# eBPF + Observability Day

CSCI 2952-F – Distributed Systems at Scale: Microservices Management

Palani Kodeswaran    Sayandeep Sen

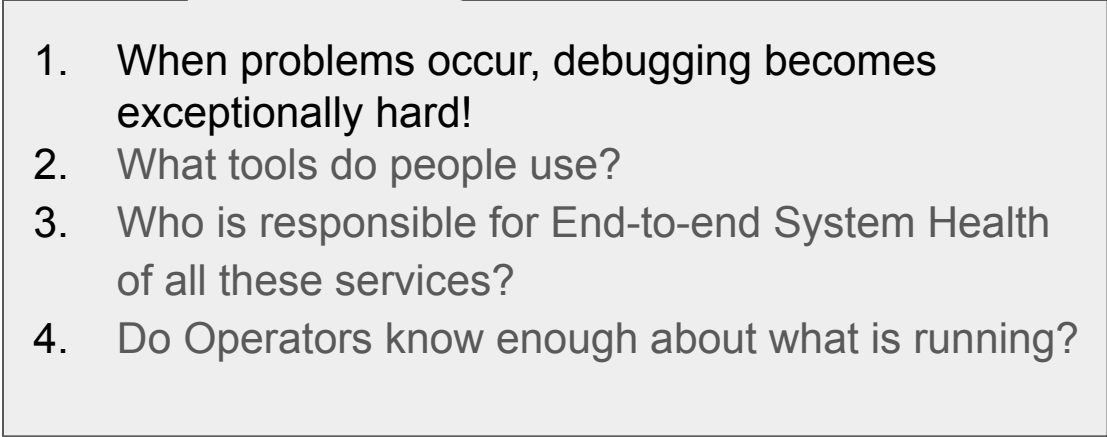IBM Research

# Microservices: What's awesome ?

# Microservices: (not-so good parts)

1. Too many layers [Isolation, security and privacy]
2. Workloads move all over [Scalability, Elasticity]
3. No-organization has full end-to-end visibility [Business secrets]
4. Distance between application developer and deployment increases [Separation of concern]

# Microservices: (not-so good parts)

1. Too many layers [Isolation, security and privacy]
2. Workloads move all over [Scalability, Elasticity]
3. No-organization has full end-to-end visibility [Business secrets]
4. Distance between application developer and deployment increases [Separation of concern]

1. When problems occur, debugging becomes exceptionally hard!
2. What tools do people use?
3. Who is responsible for End-to-end System Health of all these services?
4. Do Operators know enough about what is running?

# Debuggability

# Debuggability

*Figure 2-1. Adding features to the kernel (cartoon by Vadim Shchekoldin, Isovalent)*

# Enter eBPF

**Google** uses eBPF for security auditing, packet processing, and performance monitoring

VIDEO 1 · VIDEO 2 · TALK 1 · TALK 2

**Netflix** uses eBPF at scale for network insights

BLOG

**Cloudflare** uses eBPF through XDP for DDoS mitigation and load balancing

BLOG · TALK

**Meta** uses eBPF to process and load balance every packet coming into their data centers

VIDEO · BLOG 1 · BLOG 2 · TALK 1 · TALK 2

**Ikea** uses eBPF through Cilium for networking and load balancing in their private cloud

VIDEO

**Android** uses eBPF to monitor network usage, power, and memory profiling

DOCS

**Datadog** uses eBPF for networking and security in their SaaS product

VIDEO

**Alibaba** uses eBPF through Cilium to provide networking in their cloud

BLOG

**Seznam** uses eBPF for load balancing

BLOG

# eBPF & Kubernetes - the future

**"The Linux kernel continues its march towards becoming BPF runtime-powered microkernel."**

**Tiny core kernel** with user definable kernel functionality in BPF (instead of kernel modules)

Less security bugs & kernel crashes due to **smaller attack surface** and **safety-verified code**

Drastic reduction of 'static' feature creep for **better resource efficiency**

**Kubernetes** then ships **custom BPF**-tailored extensions to optimize needs for user workloads

Today's kube-proxy replacement through BPF is just a tiny dot in that universe ...

**Steven Rostedt**
@srostedt

BPF will replace Linux #kr2019

11:06 am · 26 Sep 2019 · Twitter for Android

2013   2014   2015   2016   2017   2018   2019   2020                    ?

Toke Høiland-Jørgensen, https://lwn.net/ml/bufferbloat/87bls8bnsm.fsf@toke.dk/

# eBPF: Full System Visibility (and much much more !!)

- Trace and debug of third party software!!
- Visibility of whole system deployment
- Bolt On (Application Transparent)

Besides applications in Networking, Security and Storage …

Linux bcc/BPF Tracing Tools

# Case Study : Slack Keeps Crashing

# Debugging Slack

## Slack's Secret STDERR Messages

27 Aug 2021

*These are rough notes.*

I run the Slack messaging application on Ubuntu Linux, and it recently started mysteriously crashing. I'd Alt-Tab and find it was no longer there. No error message, no dialog, just gone. It usually happened when locking and unlocking the screen. A quick internet search revealed nothing.

From here:

# 1. Enabling core dumps

I'm guessing it's core dumping and Ubuntu's apport is eating them. Redirecting them to the file system so I can then do core dump analysis using gdb(1), as root:

```
# cat /proc/sys/kernel/core_pattern
|/usr/share/apport/apport %p %s %c %d %P
# mkdir /var/cores
# echo "/var/cores/core.%e.%p.%h.%t" > /proc/sys/kernel/core_pattern
[...another crash...]
# ls /var/cores
#
```

This didn't work: No core file showed up. I may need to increase the core file size ulimits for Slack, but that might mean mucking around with its startup scripts; I'll try some other tracing first.

From here:   Slack's Secret STDERR Messages (brendangregg.com)

# 2. exitsnoop

Using an eBPF/bcc tool to look for exit reasons:

```
# exitsnoop -t
TIME-AEST       PCOMM                PID     PPID     TID      AGE(s)   EXIT_CODE
13:51:19.432 kworker/dying        3663305 2        3663305 1241.59 0
13:51:30.948 kworker/dying        3663626 2        3663626 835.76  0
13:51:33.296 systemd-udevd        3664149 2054939 3664149 3.55       0
13:53:09.256 kworker/dying        3662179 2        3662179 2681.15 0
13:53:25.636 kworker/dying        3663520 2        3663520 1122.60 0
13:53:30.705 grep                 3664239 6009     3664239 0.08       0
13:53:30.705 ps                   3664238 6009     3664238 0.08       0
13:53:40.297 slack                3663135 1786     3663135 1459.54 signal 6 (ABRT)
13:53:40.298 slack                3663208 3663140 3663208 1457.86 0
13:53:40.302 slack                3663140 1786     3663140 1459.18 0
13:53:40.302 slack                3663139 1786     3663139 1459.18 0
13:53:40.303 slack                3663171 1786     3663171 1458.22 0
13:53:40.317 slack                3663197 1786     3663197 1458.03 0
13:53:44.827 gdm-session-wor      3664269 1778     3664269 0.02       0
[...]
```

This traced a Slack SIGABRT which happened around the same time as a crash. A strong lead.

From here:  Slack's Secret STDERR Messages (brendangregg.com)

# 5. signals.bt

The signals.bt bpftrace tool from my BPF book traces the signal:signal_generate tracepoint, which should catch every type of generated signal, including tgkill(2). Trying it out:

```
# bpftrace /home/bgregg/Git/bpf-perf-tools-book/originals/Ch13_Applications/signals.bt
Attaching 3 probes...
Counting signals. Hit Ctrl-C to end.
^C
@[SIGNAL, PID, COMM] = COUNT

@[SIGPIPE, 1883, Xorg]: 1
@[SIGCHLD, 1797, dbus-daemon]: 1
@[SIGINT, 3665167, bpftrace]: 1
@[SIGTERM, 3665198, gdm-session-wor]: 1
@[SIGCHLD, 3665197, gdm-session-wor]: 1
@[SIGABRT, 3664940, slack]: 1
@[SIGTERM, 3665197, gdm-session-wor]: 1
@[SIGKILL, 3665207, dbus-daemon]: 1
@[SIGWINCH, 859450, bash]: 2
@[SIGCHLD, 1778, gdm-session-wor]: 2
@[, 3665201, gdbus]: 2
@[, 3665199, gmain]: 2
@[SIGWINCH, 3665167, bpftrace]: 2
@[SIGWINCH, 3663319, vi]: 2
@[SIGCHLD, 1786, systemd]: 6
@[SIGALRM, 1883, Xorg]: 106
```

Ok, there's the SIGABRT for slack. (There's a new sigsnoop(8) tool for bcc that uses this tracepoint as well.)

From here:  Slack's Secret STDERR Messages (brendangregg.com)

```
# egrep -i 'error|fail' webapp-console.log
[08/25/21, 16:07:13:051] info: [DESKTOP-SIDE-EFFECT] (TKZ41AXQD) Reacting to  {"type":"[39] Set a
[08/25/21, 16:07:13:651] info: [DESKTOP-SIDE-EFFECT] (T7GLTMS0P) Reacting to  {"type":"[39] Set a
[08/25/21, 16:07:14:249] info: [DESKTOP-SIDE-EFFECT] (T0DS04W11) Reacting to  {"type":"[39] Set a
[08/25/21, 16:07:14:646] info: [DESKTOP-SIDE-EFFECT] (T0375HBGA) Reacting to  {"type":"[39] Set a
[...]
# egrep -i 'error|fail' browser.log
[07/16/21, 08:18:27:621] error: Cannot override webPreferences key(s): webviewTag, nativeWindowOp
[07/16/21, 08:18:27:653] error: Failed to load empty window url in window
   "error": {
     "stack": "Error: ERR_ABORTED (-3) loading 'about:blank'\n    at rejectAndCleanup (electron/js
[07/16/21, 08:18:31:355] error: Cannot override webPreferences key(s): webviewTag, nativeWindowOp
[07/16/21, 08:18:31:419] error: Cannot override webPreferences key(s): webviewTag, nativeWindowOp
[07/24/21, 09:00:52:252] error: Failed to load calls-desktop-interop.WindowBorderPanel
   "error": {
     "stack": "Error: Module did not self-register: '/snap/slack/42/usr/lib/slack/resources/app.as
[07/24/21, 09:00:52:260] warn: Failed to install protocol handler for slack:// links
[07/24/21, 09:00:52:440] error: Cannot override webPreferences key(s): webviewTag
[...]
```

I browsed the logs for a while but didn't see a smoking gun. Surely it spits out some error message when crashing, like to STDERR...

From here:  Slack's Secret STDERR Messages (brendangregg.com)

# 8. STDERR Tracing

Where is STDERR written?

```
# lsof -p 3666477
[...]
slack    3666477 bgregg   mem     REG                    7,16      141930     7165 /snap/slack/44/usr/lib/
slack    3666477 bgregg   mem     REG                    7,16      165680     7433 /snap/slack/44/usr/lib/
slack    3666477 bgregg    0r     CHR                    1,3          0t0        6 /dev/null
slack    3666477 bgregg    1w     CHR                    1,3          0t0        6 /dev/null
slack    3666477 bgregg    2w     CHR                    1,3          0t0        6 /dev/null
slack    3666477 bgregg    3r     FIFO                   0,12         0t0 29532192 pipe
slack    3666477 bgregg    4u     unix 0x00000000134e3c45            0t0 29526717 type=SEQPACKET
slack    3666477 bgregg    5r     REG                    7,16    10413488     7167 /snap/slack/44/usr/lib/
[...]
```

/dev/null? Like that's going to stop me. I could trace writes to STDERR, but I think my old shellsnoop(8) tool (another from eBPF/bcc) already does that:

From here:   Slack's Secret STDERR Messages (brendangregg.com)

/dev/null? Like that's going to stop me. I could trace writes to STDERR, but I think my old shellsnoop(8) tool (another from eBPF/bcc) already does that:

```
# shellsnoop 3666477
[...]
[08/27/21, 14:46:36:314] info: [DND] (TKZ41AXQD) Will check for changes in DND status again in
5 minutes
[08/27/21, 14:46:37:337] info: [FOCUS-EVENT] Client window blurred
[08/27/21, 14:46:40:022] info: [RTM] (T029N2L97) Processed 1 user_typing event(s) in
channel(s) C0S928EBE over 0.10ms
[08/27/21, 14:46:40:594] info: [RTM] (T029N2L97) Processed 1 message:message_replied event(s)
in channel(s) C0S928EBE over 2.60ms
[08/27/21, 14:46:40:595] info: [RTM] Setting a timeout of 37 ms to process more rtm events
[08/27/21, 14:46:40:633] info: [RTM] Waited 37 ms, processing more rtm events now
[08/27/21, 14:46:40:653] info: [RTM] (T029N2L97) Processed 1 message event(s) in channel(s)
C0S928EBE over 18.60ms


[08/27/21, 14:46:44:938] info: [RTM] (T029N2L97) Processed 1 user_typing event(s) in
channel(s) C0S928EBE over 0.00ms

(slack:3666477): Gtk-WARNING **: 14:46:45.525: Could not load a pixbuf from icon theme.
This may indicate that pixbuf loaders or the mime database could not be found.
**
Gtk:ERROR:../../../../gtk/gtkiconhelper.c:494:ensure_surface_for_gicon: assertion failed
(error == NULL): Failed to load /usr/share/icons/Yaru/16x16/status/image-missing.png: Unable
to load image-loading module: /snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-
2.0/2.10.0/loaders/libpixbufloader-png.so: /snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-
2.0/2.10.0/loaders/libpixbufloader-png.so: cannot open shared object file: No such file or
directory (gdk-pixbuf-error-quark, 5)
```

From here:   Slack's Secret STDERR Messages (brendangregg.com)

It's the .so file that is missing, not the .png:

```
# ls -lh /usr/share/icons/Yaru/16x16/status/image-missing.png
-rw-r--r-- 1 root root 535 Nov  6  2020 /usr/share/icons/Yaru/16x16/status/image-missing.png
# ls -lh /snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.10.0/loaders/libpixbufloader-pr
ls: cannot access '/snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.10.0/loaders/libpixbu
```

But there is a .so file with a similar path:

```
# ls -lh /snap/slack/
total 0
drwxrwxr-x 8 root root 123 Jul 14 02:49 43/
drwxrwxr-x 8 root root 123 Aug 18 10:27 44/
lrwxrwxrwx 1 root root   2 Aug 24 09:48 current -> 44/
# ls -lh /snap/slack/44/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.10.0/loaders/libpixbufloader-pr
-rw-r--r-- 1 root root 27K Aug 18 10:27 /snap/slack/44/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.
```

Hmm, I wonder...

From here:  [Slack's Secret STDERR Messages (brendangregg.com)](Slack's Secret STDERR Messages (brendangregg.com))

# 9. Workaround

This is obviously a hack and is not guaranteed to be safe:

```
# cd /snap/slack
# ln -s current 42
# ls -lh
total 0
lrwxrwxrwx 1 root root    7 Aug 27 15:01 42 -> current/
drwxrwxr-x 8 root root  123 Jul 14 02:49 43/
drwxrwxr-x 8 root root  123 Aug 18 10:27 44/
lrwxrwxrwx 1 root root    2 Aug 24 09:48 current -> 44/
# ls -lh /snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.10.0/loaders/libpixbufloader-pr
-rw-r--r-- 1 root root 27K Aug 18 10:27 /snap/slack/42/usr/lib/x86_64-linux-gnu/gdk-pixbuf-2.0/2.
```

I don't know why Slack was looking up this library via the old directory version, but linking the new version to the old path did the trick. Slack has stopped crashing!

I'm guessing this is a problem with how the snap is built. Needs more debugging.

From here:  Slack's Secret STDERR Messages (brendangregg.com)

# Summary

eBPF has super powers that gives it unprecedented flexibility and visibility to diagnose problems in even third party code!!!

# eBPF Internals

# What is eBPF



From here: What is eBPF? An Introduction and Deep Dive into the eBPF Technology

userspace

app

syscalls

kernel

eBPF program

@lizrice

aqua

# ● man bpf

eBPF programs can be written in a restricted C that is compiled (using the clang compiler) into eBPF bytecode.  Various features are omitted from this restricted C, such as loops, global variables, variadic functions, floating-point numbers, and passing structures as function arguments.

aqua

# Helper Calls

eBPF programs cannot call into arbitrary kernel functions. Allowing this would bind eBPF programs to particular kernel versions and would complicate compatibility of programs. Instead, eBPF programs can make function calls into helper functions, a well-known and stable API offered by the kernel.



The set of available helper calls is constantly evolving. Examples of available helper calls:

- Generate random numbers
- Get current time & date
- eBPF map access
- Get process/cgroup context
- Manipulate network packets and forwarding logic

From here: What is eBPF? An Introduction and Deep Dive into the eBPF Technology

# Maps

A vital aspect of eBPF programs is the ability to share collected information and to store state. For this purpose, eBPF programs can leverage the concept of eBPF maps to store and retrieve data in a wide set of data structures. eBPF maps can be accessed from eBPF programs as well as from applications in user space via a system call.



From here: What is eBPF? An Introduction and Deep Dive into the eBPF Technology

# Hook Overview

eBPF programs are event-driven and are run when the kernel or an application passes a certain hook point. Pre-defined hooks include system calls, function entry/exit, kernel tracepoints, network events, and several others.



```
int syscall__ret_execve(struct pt_regs *ctx)
{
        struct comm_event event = {
                .pid = bpf_get_current_pid_tgid() >> 32,
                .type = TYPE_RETURN,
        };

        bpf_get_current_comm(&event.comm, sizeof(event.comm));
        comm_events.perf_submit(ctx, &event, sizeof(event));

        return 0;
}
```

From here: What is eBPF? An Introduction and Deep Dive into the eBPF Technology

# User Land probes & kernel probes

If a predefined hook does not exist for a particular need, it is possible to create a kernel probe (kprobe) or user probe (uprobe) to attach eBPF programs almost anywhere in kernel or user applications.



From here:

# Tool chain

Restricted C

Clang    LLVM    → BPF bytecode

BPF() syscall

Kernel

Verifier → JIT Compiler → ELF Object File

- The process loading the eBPF program holds the required capabilities (privileges). Unless unprivileged eBPF is enabled, only privileged processes can load eBPF programs.
- The program does not crash or otherwise harm the system.
- The program always runs to completion (i.e. the program does not sit in a loop forever, holding up further processing).

From here: life-bpf-program.png

ELF object file
- eBPF opcodes
- eBPF maps

user space

bpf() system calls

BPF_PROG_LOAD
BPF_MAP_CREATE

kernel

verifier

BPF vm

maps

@lizrice

aqua

**ELF object file**

- eBPF opcodes
- eBPF maps

**user space**

bpf() system calls

BPF_PROG_LOAD
BPF_MAP_CREATE

*Attach BPF program to event*

**kernel**

verifier

BPF vm

maps

@lizrice

aqua

# Hands ON

# eBPF is Everywhere!!



Linux bcc/BPF Tracing Tools

https://github.com/iovisor/bcc#tools 2019

# Tracing: Syscall Count

```
# Syscall count by program
  bpftrace -e 'tracepoint:raw_syscalls:sys_enter { @[comm] = count(); }'
```

```
root@ebpf1:~/tutorial_examples# bpftrace -e 'tracepoint:raw_syscalls:sys_enter { @[comm] = count(); }'
Attaching 1 probe...
^C

@[falcon-sensor]: 3
@[multipathd]: 12
@[AccessLogFlush]: 15
@[kiali]: 22
@[bpftrace]: 25
@[tmux: server]: 32
@[sshd]: 39
@[dockerd]: 41
@[grpc_global_tim]: 60
@[Executor Servic]: 80
@[wrk:worker_12]: 106
@[wrk:worker_3]: 106
@[wrk:worker_46]: 106
@[wrk:worker_5]: 106
```

# Tracing: Software Faults

```
# Count page faults by process

  bpftrace -e 'software:faults:1 { @[comm] = count(); }'
```

```
root@ebpf1:~/tutorial_examples# bpftrace -e 'software:faults:1 { @[comm] = count(); }'
Attaching 1 probe...



^C

@[pilot-discovery]: 1
@[kindnetd]: 5
@[coredns]: 22
@[kubelet]: 881
@[dpkg]: 1257
@[containerd-shim]: 1282
@[BESClient]: 1876
```

# Tracing: Read Bytes By Process

```
# Read bytes by process:

bpftrace -e 'tracepoint:syscalls:sys_exit_read /args->ret/ { @[comm] = sum(args->ret); }'
```

```
root@ebpf1:~/tutorial_examples# bpftrace -e 'tracepoint:syscalls:sys_exit_read /args->ret/ { @[comm] = sum(args->ret); }'
Attaching 1 probe...
^C

@[local-path-prov]: 5
@[kube-proxy]: 28
@[envoy]: 40
@[wrk:worker_5]: 96
@[wrk:worker_18]: 96
@[wrk:worker_17]: 96
@[wrk:worker_19]: 96
@[wrk:worker_40]: 96
@[wrk:worker_28]: 96
@[wrk:worker_46]: 96
@[wrk:worker_42]: 96
@[wrk:worker_11]: 96
@[wrk:worker_33]: 96
@[wrk:worker_34]: 96
@[wrk:worker_22]: 96
@[wrk:worker_32]: 96
@[wrk:worker_45]: 96
```

# Tracing: Disk Size Read By Process

```
# Trace disk size by process

bpftrace -e 'tracepoint:block:block_rq_issue { printf("%d %s %d\n", pid, comm, args->bytes); }'
```

```
root@ebpf1:~/tutorial_examples# bpftrace -e 'tracepoint:block:block_rq_issue { printf("%d %s %d\n", pid, comm, args->bytes); }'
Attaching 1 probe...
3727230 etcd 4096
1099 kworker/47:1H 0
3727230 etcd 4096
743 kworker/41:1H 0
3727230 etcd 4096
3727230 etcd 4096
3727230 etcd 4096
3727230 etcd 4096
703 kworker/35:1H 0
3727230 etcd 4096
703 kworker/35:1H 0
3727230 etcd 4096
782 kworker/14:1H 0
3727230 etcd 8192
668 jbd2/xvda2-8 16384
```

# Tracing : Block IO Latency

```
root@ebpf1:~/tutorial_examples# biolatency.bt
Attaching 4 probes...
Tracing block device I/O... Hit Ctrl-C to end.
^C


@usecs:
[64, 128)              2 |@@@@@                                               |
[128, 256)             7 |@@@@@@@@@@@@@@@@@@@@                                 |
[256, 512)             4 |@@@@@@@@@@@                                          |
[512, 1K)             14 |@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@             |
[1K, 2K)              18 |@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@|
```

# URetProbe : Bash ReadLine

#Print entered bash commands system wide
bashreadline.bt

# Tracing: Sample which CPUs are executing processes

# Tracing: Active TCP Connections

```
root@ebpf1:~/tutorial_examples# tcpconnect-bpfcc
Tracing connect ... Hit Ctrl-C to end
PID      COMM        IP SADDR          DADDR          DPORT
3726677 kubelet      4  10.244.0.1     10.244.0.14    3000
3729646 coredns      4  127.0.0.1      127.0.0.1      8080
3726677 kubelet      4  10.244.0.1     10.244.0.6     8080
3726677 kubelet      4  172.18.0.2     172.18.0.2     6443
3730090 coredns      4  127.0.0.1      127.0.0.1      8080
3729646 coredns      4  127.0.0.1      127.0.0.1      8080
3726677 kubelet      4  10.244.0.1     10.244.0.15    15021
3726677 kubelet      4  10.244.0.1     10.244.0.8     15021
3726677 kubelet      4  10.244.0.1     10.244.0.4     15021
3726677 kubelet      4  10.244.0.1     10.244.0.7     15021
3726677 kubelet      4  10.244.0.1     10.244.0.5     15021
3726677 kubelet      4  10.244.0.1     10.244.0.16    15021
3726677 kubelet      4  10.244.0.1     10.244.0.17    15021
3726677 kubelet      4  10.244.0.1     10.244.0.10    15021
3726677 kubelet      4  10.244.0.1     10.244.0.3     9090
3726677 kubelet      4  172.18.0.2     172.18.0.2     6443
3730090 coredns      4  127.0.0.1      127.0.0.1      8080
3726677 kubelet      4  127.0.0.1      127.0.0.1      2381
3729646 coredns      4  127.0.0.1      127.0.0.1      8080
3726677 kubelet      4  172.18.0.2     172.18.0.2     6443
```

# TCP Connection Latency

```
root@ebpf1:~/tutorial_examples# tcpconnlat-bpfcc
PID      COMM          IP  SADDR          DADDR            DPORT  LAT(ms)
3732907  kubelet       4   10.244.0.1     10.244.0.6        8080   0.12
3731867  kubelet       4   172.18.0.2     172.18.0.2       6443   0.09
3730754  coredns       4   127.0.0.1      127.0.0.1        8080   0.06
3732533  coredns       4   127.0.0.1      127.0.0.1        8080   0.07
3732816  kubelet       4   10.244.0.1     10.244.0.11     14269  0.12
3731712  kubelet       4   127.0.0.1      127.0.0.1       10257  0.09
3731960  kubelet       4   10.244.0.1     10.244.0.15     15021  0.10
3731959  kubelet       4   10.244.0.1     10.244.0.8      15021  0.09
3736066  kubelet       4   10.244.0.1     10.244.0.4      15021  0.09
3726870  kubelet       4   10.244.0.1     10.244.0.16     15021  0.04
3736066  kubelet       4   10.244.0.1     10.244.0.7      15021  0.08
3731960  kubelet       4   10.244.0.1     10.244.0.5      15021  0.13
3731959  kubelet       4   10.244.0.1     10.244.0.17     15021  0.06
3731959  kubelet       4   10.244.0.1     10.244.0.10     15021  0.05
3936549  kubelet       4   10.244.0.1     10.244.0.13      8181   0.09
3731959  kubelet       4   172.18.0.2     172.18.0.2       6443   0.10
3730754  coredns       4   127.0.0.1      127.0.0.1        8080   0.05
3732686  coredns       4   127.0.0.1      127.0.0.1        8080   0.07
3936549  kubelet       4   10.244.0.1     10.244.0.3       9090   0.11
3727907  kubelet       4   10.244.0.1     10.244.0.3       9090   0.12
3731959  kubelet       4   172.18.0.2     172.18.0.2       6443   0.09
```

# Tracing SSL Data

sslsniff-bpfcc

```
READ/RECV      0.579548718         curl              41531  1256
----- DATA -----
<!doctype html>
<html>
<head>
    <title>Example Domain</title>

    <meta charset="utf-8" />
    <meta http-equiv="Content-type" content="text/html; charset=utf-8" />
    <meta name="viewport" content="width=device-width, initial-scale=1" />
    <style type="text/css">
    body {
        background-color: #f0f0f2;
        margin: 0;
        padding: 0;
        font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sans", "Helvetica Neue"
----- END DATA (TRUNCATED, 792 bytes lost) -----
```
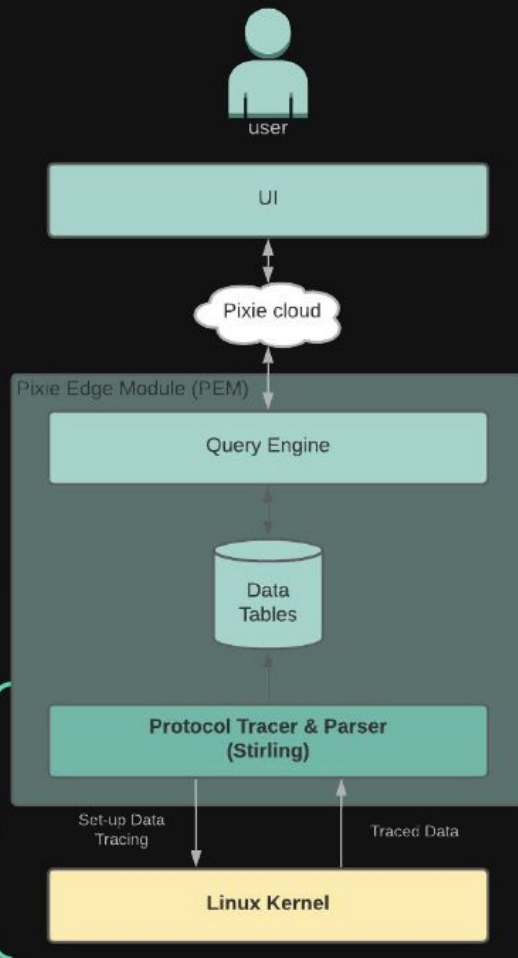
# Tracing UserLand Code

# Tracing With Uprobes: C code

```c
#include<stdio.h>
 int check(char* ip, int v){
         printf("%s\n",ip);
         return v;
 }
int main(){
     ▮printf("Hello\n");
       check("Hi", 10);
}
~
~
```

# Tracing With Uprobes:

# Tracing TLS Connections

# On SSL Function Entry

```c
// Function signature being probed:
// int SSL_write(SSL *ssl, const void *buf, int num);
int probe_entry_SSL_write(struct pt_regs* ctx) {
  uint64_t current_pid_tgid = bpf_get_current_pid_tgid();
  uint32_t pid = current_pid_tgid >> 32;

  if (pid != TRACE_PID) {
    return 0;
  }

  const char* buf = (const char*)PT_REGS_PARM2(ctx);
  active_ssl_write_args_map.update(&current_pid_tgid, &buf);

  return 0;
}
```

# On SSL Function Return

```
int probe_ret_SSL_write(struct pt_regs* ctx) {
  uint64_t current_pid_tgid = bpf_get_current_pid_tgid();
  uint32_t pid = current_pid_tgid >> 32;

  if (pid != TRACE_PID) {
    return 0;
  }

  const char** buf = active_ssl_write_args_map.lookup(&current_pid_tgid);
  if (buf != NULL) {
    process_SSL_data(ctx, current_pid_tgid, kSSLWrite, *buf);
  }

  active_ssl_write_args_map.delete(&current_pid_tgid);
  return 0;
}
```

# Tracing OpenSSL with Uprobes

# Use Cases

1. K8s Observability
2. Network Acceleration

# Use Cases

1. **K8s Observability**
2. Network Acceleration

# Pixie : Monitoring Kubernetes Clusters

# Overview

No instrumentation + low overhead ⇨ **eBPF**.

General approach:

- Capture data in kernel-space with eBPF.
- Process data in user-space (protocol parsing).
- Store data into tables for querying by user.

script: px/net_flow_graph ▾    namespace*: px-sock-shop ▾    from_entity_filter: ▾    to_entity_filter: ▾    throughput_filter: 0.0 ▾    start_time: -5m ▾

## Net Flow Graph

rabbitmq.px-sock-shop.svc.cluster.local

px-sock-shop/rabbitmq-6f1d55c844f-qf8lh

kube-dns.kube-system.svc.cluster.local

user.px-sock-shop.svc.cluster.local

**Automatic Service Maps**

px-sock-shop/carts-5fc45568c4-bvwbs

shipping.px-sock-shop.svc.cluster.local

bytes_total: 103.5 KB/s
bytes_sent: 40.1 KB/s
bytes_recv: 63.4 KB/s

px-sock-shop/orders-77c57c89dc-p47gw

carts.px-sock-shop.svc.cluster.local

carts-db.px-sock-shop.svc.cluster.local

px-sock-shop/queue-master-bd5f62d5-45t4k

px-sock-shop/shipping-745b0d8755-g8h2k

10.109.136.241

ENABLE HIERARCHY

## Table

| FROM_ENTITY ^ | TO_ENTITY ^ | BYTES_SENT ^ | BYTES_RECV ^ | BYTES_TOTAL ^ |
|---|---|---|---|---|
| px-sock-shop/carts-5fc45568c4-bvwbs | carts-db.px-sock-shop.svc.cluster.local | 40.1 KB/s | 63.4 KB/s | 103.5 KB/s |
| px-sock-shop/carts-5fc45568c4-bvwbs | kube-dns.kube-system.svc.cluster.local | 38.3 B/s | 91.7 B/s | 130 B/s |
| px-sock-shop/orders-77c57c89dc-p47gw | shipping.px-sock-shop.svc.cluster.local | 3.8 KB/s | 2.3 KB/s | 6.1 KB/s |

# Where to Trace the Data?

Many options in the software stack:

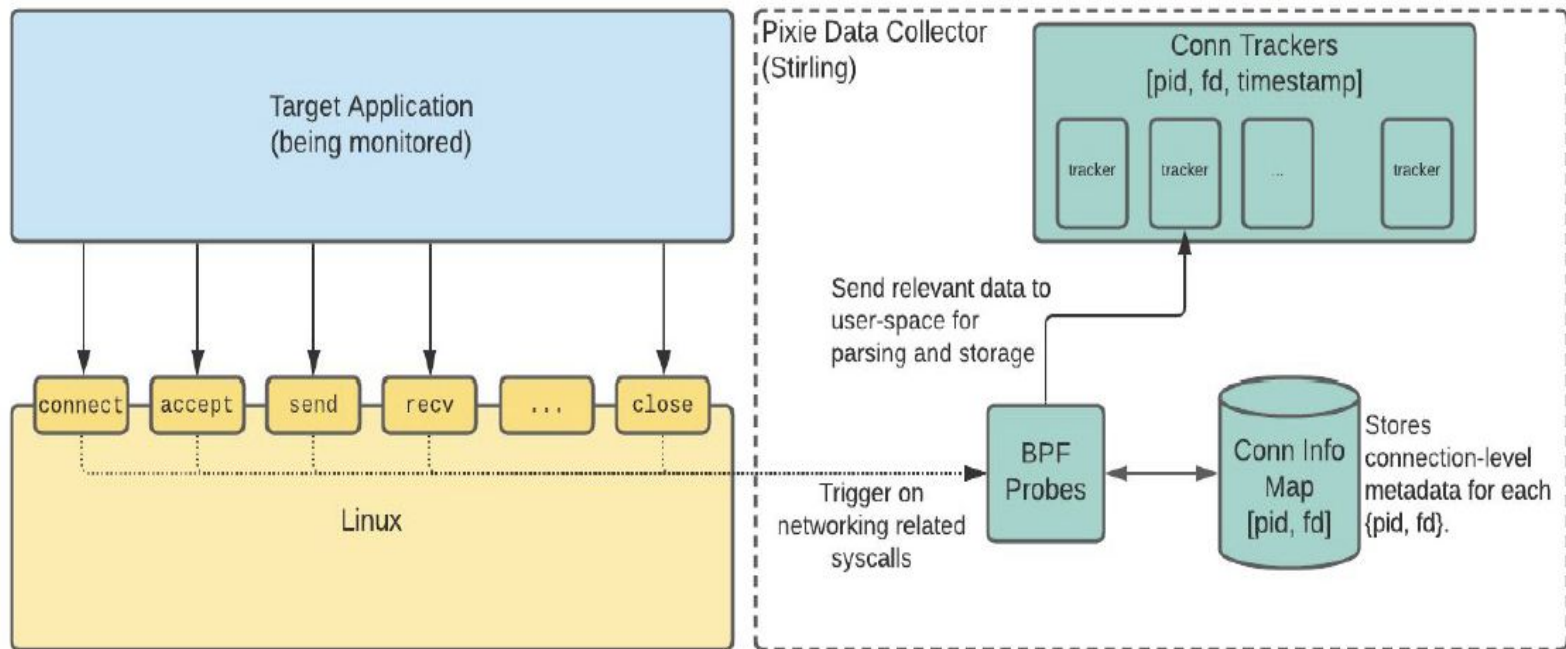We preferred tracing as close to the application layer as possible.

Tracing

| | |
|---|---|
| **Application** | |
| ↓ | |
| **HTTP Library** | ◁┈┈┈┈ Protocol libra |
| ↓ | |
| **SSL** | |
| ↓ | |
| **Linux API** | ◁┈┈┈┈ Syscall kprob |
| **Socket Layer** | |
| **TCP Layer** | |
| **IP Layer** | ◁┈┈┈┈ libpcap probe |
| **Network Driver** | ◁┈┈┈┈ XDP probes |

# Approaches Compared

| | protocol library uprobes | syscall kprobes | libpcap/XDP |
|---|---|---|---|
| Tracing overhead | Low | Low | Low |
| Scalability & Stability | Uprobes per library, Probe targets may change | High | High |
| Parsing effort | None | Protocol parsing | Packet processing & protocol parsing |
| SSL tracing | Cleartext available | Data encrypted | Data encrypted |

We chose to use syscall kprobes on functions such as send() & recv().

- Rationale: close to the application layer, but stable API.

# Architecture

# Architecture

## 1 - Setup probes on network related syscalls.

# Architecture

## 2 - Record connection metadata in BPF maps.

# Architecture

**3 - Infer protocol with basic rule-based classification as a simple filter.
Transfer connection information and data through two perf buffers.**

# Architecture

**4 - Track connections in user-space with ConnTrackers.**
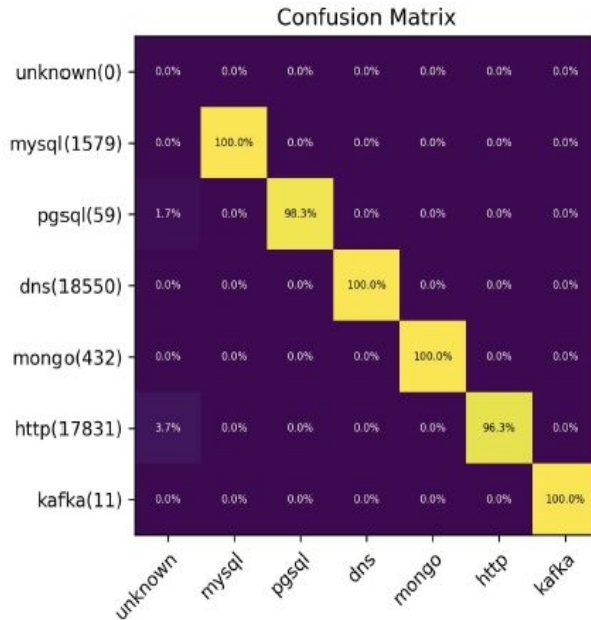   **Parse ConnTracker data into structured messages.**

# eBPF-Side Protocol Inference

To filter data transfers to user-space, we apply protocol inference in BPF.

- Just a filter: False positives are okay.

- Example for HTTP:

```
static __inline enum MessageType infer_http_message(const char* buf, size_t count) {
  ...

if (buf[0] == 'H' && buf[1] == 'T' && buf[2] == 'T' && buf[3] == 'P') {
        return kResponse;
}
  ...
```

Likelihood that our inference eventually identifies the right protocol



Confusion Matrix

# Pluggable Protocol Parsers

## Architecture consists of pluggable protocol parsers



### Supported Protocols List

HTTP
MySQL
Postgres
Redis
Cassandra
Kafka
NATS
DNS
gRPC*

*gRPC is traced with dedicated uprobes
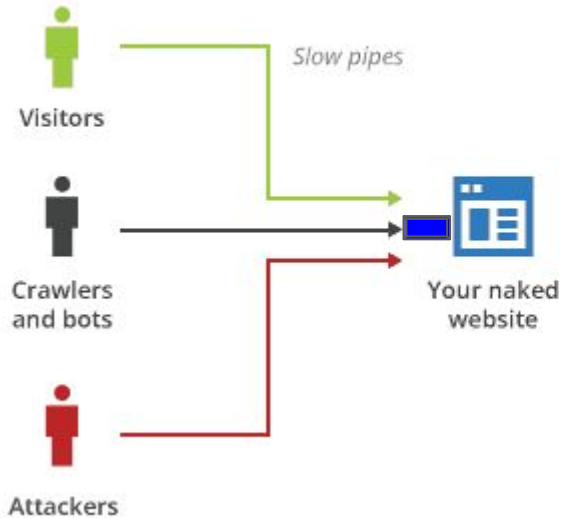
# Bonus Material

# Use Cases

1. K8s Observability
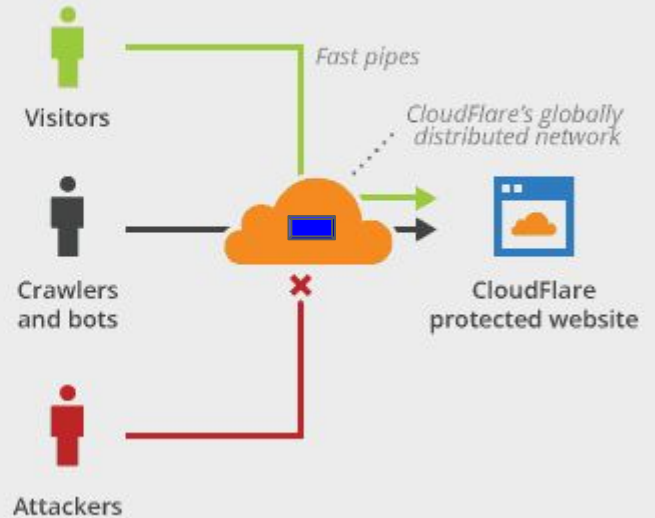2. **Network Acceleration**

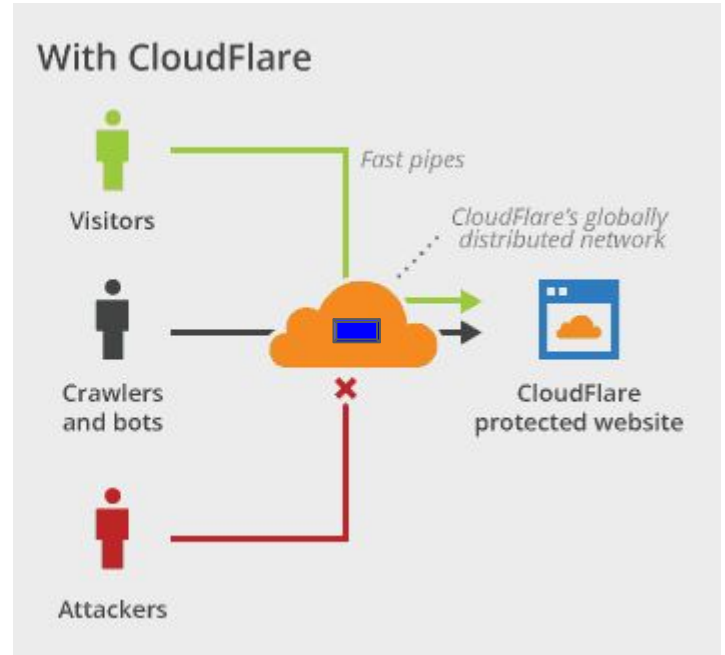# Content Distribution Networks

# Cloudflare CDN

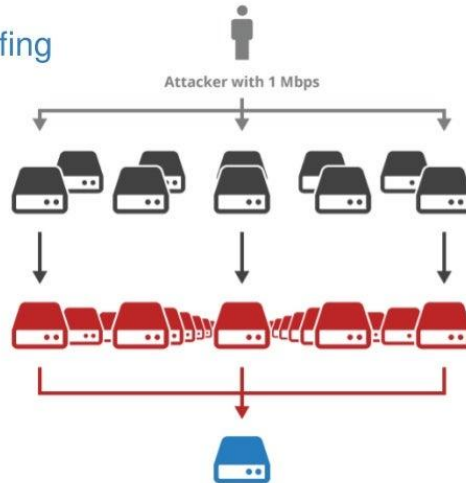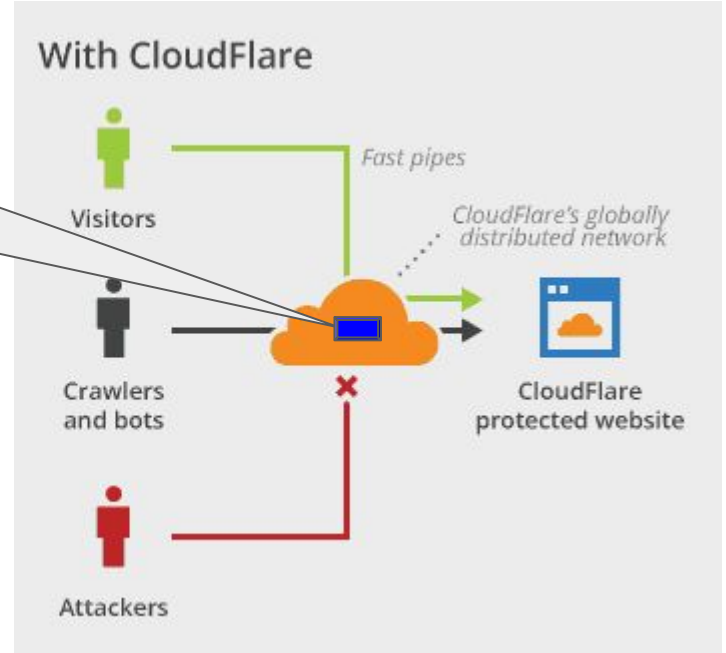# Cloudflare CDN

## 300Gbps+ of DDoS attack traffic
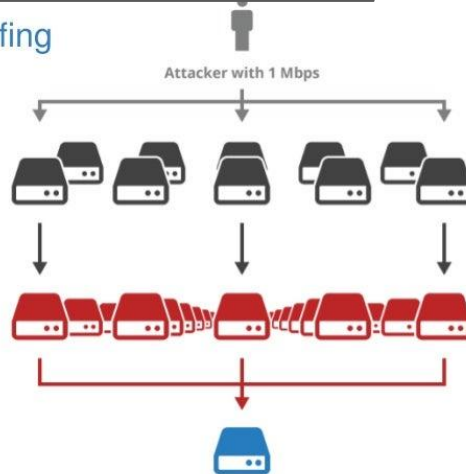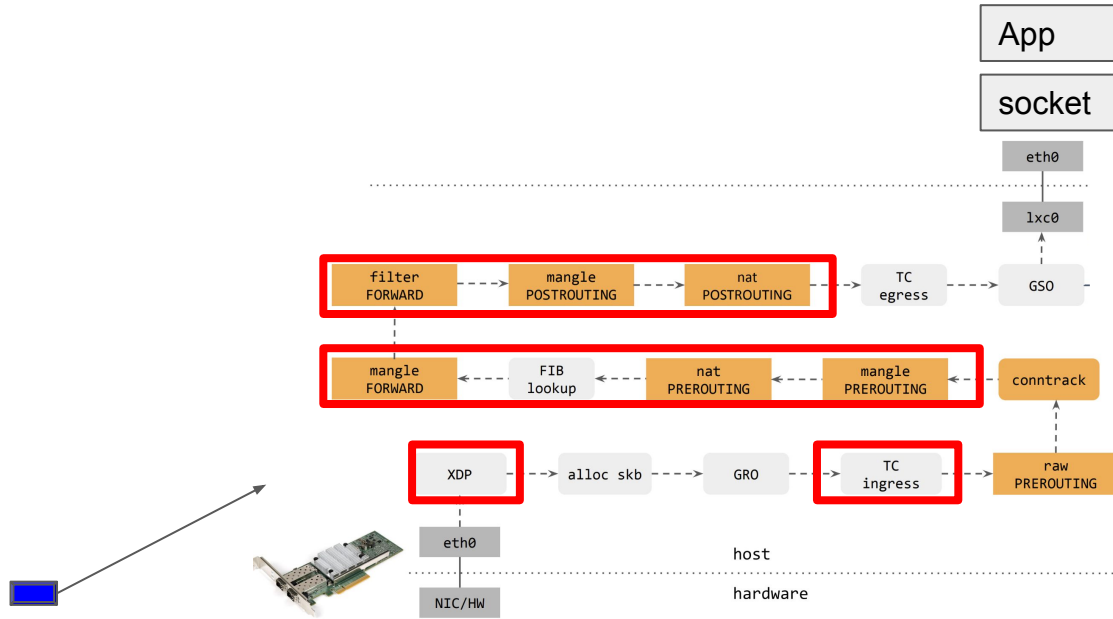
1   laptop
+   5-7 compromised servers
+   3 networks which allow spoofing
+   9Gbps of DNS requests to
+   0.1% of all open resolvers
-----------------------------------------

=   300Gbps of DDoS traffic

Attacker with 1 Mbps

### With CloudFlare

Visitors

Fast pipes

CloudFlare's globally distributed network

Crawlers and bots

CloudFlare protected website

Attackers

# Cloudflare CDN

300Gbps+ ~~~~~~ ffic

1. Edge of infra.
2. HW/SW ?
3. Sophisticated Detection
4. Drop pkts.

1   laptop
+   5-7 compromised ~~~~~~
+   3 networks which allow spoofing
+   9Gbps of DNS requests to
+   0.1% of all open resolvers
----------------------------------------

=   300Gbps of DDoS traffic

Attacker with 1 Mbps

**With CloudFlare**

Visitors

Fast pipes

CloudFlare's globally distributed network

Crawlers and bots
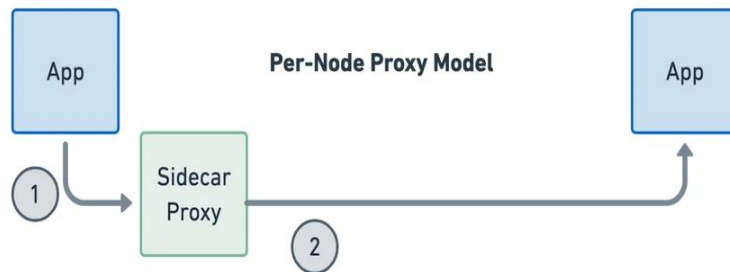
CloudFlare protected website

Attackers

CLOUDFLARE

# Packet Dropping for DDoS Mitigation
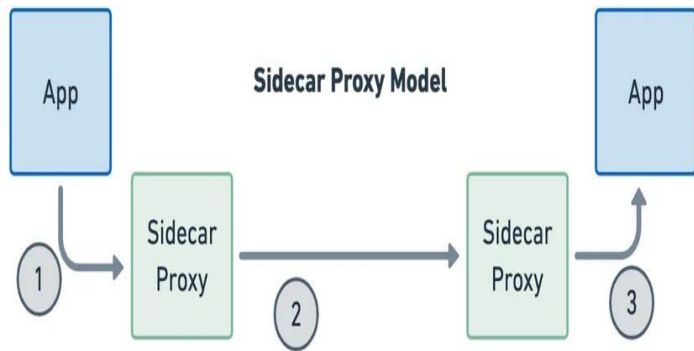


From Here: How to drop 10 million packets per second (cloudflare.com)

Packet dropping performance

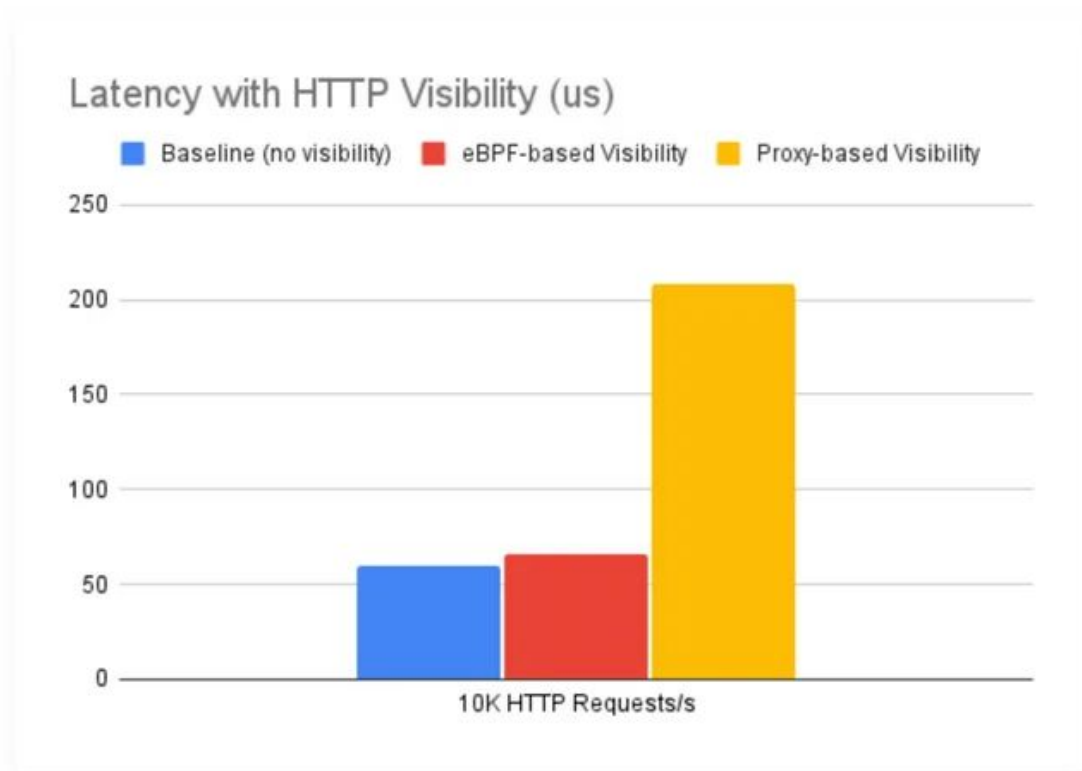From Here: [How to drop 10 million packets per second (cloudflare.com)](https://cloudflare.com)

# De-duplicate Network Stack Traversal with eBPF SockMap

# Low Cost Service Mesh with eBPF

# Conclusion

a) Where else is eBPF used : Storage, Security, DBMS,Scheduling etc.

b) What next -- eBPF startups ma

**Isovalent Raises $40M in Series B Funding**

## Groundcover lands $20M to he
companies monitor app perfor

USA
Published on September 7, 2022

Kyle Wiggers  @kyle_l_wiggers  /  6:30 PM GMT+5:30 • September 14, 2022

eBPF-powered Cilium Networking

Isovalent, a Mountain View, CA-based company behind open source technologies Cilium and eBPF, closed a $40M Series B funding round.

The round was led by Thomvest Ventures with participation from M12 (Microsoft's Venture Fund) and Grafana Labs, which joined Google and Cisco as existing strategic investors in the company, as well as Andreessen Horowitz, Mango Capital, and Mirae Asset Capital.

# OPENED Tool for Managing eBPF Heterogeneity

[Microservices Observatory (microserviceobservatory.github.io)](microserviceobservatory.github.io)

Theophilus A. Benson      Brown University    tab@cs.brown.edu
Palanivel Kodeswaran    IBM Research    palani.kodeswaran@in.ibm.com
Sayandeep Sen      IBM Research    sayandes@in.ibm.com

# eBPF Programs are Monoliths

## One Off Programs



## Observability

## Complex Systems(s)



## Network Functions
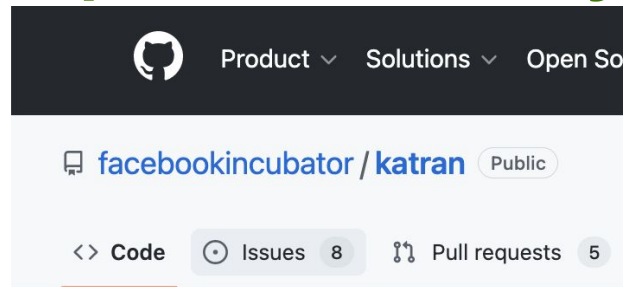
*Code from a Katran function

# Implications of Monolith on Developer Productivity

Developing a new program

Find sub functionality on GitHub
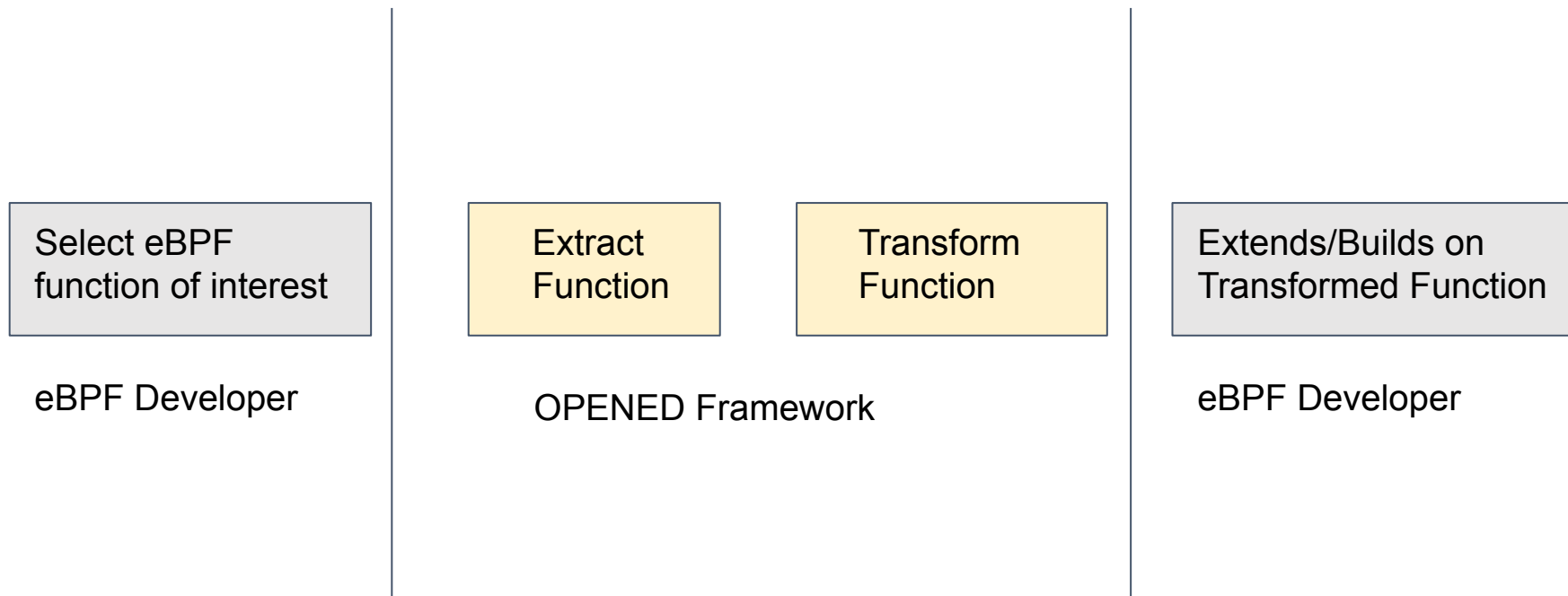
Extracting and reusing functionality is non-trival

Product ✓   Solutions ✓   Open So

🗔 facebookincubator / katran   Public

<> Code   ⊙ Issues   8   ⇄ Pull requests   5

Step 1: Extract lines

Step 2: Identify + Extract Deps

Surprise Step 3: Rewrite for your target hookpoint

# The OPENED Vision

# **OPENED Vision:** Reduce time to new functionality development

- Automated extraction of relevant code

- Automated transformation of code
  - Enable moving code between hook-points
  - Enable moving code between programs

- Developer-first automation
  - Extraction + Transformation guided by developer choices