

Super Fast Packet Filtering with eBPF and XDP

Helen Tabunshchyk, Systems Engineer, Cloudflare

@advance_lunge



Agenda

- 1. Background.
- 2. A tiny bit of theory about routing.
- 3. Problems that have to be solved.
- 4. Overview of existing solutions.
- 5. DDoS mitigation pipeline.
- 6. eBPF and XDP.
- 7. Bonus part.

A bit of context about the work I do



What does Cloudflare do?



CDN

- Moving content physically closer to visitors with our CDN
- Intelligent caching
- Unlimited DDOS mitigation
- Unlimited bandwidth at flat pricing with free plans



Website Optimisation

- Making web fast and up to date for everyone.
- TLS 1.3 (with 0-RTT)
- HTTP/2 + QUIC
- Server push
- AMP
- Origin load-balancing
- Smart routing
- Workers
- Post quantum crypto
- Many more



DNS

- Cloudflare is the fastest managed DNS providers in the world.
- 1.1.1.1
- · 2606:4700:4700::1111
- DNS over TLS

Scale

- 154 data centres in 74 countries
- More than 10 million domains
- 10% of all Internet requests
- 7.5M requests per second on average, 10M at peak
- 1.6M DNS queries per second
- 2.8 billion people served each month
- Biggest DDoS attack 942 Gbps
- 20 Tbps network capacity and growing





Life of a packet

A long time ago in a galaxy far, far away...





The OPTE Project Internet 2015 Map



North America (ARIN)

Europe (RIPE)

Latin America (LACNIC)

Asia Pacific (APNIC)

Africa (AFRINIC)

"Backbone" (highly connected networks)

http://www.opte.org

Load Balancing Between Data Centres

- Locality and congestion control
- DNS
- BGP
- Anycast



https://www.cloudflare.com/learning/dns/what-is-dns/





THE PACKET 5 6 1 7 8 6 5 5 S C V S L C T I F F

BGP — A Tale of Two Napkins

At an Internet Engineering Task Enco (IETF) conferences had January, Sitk Laugheed and Ice. Broad of these and Walton Bechief. of 1000 sat closes in the meeting hall cafeconic and average needs controppedatest. What has since become REC 1025, the Bordes. Barray Protocol (BCP), is will known to some as the "Two-Nackin Partners," in reference to the

PTC 1105. The Booker Generater Franks of in still. known to some as the Teo-Nepillo Peotocol

handy media is a post which the malmenty first deal of 11.

Arconding to Longhead, riseo?). director of software explored by BGP developed as a solution to the deficiencies of DGP. The peaks Icm everyoil with the conservation increases in the number of increases. hasts, and with its expanding. upplear, "The internet Present a its somewhall beyond any only. expectations," Longiered copicing, "Diff" and single out chromed to handly networks of this size," With the internet's divergification and expanding ranting damages, at belong none expension and ditaexecute score montrol over their resources by introducing different. tapes of new policies. IVIP could so provisions for such policies. Nor ... Improve taken before transmission. Od 11 years to happy runnions of

networks. The networking comman AS heating to express a degree of concern that the core renting system would simply fail at some solat. Morecret, EGF. showed furthers since of mostnext as increasingly large non-ting updates were such over the Interact. Detagrams contain ing these updates outpress the ARPANET's meximum transport. size of 1008 listes, thus requiring applying of a la

cisco Makes Bold Entry to OSI Marketplace, Designs Largest OSI Network Demo to Date

The meet complice OSI network man assessibled run throughout the Interpy 59 madedhow, this year in the San Jase Convention Center, Sorthern California, All togethers, about 14 condex supporting the 051 setwork protocol accountaily. interconnected their systems to

form the Interop (ISI demo notwork.

timo played a major role in the triamph of the UEI demo. Houses from doos - suming the ISO. CDNS (Connectionless Network) benear located (colored continued on p. 3.



Types of Routing







Okay, our little packet is inside the DC



1. Uneven load





2. Different kinds of traffic





3. Per packet load balancing





4. Heterogenous hardware





5. Locality (e.g. for cache) and transport affinity





6. DDoS





Types of DDoS Attacks

	Volumetric Attack	Protocol Attack	Application Attack
What is it?	Saturating the bandwidth of the target.	Exploiting a weakness in the Layer 3 and Layer 4 protocol stack.	Exploiting a weakness in the Layer 7 protocol stack.
How does it cripple the target?	Blocks access to the end-resource.	Consume all the processing capacity of the attacked-target or intermediate critical resources.	Exhaust the server resources by monopolising processes and transactions.
Examples	NTP Amplification, DNS Amplification, UDP Flood, TCP Flood, QUIC HelloRequest amplification	Syn Flood, Ping of Death, QUIC flood	HTTP Flood, Attack on DNS Services



7. Group change





8. Graceful connection draining





Load balancing techniques

ECMP

ID (packet) mod N,

- ID some function that produces connection ID, e.g. 5-tuple flow;
- N the number of configured backends.

Uneven load

Different kinds of traffic

Per packet load balancing

Heterogenous hardware

Locality

DDoS

Group change

Graceful connection draining





Payload

ECMP-CH

populating the ECMP table not simply with next-hops, but with a slotted table that's made up of redundant next-hops

Uneven load Different kinds of traffic Per packet load balancing Heterogenous hardware Transport affinity DDoS Group change Graceful connection draining





Stateful Load Balancing

Uneven load Different kinds of traffic Per packet load balancing Heterogenous hardware Transport affinity DDoS Group change Graceful connection draining





Google Maglev





Daisy Chaining a.k.a Beamer

- Beamer muxes do not keep per-connection state; each packet is forwarded independently.
- When the target server changes, connections may break.
- Beamer uses state stored in servers to redirect stray packets.









• Packets contain previous server and time of reassignment





• New connections are handled locally





• Daisy chained connections die off in time


Daisy Chaining a.k.a Beamer

Uneven load

Different kinds of traffic

Per packet load balancing

Heterogenous hardware

Transport affinity

DDoS

Group change

Graceful connection draining

Performance



https://www.usenix.org/conference/nsdi18/presentation/olteanu https://github.com/Beamer-LB



Also FPGA Packet Processing

- Early experiments with FPGAs
- Smart NICs
- P4 language



Fun (?) Facts

https://www.fastly.com/blog/anatomy-an-iot-botnet-attack

An average IoT device gets infected with malware and launches an attack within 6 minutes of being exposed to the internet. Over the span of a day an average of over 400 login attempts per device; 66 percent of them on average are successful.

Over the span of a day, IoT devices are probed for vulnerabilities 800 times per hour.



DDoS Mitigation

Disclaimer



r

Jérôme Petazzoni @jpetazzo

Follow V

OH: "In any team you need a tank, a healer, a damage dealer, someone with crowd control abilities, and another who knows iptables"

10:41 AM - 27 Jun 2015 from Kansas City, MO





Disclaimer

BPF and eBPF

- Low overhead sandboxed user-defined bytecode running in kernel
- Written in a subset of C, compiled by clang llvm
- It can never crash, hang or interfere with the kernel negatively
- If you run Linux 3.15 or newer, you already have it
- Great intro from Brendan Gregg: http://www.brendangregg.com/ebpf.html

opessnoop c* tava* node* myseld eslower gethostlatency Other: filetop filelife fileslowsr statsnoop php* gythan* hashreadlise nenleak espuble EUDY* synesnoop seleniff vfsocunt vfsstat monlin wflow syscost cachestat cachetos use upbjnew killsnoop dostat dosnoop untat othrands exectnoop mountaneop pidpersed Applications coudist. System Libraries rusglat runglen trace deaclock detector anglist System Call Interface counclaimed functiont VES funcalower Gockets Scheduler offcputime funclatency File Systems TCP/UDP wahauphing starscoutt. offwaketize P profile Volume Manager Virtual Memory 7 softires Block Device Interface Ethernel pomkill memleak **Device Drivers** mdflueb slabratetop hardirgs ttysnoop btrfedist btrfsslowe: DRAM coptop toplife coptracer extédist estésiower tepcomect topaccept xfediat xfuslower llestat tesconslat tegretrans afediat afsolower CPU profile biotop hiossoop bielatency bitesize

Linux bcc/BFF Tracing Tools





Classical BPF Machine

Extended BPF Machine

End of disclaimer

DDoS Mitigation Pipeline

Gatebot





iptables

- Initially it was the only tool to filter traffic
- Leveraged modules ipsets, hashlimit, connlimit
- With the xt_bpf module it was possible to specify complex filtering rules
- But we soon started experiencing IRQ storms during big attacks
- All CPUs were busy dropping packets, userspace applications were starving of CPU





Userspace Offload a.k.a. Kernel Bypass

- Network traffic is offloaded to userspace before it hits the Linux network stack
- Allows to run BPF in userspace
- An order of magnitude faster than iptables (5M pps)



- Requires one or more CPUs to busy poll the NIC event queue
- Reinjecting packets in the network stack is expensive
- Hardware dependant





XDP to the rescue!

XDP Packet Processing Overview



+AF_XDP since 4.19



```
1 SEC("xdp1")
     int xdp_prog(struct xdp_md * ctx)
      void * data = (void *)(long)ctx->data;
       void * data_end = (void *)(long)ctx->data_end;
       int ret;
       ret = rule_1(data, data_end);
       if (ret != XDP_PASS) // multiple if statements
         return ret;
12
13
14
       ret = rule_2(data, data_end); // one for each rule
      if (ret != XDP_PASS)
15
17
         return ret;
21
       return XDP_PASS; // if none of them match - the packet is accepted
22
     £
```



```
1 static inline int rule_1 (void * data, void * data_end)
 2 {
     if (!condition_1)
 4
       return XDP_PASS;
     if (!condition_2) // multiple if conditions
       return XDP_PASS;
10
11
12
     . . .
13
    update_rule_counters(1); // do additional work
14
     sample_packet(data, data_end);
15
16
17
     return XDP_DROP;
18 }
```

```
2 struct bpf_map metrics_map __section("maps") = {
       .type = BPF_MAP_TYPE_PERCPU_ARRAY,
      .key_size = sizeof(unsigned int),
      .value_size = sizeof(rule_metrics_t),
       .max_entries = STATS_MAP_SIZE,
7 };
 9
10 static inline void update_rule_counters(int rule_id)
11 {
12
     long * value = bpf_map_lookup_elem(&metrics_map, &rule_id);
13
    if (value)
14
15
16
      *value += 1;
17
    ·}
18 }
```



Packet dropping performance

A perfect match: XDP both for load balancing and DDoS mitigation <3

XDP L4LB with daisy chaining using encapsulation

Uneven load Different kinds of traffic Per packet load balancing Heterogenous hardware Transport affinity DDoS Group change Graceful connection draining Performance





And They Lived Happily Ever After

But... DPDK?

Advantages of XDP over DPDK

- Allows option of busy polling or interrupt driven networking
- No need to allocate huge pages
- Dedicated CPUs are not required, user has many options on how to structure the work between CPUs
- No need to inject packets into the kernel from a third party user space application
- No special hardware requirements
- No need to define a new security model for accessing networking hardware
- No third party code/licensing required



Bonus part



(Quick UDP Internet Connections)





HTTP Request Over TCP + TLS



HTTP Request Over QUIC





L4LB to the rescue!
Thank you!

