

DDoS Mitigation at CloudFlare

Martin J. Levy CEE Peering Day 2015 – Bratislava Slovakia 19 March 2015

CloudFlare

What is CloudFlare?

CloudFlare makes websites faster and safer using our globally distributed network to deliver essential services to any website

- Performance
- Content
- Optimization
- Security
- 3rd party services
- Analytics





CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy

3

How does CloudFlare work?

CloudFlare works at the network level

- Once a website is part of the CloudFlare community, its web traffic is routed through CloudFlare's global network of 30+ (and growing) data centers.
- At each edge node, CloudFlare manages DNS, caching, bot filtering, web content optimization and third party app installations.





CloudFlare works globally

CloudFlare protects globally

• DDoS attack traffic is localized and lets other geographic areas continue to operate





CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy

IPv4/IPv6 - automagically enabled

CloudFlare offers an automatic IPv6 gateway seamlessly bridging the IPv4 and IPv6 networks

- For most businesses, enabling IPv6 is costly and time consuming
- CloudFlare's solution requires NO hardware, software, or other infrastructure changes by the site owner or hosting provider
- Enabled via the flip of a switch on the site owner's CloudFlare dashboard
- Users can choose two options: (FULL) which will enable IPv6 on all subdomains that are CloudFlare Enabled, or (SAFE) which will automatically create specific IPv6-only subdomains (e.g. www.ipv6.yoursite.com)





IPv6 - automagically enabled

Automatic IPv6

Enable IPv6 support. Learn more...

Automatic IPv6 Enable IPv6 support. Learn more...



7

v

Off

enabled by default!



CloudFlare has customers globally



Nearly two million websites



Anycast CDN

Anycast CDN – equally IPv4 and IPv6

Anycast prefixes

- Same IP prefixes (IPv4 & IPv6) advertised in each of the 30+ sites around the world (and growing)
- Unicast (from separate site-specific prefixes) used to pull traffic from "origin" web source

Traffic Control

- Eyeball ISPs (should) route to closest node, resulting in a very low latency to our services from everywhere in the world
- If ISP A routes to CloudFlare in Germany then traffic will be served from Frankfurt or Düsseldorf
- If ISP B routes to CloudFlare in Texas then traffic will be served from Dallas

This results in a reasonable distribution of attack traffic between our sites

• Easier to mitigate 10 sites receiving a ~50Gbit DDoS than 1 site receiving 1,500Gbit DDoS



Anycast CDN

How does it work?

- DNS Query to anycast DNS address
- DNS result returned with "Anycast" IP
- Client makes connection to closest server
- CloudFlare replies session established

What happens in the event of an outage?

- Traffic re-routes to next closest DC
 - TCP session resets at this point





DDoS Mitigation

Hundreds of millions of packets per second







The Evolving Landscape of DDoS Attacks

14

Layer 7 Attack methodology



Attackers use millions of compromised machines to launch a sophisticated attack that mimics real users and overloads the slow points in your web property.

As infected machines (both end user boxes and datacenter/hosting boxes) become IPv6 enabled, these Layer 7 attacks are showing up as IPv6 attacks



Layer 7 Attack methodology (protected)



Protecting at Layer 7 can be independent of IPv4 or IPv6 usage



DNS Infrastructure Attacks



Protecting DNS can be independent of IPv4 or IPv6 usage



DNSSEC and DNS attacks and open resolvers



CLOUDFLARE

CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy

18

Layer 7 – Malicious Payload

- Request sent to exploit vulnerability on server
 - WAF on CloudFlare blocks 1.2 billion request per day
- Shellshock
 - 10 to 15 attacks per second during the first week
 - Top countries: France (80%), US (7%), Netherlands (7%)





DDoS mitigation

- How can we solve this?
- Detection
- Mitigation
 - Server Farms
 - Network
- Scaling to respond
- IPv6 challenges



How can we solve this?

How can we solve this?

- At this scale no one solution is going to be a "silver bullet"
- The cost per Gbit of scrubbing hardware is insane and not really an option.
- However.... there is a lot you can do with commodity servers and vendor routers.
- Can filter any criterion that can be matched up to layer4 in hardware on any major router vendor's hardware.
- Can do a bunch of tricks to reduce the burden of filtering traffic on your server metal.
- Plus a bunch more.



Detection

How do we even work out what to mitigate

Detection - how do we do it?

If I asked you to tell me what was DDoSing you, without expensive vendor hardware how would you do it?

- tcpdump(1)?
- Some other packet sniffer

Servers under load during attacks (CPU, RAM, etc), despite great scale.

tcpdump attempts to find a large block of contiguous free RAM, then times out if this is not possible, leaving it often useless until the attack is over.

It is also very resource intensive to start sniffing all traffic on a server.



Detection - how do we do it?

So how do we do this?

- Taking the burden of detection away from the device being attacked can be very helpful
 - Export NetFlow records from the edge routers
 - Export sFlow from the switches in our datacenters
 - Automating this process has helped considerable
- Reading data from the application
 - NGINX logs tell you a lot that is useful.
- Sometimes calming the attack down to a manageable level with blunt rules (rate-limit all traffic from these 5x /16s to this single /32) can help to be able to then do deeper inspection and fine-tune the rules we implement to mitigate



Mitigation - on our server farms

CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy

26

DNS - BPF tools + lots and lots of DNS IPs

DNS attacks have a number of unique solutions;

• CloudFlare have many many thousands of DNS servers

\$ host bob.ns.cloudflare.com
bob.ns.cloudflare.com has address 173.245.59.104
bob.ns.cloudflare.com has IPv6 address 2400:cb00:2049:1::adf5:3b68

- o Allows us to distribute the attack more effectively
- o Can null route specific DNS server IPs with minimal impact
- BPF (Berkeley Packet Filter) tools
 - High performance pattern matching driven filtering
 - o Allows us to filter out DNS attack traffic using far less CPU resource
 - http://blog.cloudflare.com/introducing-the-bpf-tools/
 - *https://github.com/cloudflare/bpftools*



Hashlimits

Enforce "no more than X connection attempts per minute for this hash", otherwise blacklist

Hash is made up from whatever criterion you want, but for our purposes combo of src + dest IPs

Fairly effective method of easily detecting "ddos-like" traffic.

Trick is preventing false detections.

- Customer with many millions of users released an application update causing the application to regularly perform JSON queries against their application.
- Users behind a CG-NAT appeared as if they were coming from a single IP.
- Triggered enforcement on non-malicious traffic.



"I'm under attack" mode

Customer enabled mode that forces users to a challenge page.

Uses an early code path in NGINX.

Significantly less CPU required to process requests than going through the full process of serving their request.



Connection tracking to validate sessions

Another way to perform non-cpu-intensive processing of malicious traffic.

Using conntrack to inspect all sessions for certain source/destination IP combinations can be incredibly useful, but it has it's risks as well.

While maintaining a session table can be incredibly helpful and allow you to mitigate malicious traffic early on, it also presents another attack surface (session table flooding).

Used incredibly selectively, particularly for large ACK attacks.



ECMP to distribute traffic between servers

Allows us to ensure no one server bears the entire brunt (for traffic coming into a given site) of the attack load aimed at a single IP. 16 servers can more easily mitigate an attack than 1.

All our servers speak BGP to our routing infrastructure, so this is not particularly difficult to implement.

By default, ECMP hashes will be re-calculated every time there is a next-hop change.

- Causes flows to shift between servers
 - TCP sessions reset
- Can solve this with consistent ECMP hashing
 - Available in Junos from 13.3R3 for any trio based chipset
 - Only works for up to 1k unicast prefixes, so struggles to scale



Solarflare cards and OpenOnload

In our latest generation of server hardware we;

- Made the move to 2x10Gbit per server (from 6x1Gbit LAGs)
- Did this with NICs from Solarflare.

SolarFlare NICs have very cool abilities to pre-process traffic on-board before handing to the CPU (OpenOnload).

Can identify certain types of traffic and assign it to cores based on rules pushed in the cards.

Can handle certain requests in userspace without creating CPU interrupts



Cloudflare have been helping the SolarFlare develop this functionality for their cards.

http://blog.cloudflare.com/a-tour-inside-cloudflares-latest-generation-servers/

Mitigation - in the network

Null route and move on

When an attacker targets a website or a service, while they may want to take this website/service down, they target the IP address in order to do this.

First order of business can be to update the DNS A/AAAA record and move on.

If the attacker follows, keep doing this.

Easy to automate, requires an attacker to continually change the attack to follow.

Depends on rDNS service operators honouring our TTLs



FlowSpec (RFC 5575)

Important to understand from the outset that ALL flowspec does is automate the provisioning of a backplane-wide firewall filter on multiple devices. Having said that, **it does this really well**.

Can use most "from" and "then" actions available in Juniper firewall filters in FlowSpec. While Juniper have been an early adopter, other vendors have struggled to get this into their code. Even Juniper has only recently implemented IPv6 support for FlowSpec.

Being able to match "TCP packets from this /24, to this /32, with SYN but no ACK and a packet length of 63 bytes" and "rate-limit to 5Mbit" per edge router is incredibly useful.

Being able to configure this in one place and have it push to the entire network is awesome!



Regional enforcement

Under certain circumstances, it makes sense to enforce regionally

- Seeing 300Gbit of traffic targeted at AMS, LHR, FRA, CDG for a website with 99% of legitimate traffic being served into HKG and SIN
 - Can implement strict flowspec enforcement in sites targeted, while no enforcement needed in sites traffic is legitimately needed in.
 - Take advantage of any opportunity presented

Regional null routing can also be worthwhile at times

- Want to move site to new IPs and move on.
 - Null route in only the regions that are being targeted.

Have your transit provider configure firewall filters in their network to filter certain packet types / lengths / src-IPs / dst-IPs / etc upstream in one region only to help filter malicious traffic.



Dealing with attacks on infrastructure IPs

Relatively easy to mitigate attacks on Anycast IP space.

- Multiple hundred gig attack on an anycast IP
 - o Distributed over 28 sites
 - Multiple tens of gigs per site
- Vs:
 - Multiple hundred gig attack on an IP specific to a single router, link or DC
 - Very hard to mitigate
 - Multiple hundred gig attack traffic > 100Gbit link

How do we prevent this from happening?

What can we do about it? What gain do you get from exposing this?



Attacks on Infrastructure - obfuscation of IPs

Traceroutes that show you the full path are nice... but... at what expense?

- Reveals a lot of the IP addressing information of your infrastructure to the entire internet
 - Becomes easy to figure out what to attack.
 - Makes every linknet, loopback, and infrastructure IP a target

Worth at least considering obscuring some of your infrastructure

- Stop responding to ICMP and UDP ttl expired
- Avoid ICMP-Packet-Too-Big in IPv6
 - Killing this can cause serious problems.`



Attacks on Infrastructure - kill routability to IPs

Can take the next step and kill reachability entirely.

Make your linknet IPs non-routable;

- Take all your linknet IPs from a /24 that is not advertised on the internet
- Use RFC1918 space
- Blackhole all your linknets
 - Don't forget to blackhole the provider side also!

This can make debugging significantly harder!

A lot of work will need to be done in the pre-sales stage with transit providers to ensure that one of these options is possible.

Peering exchanges should not be reachable on the internet anyway



Scaling the network to respond

CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy 4

40

Scaling the network - capacity

Ultimately, this is all a capacity game.

If you are seeing attacks roughly equivalent to your transit capacity you'll struggle to mitigate it, however if your transit capacity is 10x or 100x the size of the attack you'll struggle a lot less

Attack traffic will often be very very small packets (<80bytes) so need to ensure that any routers, line cards, are specced for capacity based on those numbers.

- As an example, MPC4E cards perform at nowhere near line rate when challenged with hundreds of gigs of 64 byte packets
- How oversubscribed is your backplane
- Measurements should always be in PPS rather than Gbit/sec

As you scale up your routers, you may discover that PPS bottlenecks simply move to your transit providers.



Peering vs Transit

Far easier to mitigate a DDoS coming in on an IX than a DDoS coming in via a transit provider.

- Can negotiate with transit provider for features such as RTBH, NOC implementing firewall filters for you, etc
- Peering exchanges generally don't have these features.

Peering exchanges are also surprisingly expensive to scale up for DDoS. Generally will be more expensive to order more 10Gbit ports at an IX vs additional handovers to a transit provider.

Often end up de-peering a network sourcing large amounts of attack traffic to force them onto a transit provider where you have more control.

This seems broken - surely there is a better way to ingest this traffic?



The solution - lots and lots of PNIs

PNI with networks sourcing DDoS traffic in multiple locations

• Limits the scope of impact to the network sourcing the attack, or their upstream who you are peering with

Easy to scale very cheaply

• Only costs are router interfaces, DC cross connects

Can easily filter traffic from that specific network on your router interfaces



Scaling with more sites

The more we scale, the more we will distribute attack traffic, and the greater attack mitigation capacity we will have.



44

Scaling using caches

A step beyond scaling up the sites. Ability to place caches in hundreds, if not thousands of locations within ISP networks.

Distribute attack traffic significantly.





IPv6 is a challenge

IPv6 changes this all... a lot

For all network operators, IPv6 presents significant security challenges

- How good is your IPv6 security compared to your IPv4 security
- Obfuscation via NAT is no longer an option
- Are your IPv6 ACLs written as well as your IPv4 ACLs?
- Do all the tools you use support IPv6?
- Do your traffic measurement and DDoS mitigation tools fully support IPv6

We can reasonably assume that it is more likely that reflection attacks can be launched at routers, servers, etc using IPv6, as services won't be as consistently secured

- NTP attacks, DNS attacks, SNMP?
- Who knows what else?



IPv6 changes this all... a lot

We are currently seeing over ¼ of all DNS attacks coming in via IPv6. Balance of malicious vs legitimate DNS traffic about the same on v4 and v6.

No CG-NAT on IPv6 means

- Not going to falsely identify legitimate traffic as malicious due to many hosts behind a CG-NAT
- One bot is not going to cause us to challenge thousands of users because their IP is in a blacklist

Scaling issues on routers

- IPv6 addresses are 400% the size of IPv4 address
 - At scale, this may long term present challenges as to the number of null routes, flowspec rules, etc that we can push to our router silicon
 - At this point, should we start doing enforcement on a per /64 basis?



CEE Peering Day 2015 - CloudFlare DDoS Mitigation - Martin J Levy 48

IPv4

ECMP & ICMP & IPv6

ECMP (Equal Cost Multi Path routing) & ICMP

TCP it hashes a tuple (src-ip, src-port, dst-ip, dst-port)

- Uses this hash to chose destination server. This guarantees that packets from one "flow" will always hit the same server
- No special knowledge of ICMP. It hashes only (src-ip, dst-ip)
- ICMP src-ip is most likely the IP of a intermediary router on the Internet
- IPv6 requires PMTU discovery to work; ECMP breaks it
- Solution: Broadcast ICMP PMTU messages to all servers
- Posted source code: https://github.com/cloudflare/pmtud





Questions?

Thank you!

