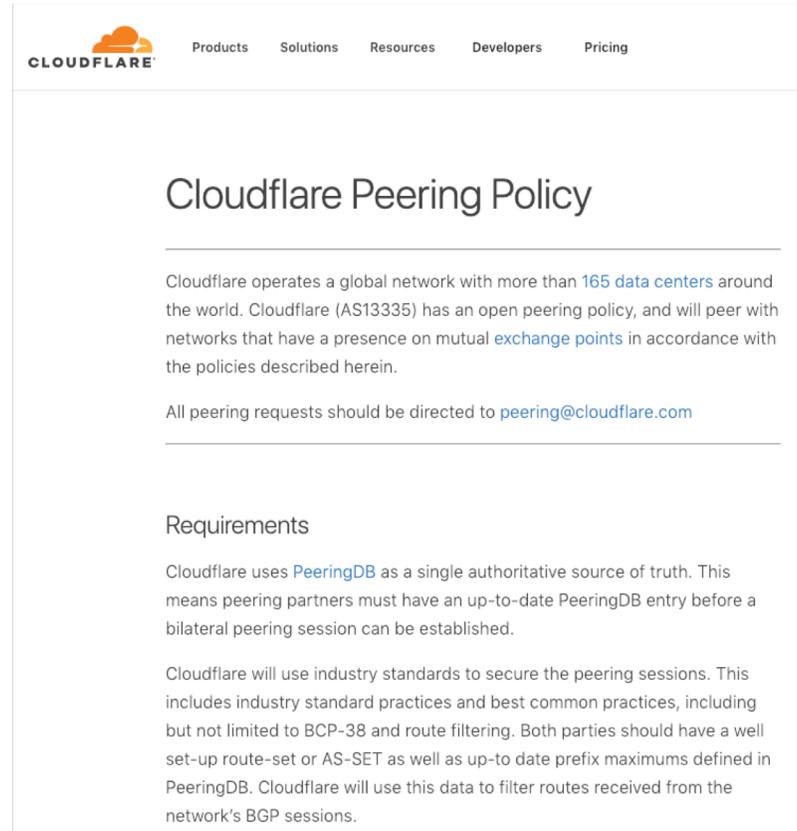
Reject/Drop by 03/2019

October 2018

[https://pc.nanog.org/static/published/meetings/NANOG74/1760/20181003\\_Tzvetanov\\_Security\\_Track\\_Bgp\\_v1.pdf](https://pc.nanog.org/static/published/meetings/NANOG74/1760/20181003_Tzvetanov_Security_Track_Bgp_v1.pdf)

# In the meantime...

Over at the Cloudflare ranch:



The screenshot shows the Cloudflare website's navigation bar with the logo and links for Products, Solutions, Resources, Developers, and Pricing. The main content area is titled "Cloudflare Peering Policy" and contains the following text:

Cloudflare operates a global network with more than [165 data centers](#) around the world. Cloudflare (AS13335) has an open peering policy, and will peer with networks that have a presence on mutual [exchange points](#) in accordance with the policies described herein.

All peering requests should be directed to [peering@cloudflare.com](mailto:peering@cloudflare.com)

---

### Requirements

Cloudflare uses [PeeringDB](#) as a single authoritative source of truth. This means peering partners must have an up-to-date PeeringDB entry before a bilateral peering session can be established.

Cloudflare will use industry standards to secure the peering sessions. This includes industry standard practices and best common practices, including but not limited to BCP-38 and route filtering. Both parties should have a well set-up route-set or AS-SET as well as up-to date prefix maximums defined in PeeringDB. Cloudflare will use this data to filter routes received from the network's BGP sessions.

<https://www.cloudflare.com/peering-policy/>

# What can you do today?

“Don’t let the perfect be the enemy of the good!”

- You could just wait for a complete routing security framework to be invented
- Or you could do something practical right now that might be helpful

# What can you do today?

You might want to take some steps to make routing attacks easier to detect and easier to deflect

- BCP38 filters can help
  - UDP DOS attacks are very common
- Generating ROAs can help
  - Maybe they won't help a lot today, but as more networks filter on ROAs then they will be more effective to protect against simple address hijacking
- Route Registry objects can help
  - [www.irr.net](http://www.irr.net)
  - Again this is not a complete answer, but its better than nothing
- You should really should filter your customers
  - Filter customer routing updates according to BCP38, ROAs and IRR profiles
- Consider signing up to MANRS
  - <https://www.manrs.org/> (Even spending a few minutes thinking about routing security is better than not thinking about it at all)
- DNSSEC-sign your domain name
- Validate DNS responses

**Should have bought me  
that pony**



**Thanks!**

Questions?

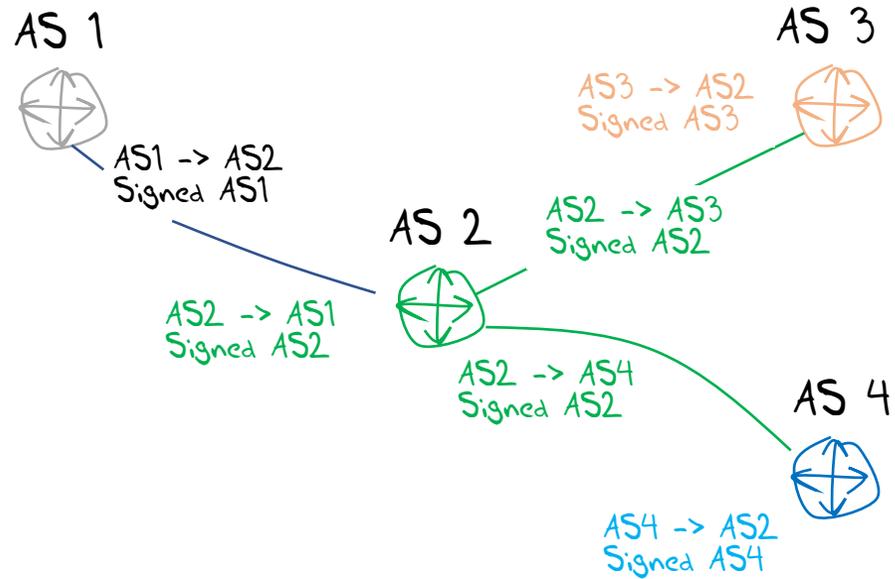
# Additional Material

1. soBGP
2. General Comments
3. Opinions

# 1. soBGP: an alternative to BGPSEC

- Instead of the high overhead of AS Path validation we can look at secure origin BGP (soBGP) from 2003
- soBGP looked at the AS Path as a topology vector composed of a number of paired AS adjacencies
  - An AS publishes a signed adjacency attestation for all of its neighbors
  - If a signing AS appeared in an AS Path then its neighbors in the AS Path must also be described in the adjacency attestation
- This replaces strict AS Path **Validation** with AS Path **Plausibility**

# 1, soBGP and AS Adjacencies



## AS Path Processing using AS Adjacency 'hints'

AS1 -> AS2 -> AS3	plausible
AS1 -> AS3 -> AS2	implausible
AS1 -> AS2 -> AS3 -> AS4	implausible

# 1. s0BGP compared to BGPSEC

- Lower crypto overhead
- Can be used in scenarios of partial adoption
- Does not prevent a network from learning false information, but prevents a network being used in a falsified AS path
  - Unless you also include the AS's peers
  - And so on
  - Incremental deployment generates incremental benefit
- Can include directionality in the AS adjacency attestation
  - As a simple “policy” filter

## 2. Generic Concerns over PKIs

Is a *trust hierarchy* the best approach to use?

- The concern here is **concentration of vulnerability**

If validation of routing information is dependent on the availability and validity of a single root trust anchor then what happens when this single digital artifact is attacked?

- But is there a viable alternative approach?

Can you successfully incorporate robust diversity of authentication of security credentials into a supposedly highly resilient secure trust framework?

This is a very challenging question about the nature of trust in a diverse networked environment!

Web trust – 1,500 CAs vs DNSSEC trust – 1 key  
which is ‘better’?

## 2. Generic Concerns over universality

A major issue here is that of *partial use and deployment*

- This security mechanism has to cope with partial deployment in the routing system
  - The basic conventional approach of “what is not certified and proved as good must be bad” will not work in a partial deployment scenario
- In BGP we need to think about both origination and the AS Path of a route object in a partial deployed environment
  - AS path validation is challenging indeed in an environment of piecemeal use of secure credentials, as the mechanism cannot tunnel from one BGPsec “island” to the next “island”
- A partially secured environment may incur a combination of high incremental cost with only marginal net benefit to those deploying BGPsec

## 2. Generic Concerns: Prevention vs Detection

Is certification the *only way* to achieve useful outcomes in securing routing?

- Is this form of augmentation to BGP to enforce “protocol payload correctness” over-engineered, and does it rely on impractical models of universal adoption?
- Can various forms of *routing anomaly detectors* adequately detect the most prevalent forms of typos and deliberate lies in routing with a far lower overhead, and allow for unilateral detection of routing anomalies?
- Or are such anomaly detectors yet another instance of “cheap security pantomime” that offer a thinly veiled placebo of apparent security that is easily circumvented or fooled by determined malicious attack?

# 3. My Opinions!

My personal view of a design compromise for secure BGP:

- Improve the robustness of RPKI certs by altering the cert validation algorithm
- Flatten the certificate hierarchy by using a single CA and distributed RAs
- Place origination signatures, ROAs and certs into the BGP protocol updates as opaque attributes
- Use AS Adjacency attestations
- Place AS Adjacency attestations into BGP protocol updates as opaque attributes
- Exploit the use of TCP in BGP to never resend already-sent certs
- Flatten parts of the CA hierarchy by using RAs rather than CA delegations
- Reduce OOB credential distribution to just TA material
  - For which you can use the DNS and DNSSEC if you really want to put all your eggs in one basket!

*Like all the other approaches, this represents a particular set of compromises about speed, complexity, cost, deployment characteristics and robustness – it has its weaknesses in terms of comprehensive robustness, but it attempts to reduce the number of distinct moving parts*