



**KubeCon**



**CloudNativeCon**

---

**Europe 2023**

---



TiKV

# Past, Present, and Future of eBPF in Cloud Native Observability

*Natalie Serrino, New Relic  
Frederic Branczyk, Polar Signals*



KubeCon

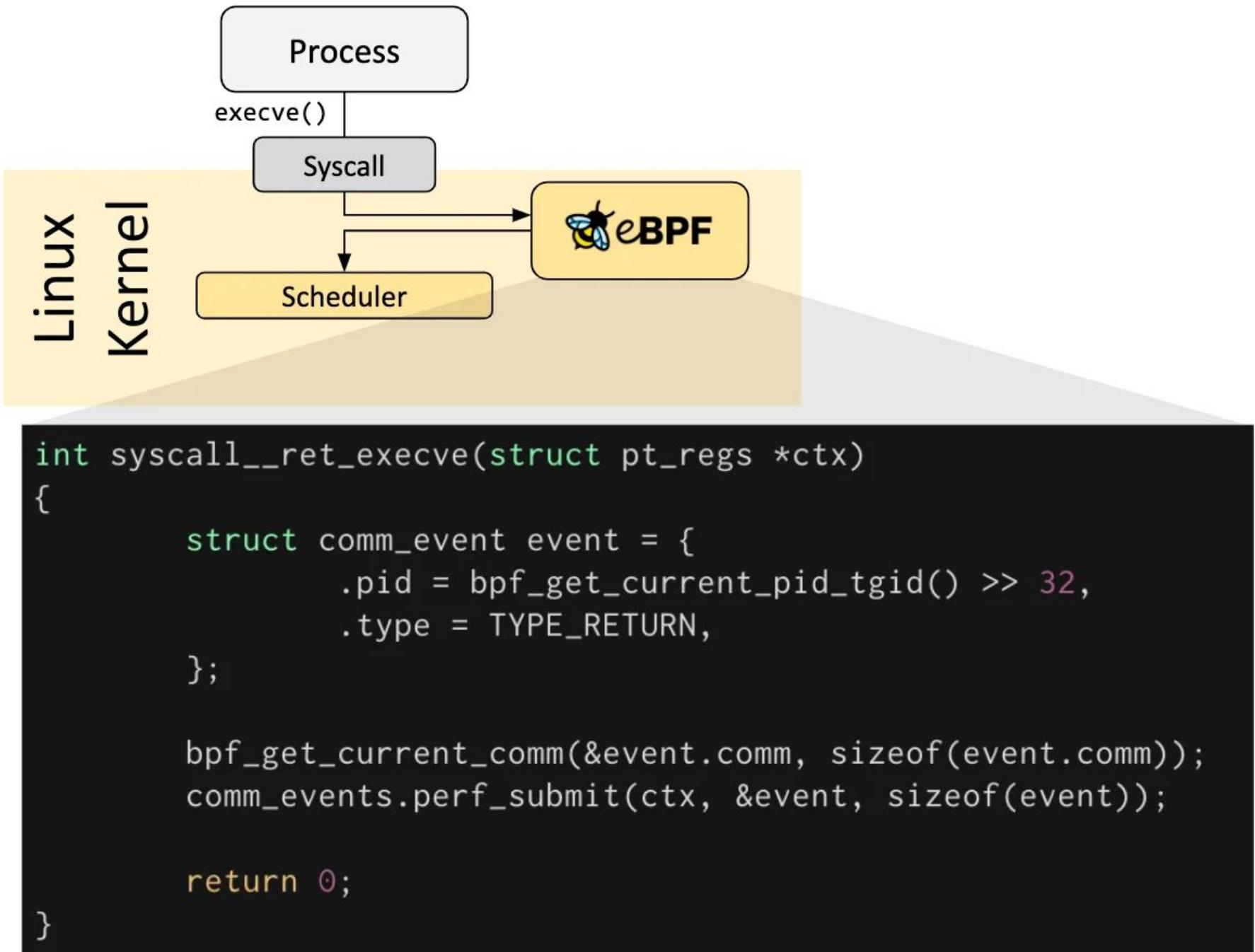


CloudNativeCon

Europe 2023



# Intro to eBPF



# Hooks

- Pre-defined (syscalls, kernel-functions, tracepoints, network events...)
- Custom hooks
  - kprobe
  - uprobe
  - perf\_event

# Compiling an eBPF program



```
clang --target=x86_64-unknown-linux-elf main.c
```

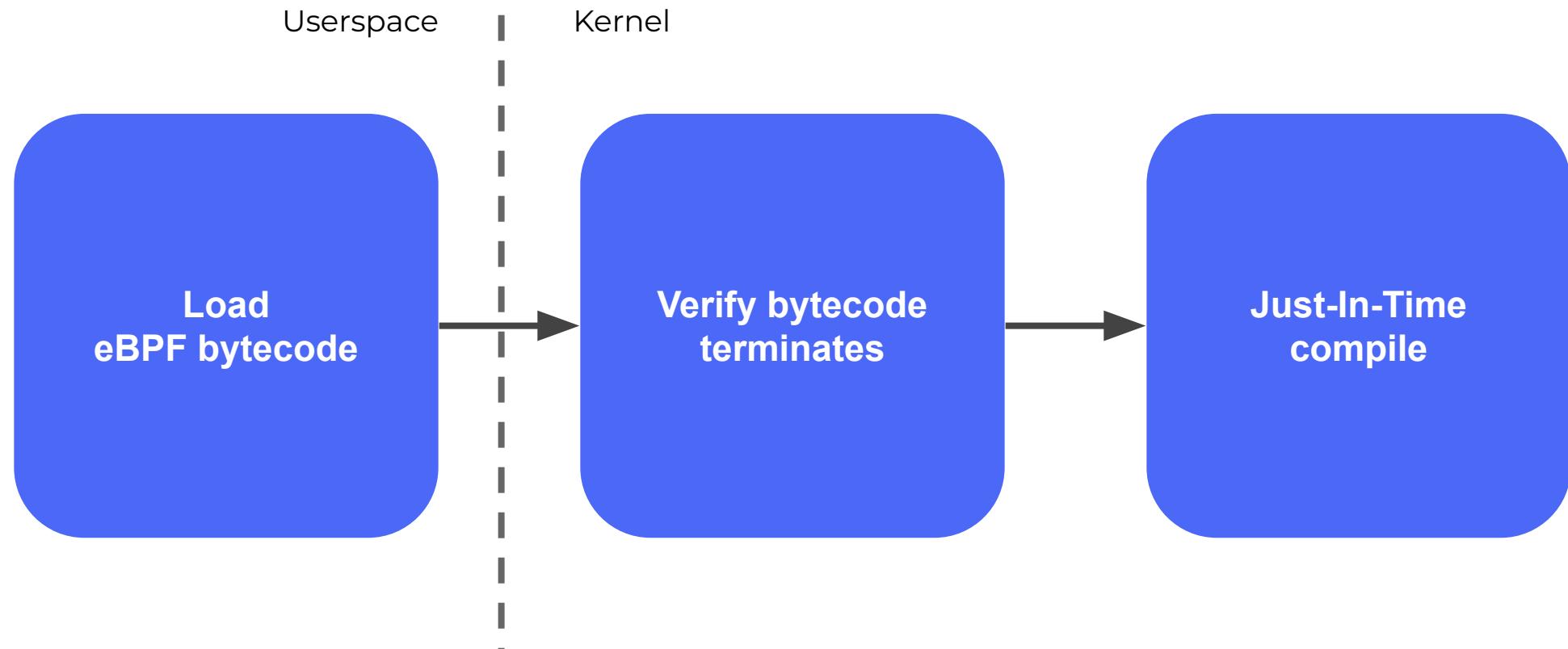
# Compiling an eBPF program

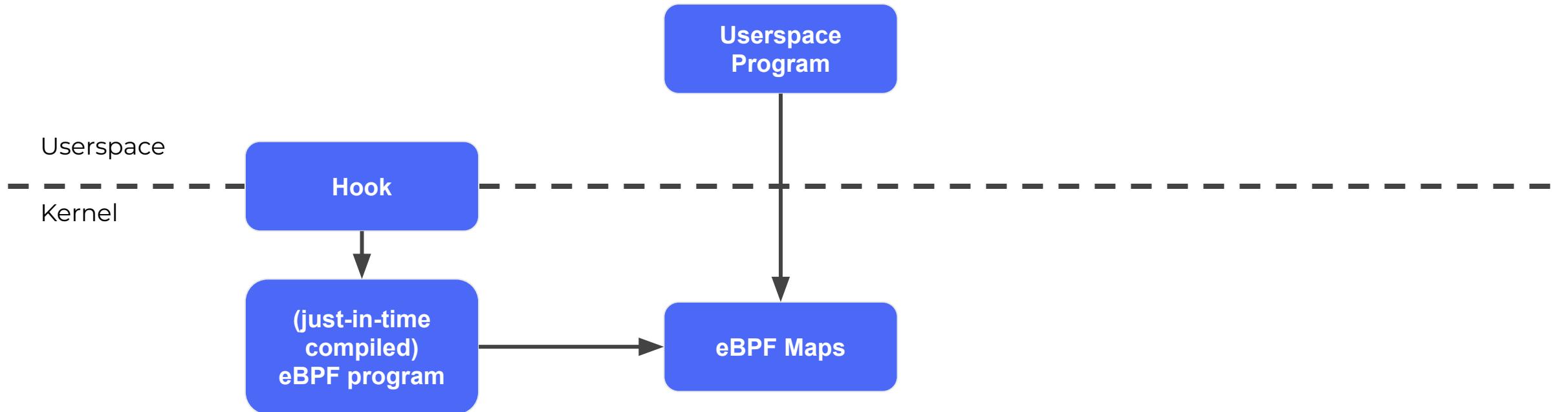
---



```
clang --target=bpf main.c
```

→ **eBPF byte code**





# eBPF has been a buzzword

Let's resolve some misconceptions

# eBPF

That networking thing right?

**FALSE**

# eBPF == all languages are supported

Demos are easy, actual wide language support is very hard, and not widespread!

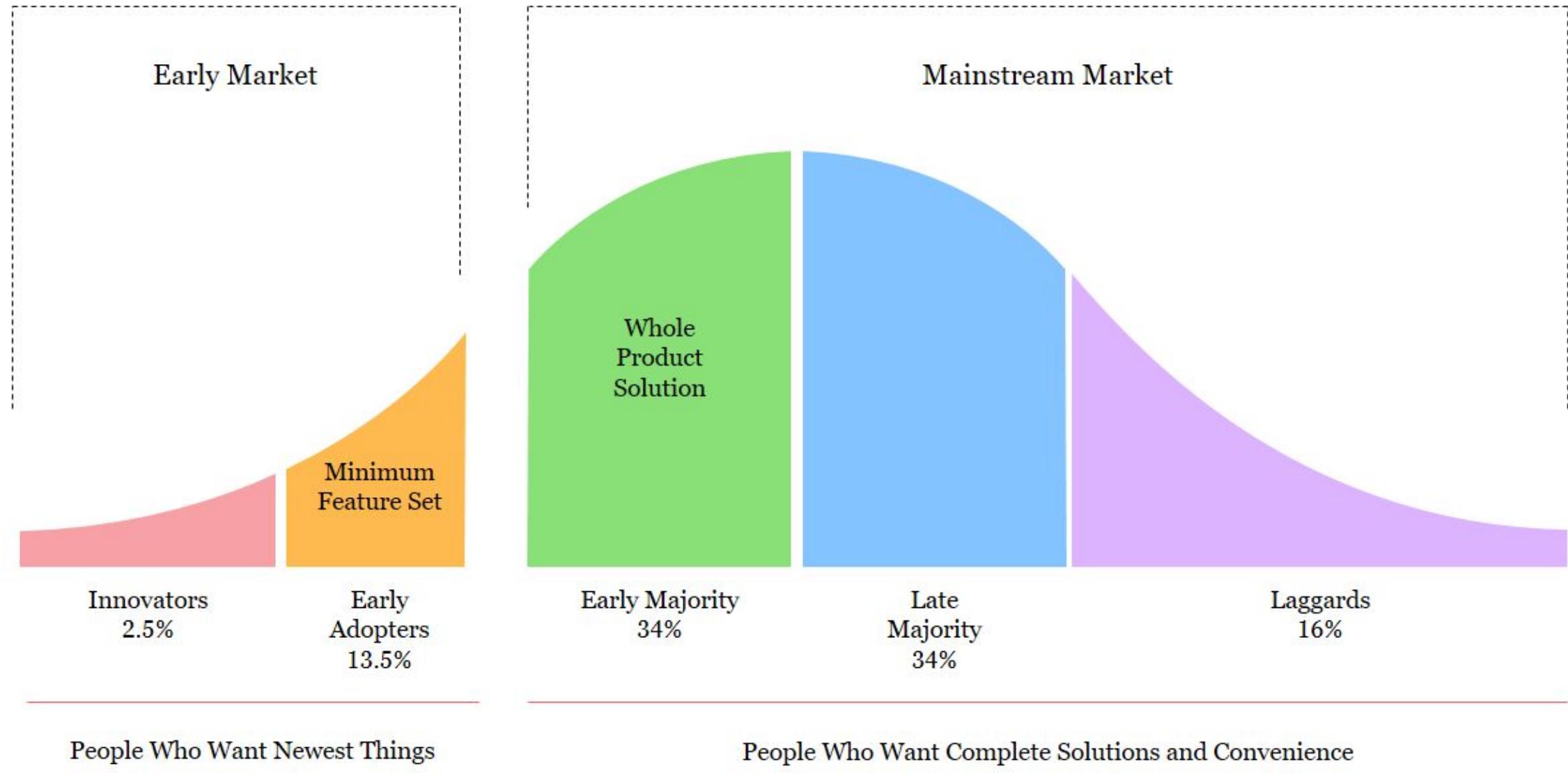
The opposite:  
eBPF means interpreted languages  
are impossible

Also incorrect. It's difficult, but possible.

# Past

History and initial use cases of eBPF

# Crossing The Chasm



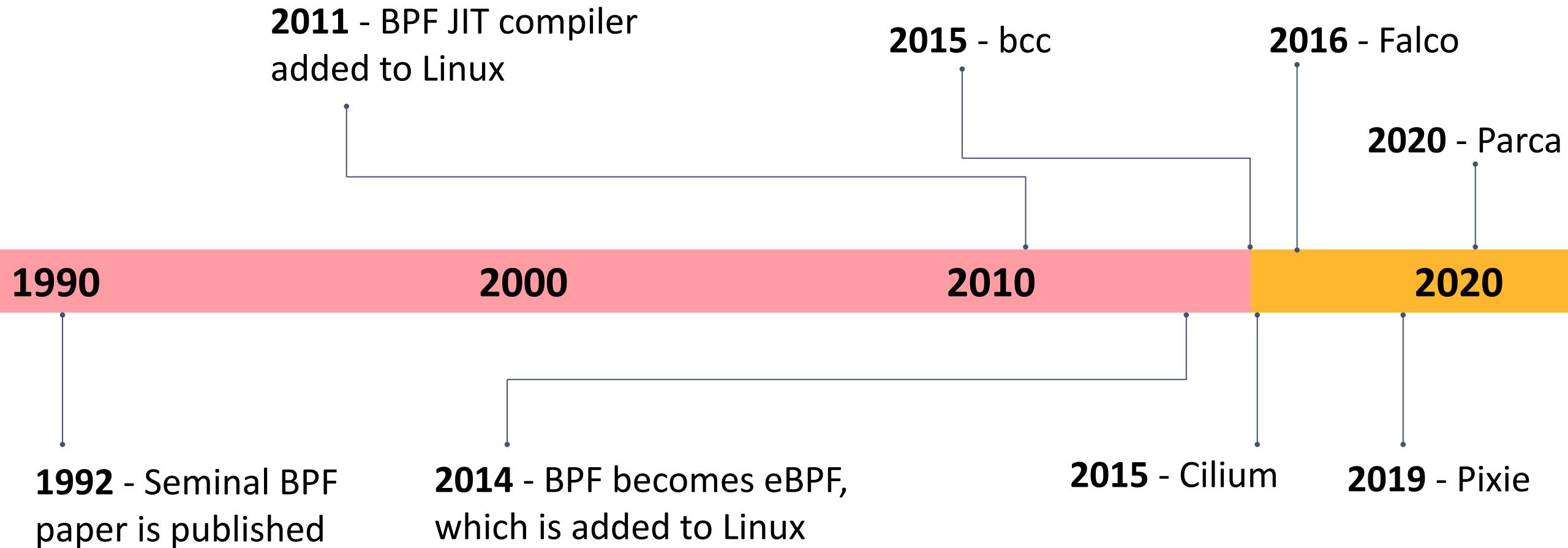
*Image credit: Geoffrey A Moore*

# Crossing The Chasm



Image credit: Geoffrey A Moore

# How did we get here?



# What was the Berkeley Packet Filter and why was it created?

**When:** Late 80s/early 90s

**Who:** Steven McCanne and Van Jacobson,  
Lawrence Berkeley Laboratory

**Use case:** Diagnosing and fixing network problems  
(packet filtering)

**Motivation:** Allow user-space programs to define rules for interacting with raw, unprocessed packets

**Problem:** Previous approaches were slow

**Solution:** a new packet filtering architecture, which allows the userspace programs to define packet rules that execute in the kernel

## The BSD Packet Filter: A New Architecture for User-level Packet Capture

Steven McCanne & Van Jacobson – Lawrence Berkeley Laboratory<sup>1</sup>

### ABSTRACT

Many versions of Unix provide facilities for user-level packet capture, making possible the use of general purpose workstations for network monitoring. Because network monitors run as user-level processes, packets must be copied across the kernel/user-space protection boundary. This copying can be minimized by deploying a kernel agent called a *packet filter*, which discards unwanted packets as early as possible. The original Unix packet filter was designed around a stack-based filter evaluator that performs sub-optimally on current RISC CPUs. The BSD Packet Filter (BPF) uses a new, register-based filter evaluator that is up to 20 times faster than the original design. BPF also uses a straightforward buffering strategy that makes its overall performance up to 100 times faster than Sun's NIT running on the same hardware.

### Introduction

Unix has become synonymous with high quality networking and today's Unix users depend on having reliable, responsive network access. Unfortunately, this dependence means that network trouble can make it impossible to get useful work done and increasingly users and system administrators find that a large part of their time is spent isolating and fixing network problems. Problem solving requires appropriate diagnostic and analysis tools and, ideally, these tools should be available where the problems are – on Unix workstations. To allow such tools to be constructed, a kernel must contain some facility that gives user-level programs access to raw, unprocessed network traffic [7]. Most of today's workstation operating systems contain such a facility, e.g., NIT[10] in SunOS, the Ultrix Packet Filter [2] in DEC's Ultrix and Snoop in SGI's IRIX.

These kernel facilities derive from pioneering work done at CMU and Stanford to adapt the Xerox Alto 'packet filter' to a Unix kernel[8]. When completed in 1980, the CMU/Stanford Packet Filter, CSPF, provided a much needed and widely used facility. However on today's machines its performance, and the performance of its descendants, leave much to be desired – a design that was entirely appropriate for a 64KB PDP-11 is simply not a good match to a 16MB Sparstation 2. This paper describes the BSD Packet Filter, BPF, a new kernel architecture for packet capture. BPF offers substantial performance improvement over existing packet capture facilities – 10 to 150 times faster than Sun's NIT and 1.5 to 20 times faster than CSPF on the

same hardware and traffic mix. The performance increase is the result of two architectural improvements:

- BPF uses a re-designed, register-based 'filter machine' that can be implemented efficiently on today's register based RISC CPU. CSPF used a memory-stack-based filter machine that worked well on the PDP-11 but is a poor match to memory-bottlenecked modern CPUs.
- BPF uses a simple, non-shared buffer model made possible by today's larger address spaces. The model is very efficient for the 'usual cases' of packet capture.<sup>2</sup>

In this paper, we present the design of BPF, outline how it interfaces with the rest of the system, and describe the new approach to the filtering mechanism. Finally, we present performance measurements of BPF, NIT, and CSPF which show why BPF performs better than the other approaches.

### The Network Tap

BPF has two main components: the network tap and the packet filter. The network tap collects copies of packets from the network device drivers and delivers them to listening applications. The filter decides if a packet should be accepted and, if so, how much of it to copy to the listening application.

Figure 1 illustrates BPF's interface with the rest of the system. When a packet arrives at a network interface the link level device driver normally sends it up the system protocol stack. But when BPF is listening on this interface, the driver first

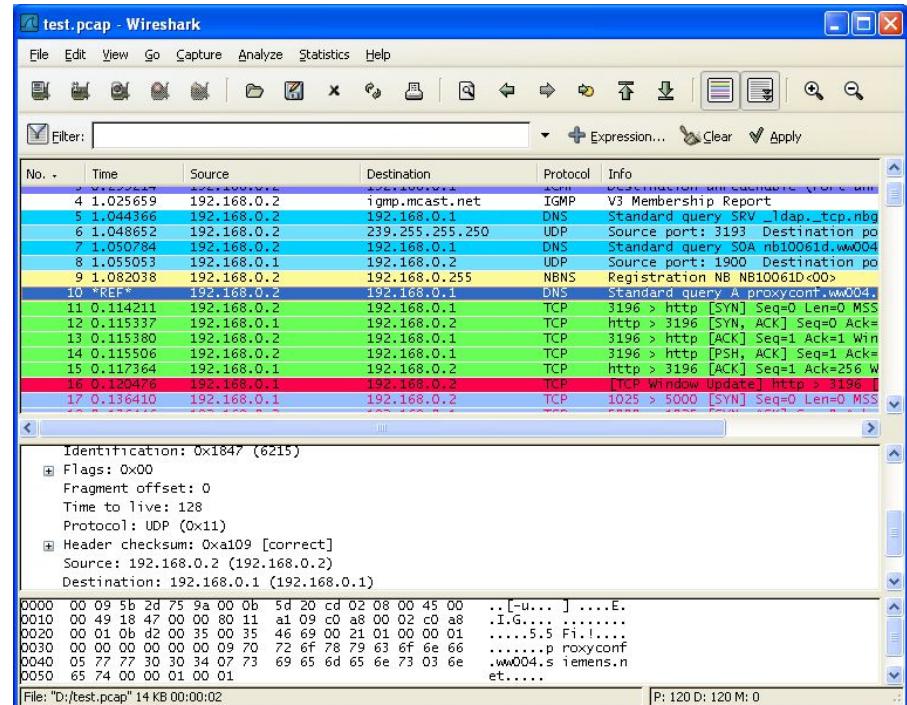
<sup>1</sup>This work was supported by the Director, Office of

<sup>2</sup>As opposed to, for example, the AT&T STREAMS

# What happened with BPF from 1992-2011?

- 🏆 State of the art for packet filtering
- 📡 Integrated into tcpdump and Wireshark
- ⌚ BPF JIT added to Linux in 2011

*Then, one day in 2014 ...*



# From BPF to eBPF

The first version of eBPF was added to Linux in 2014.

It built upon BPF in some game-changing ways.

So then all hell broke loose.



*eBPF programs are written in C*



*eBPF can invoke kernel functions*



*eBPF can view and edit raw kernel memory  
(safely, with the help of the new verifier)*

# eBPF exists. Now what?

*It's 2014, eBPF is brand new and super powerful. But there are some obstacles:*

- Writing eBPF programs is technically difficult
- It's still relatively obscure (for now)
- People mostly think of it as a packet filter



*Enter: libbpf, bpftools, and bcc*

- Writing eBPF programs is now more accessible
- More people try it
- Use cases beyond networking are gaining awareness



# An explosion of adoption

Hacker News Search ebpf algolia Settings

Search Stories by Popularity for All time 494 results (0.009 seconds)

How io\_uring and eBPF Will Revolutionize Programming in Linux ([https://www.scylladb.com/2020/05/05/how-io\\_uring-and-ebpf-will-revolutionize-programming-in-linux/](https://www.scylladb.com/2020/05/05/how-io_uring-and-ebpf-will-revolutionize-programming-in-linux/))  
709 points | harpoeder | 2 years ago | 320 comments

New Relic to open-source Pixie's eBPF observability platform (<https://blog.pixielabs.ai/pixie-new-relic/>)  
345 points | htrosli | 2 years ago | 63 comments

How Netflix uses eBPF flow logs at scale for network insight (<https://netflixtechblog.com/how-netflix-uses-ebpf-flow-logs-at-scale-for-network-insight-e3ea997dca96>)  
310 points | el\_duderino | 2 years ago | 34 comments

eBPF on Windows (<https://github.com/microsoft/ebpf-for-windows>)  
294 points | prasedym | 2 years ago | 165 comments

EBPFsnitch: An eBPF based Linux Application Firewall (<https://github.com/harpoeder/ebpfsnitch>)  
290 points | harpoeder | 2 years ago | 70 comments

eBPF Is Awesome (<https://filipnikolovski.com/posts/ebpf/>)  
268 points | filipn | 2 years ago | 44 comments

eBPF will help solve service mesh by getting rid of sidecars (<https://isovalent.com/blog/post/2021-12-08-ebpf-servicemesh>)  
237 points | tgraf | 1 year ago | 108 comments

Show HN: Credentials dumper for Linux using eBPF (<https://github.com/citronneur/pamspy>)  
235 points | citronneur | 10 months ago | 50 comments

Notes on BPF and eBPF (<https://jvns.ca/blog/2017/06/28/notes-on-bpf---ebpf/>)  
216 points | mlerner | 1 year ago | 43 comments

How to add eBPF observability to your product (<https://brendangregg.com/blog/2021-07-03/how-to-add-bpf-observability.html>)  
201 points | mrry | 2 years ago | 31 comments

eBPF – The Future of Networking and Security (<https://cilium.io/blog/2020/11/10/ebpf-future-of-networking/>)  
195 points | genbit | 2 years ago | 36 comments

eBPF is turning the Linux kernel into a microkernel ([https://docs.google.com/presentation/d/1AcB4x7JCWET0ysDr0gsX-EIdQSTyBtmi6OAW7bE0jm0/edit#slide=id.g704abb5039\\_2\\_106](https://docs.google.com/presentation/d/1AcB4x7JCWET0ysDr0gsX-EIdQSTyBtmi6OAW7bE0jm0/edit#slide=id.g704abb5039_2_106))  
194 points | youquan | 3 years ago | 89 comments

GCC eBPF for Linux port has landed (<https://gcc.gnu.org/ml/gcc-patches/2019-08/msg01987.html>)  
167 points | edelsohn | 4 years ago | 61 comments

How I ended up writing opensnoop in pure C using eBPF (<https://bolinfest.github.io/opensnoop-native/>)  
158 points | aberoram | 5 years ago | 10 comments

eBPF Can't Count? (<https://blog.cloudflare.com/ebpf-cant-count/>)  
154 points | jgraham | 4 years ago | 29 comments

Use of eBPF in CPU Scheduler (<https://linuxplumbersconf.org/event/11/contributions/954/>)  
147 points | marcodiego | 2 years ago | 26 comments

Learn EBPF Tracing: Tutorial and Examples (<https://www.brendangregg.com/blog/2019-01-01/learn-ebpf-tracing.html>)  
147 points | knoxa2511 | 4 years ago | 20 comments

Packet, where are you? – eBPF-based Linux kernel networking debugger (<https://github.com/cilium/pwru>)  
143 points | timetogo | 12 days ago | 19 comments

## Adoption [ edit ]

eBPF has been adopted by a number of large-scale production users, for example:

- Meta uses eBPF through their Katran layer 4 load-balancer for all traffic going to facebook.com[42][43][44][29]
- Google uses eBPF in GKE, developed and uses BPF LSM to replace audit and it uses eBPF for networking[27][50]
- Cloudflare uses eBPF for load-balancing and DDoS protection and security enforcement[48][49][50][51][52]
- Netflix uses eBPF for fleet-wide network observability and performance diagnosis[53][54]
- Dropbox uses eBPF through Katran for layer 4 load-balancing[55]
- Android uses eBPF for NAT46 and traffic monitoring[56][57][58]
- Alibaba uses eBPF for Kubernetes Pod load-balancing[59]
- Datadog uses eBPF for Kubernetes Pod networking and security enforcement[60][61][62]
- Trip.com uses eBPF for Kubernetes Pod networking[63][64]
- Microsoft ported eBPF and XDP to Windows[65][66][67]
- Seznam uses eBPF through Cilium for layer 4 load-balancing[68]
- CapitalOne uses eBPF for Kubernetes Pod networking[69]
- Apple uses eBPF for Kubernetes Pod security[70]
- Sky uses eBPF for Kubernetes Pod networking[71]
- Walmart uses eBPF for layer 4 load-balancing[72][73]
- Huawei uses eBPF through their DIGLIM secure boot system[74]
- Ikea uses eBPF for Kubernetes Pod networking[75]



# Mini Demo #1

Request tracing with Pixie

**It was supposed to be so easy!**

# Present

The present of eBPF

### 3. uprobes/uretprobes: Dynamic Tracing, User-Level

To list available uprobes, you can use any program to list the text segment symbols from a binary, such as `objdump` and `nm`. For example:

```
# objdump -tT /bin/bash | grep readline
0000000007003f8 g    D0 .bss  0000000000000004 Base      rl_readline_state
000000000499e00 g    DF .text  00000000000001c5 Base      readline_internal_char
0000000004993d0 g    DF .text  0000000000000126 Base      readline_internal_setup
00000000046d400 g    DF .text  000000000000004b Base      posix_readline_initialize
00000000049a520 g    DF .text  0000000000000081 Base      readline
[...]
```

This has listed various functions containing "readline" from `/bin/bash`. These can be instrumented using `uprobes` and `uretprobes`.

Examples:

```
# bpftrace -e 'uretprobe:/bin/bash:readline { printf("read a line\n"); }'
Attaching 1 probe...
read a line
read a line
read a line
^C
```

# Reality



```
1 $ nm /usr/sbin/mariadb  
2 nm: /usr/sbin/mariadb: no symbols
```

Binaries don't typically have symbols

# debuginfod



```
1 $ file /usr/sbin/mariadb
2 /usr/sbin/mariadb: ELF 64-bit LSB pie executable, ARM aarch64, version 1 (GNU/Linux),
  dynamically linked, interpreter /lib/ld-linux-aarch64.so.1,
  BuildID[sha1]=a7dc33714ad2d574f324cdb99b1dc14fc0613d8, for GNU/Linux 3.7.0, stripped
3
4 $ curl https://debuginfod.elfutils.org/buildid/a7dc33714ad2d574f324cdb99b1dc14fc0613d8/debuginfo
5 #      https://debuginfod.elfutils.org/buildid/<build-id>/debuginfo
```

# clients

debuginfod client-side support is under construction or already available in a variety of binary-related utilities. We summarize current upstream status [2023-03] below. Note that distros may lag behind upstream developments.

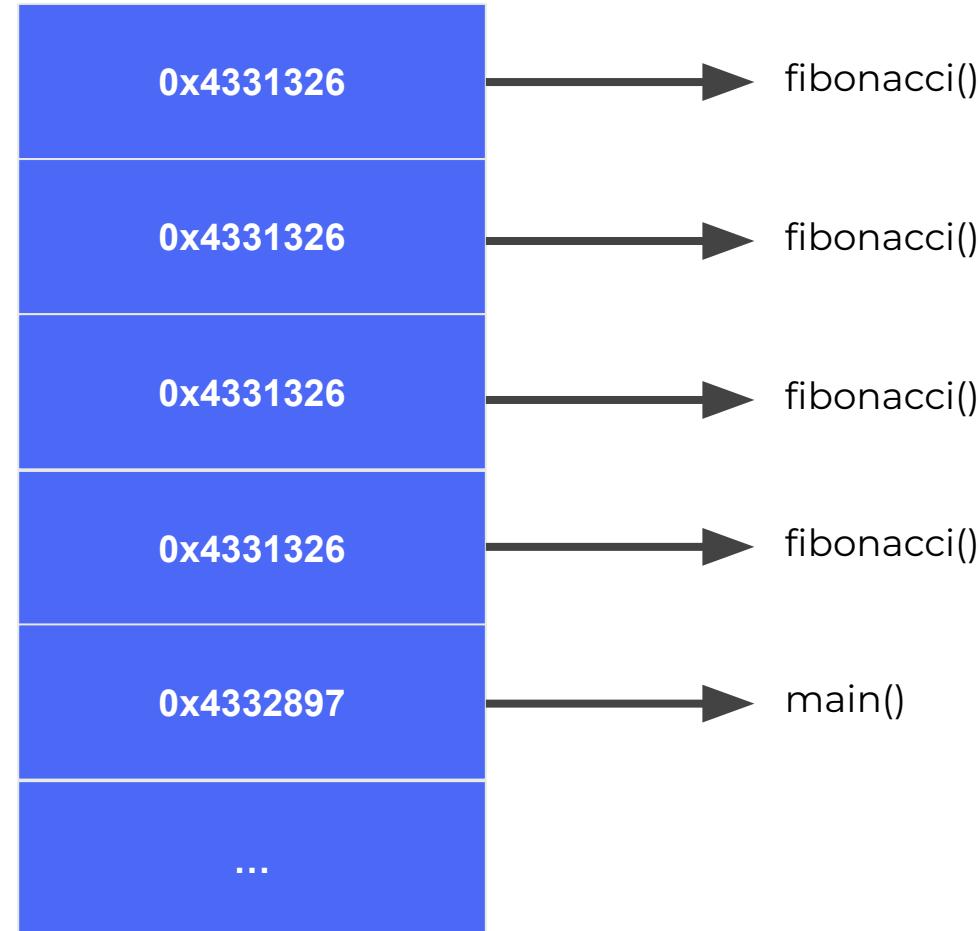
tool	status
elfutils	released in version 0.178, 2019-11
systemtap	automatic via elfutils
dwarves	automatic via elfutils
dwgrep	automatic via elfutils
ltrace	automatic via elfutils
libabigail	automatic via elfutils
binutils	released in version 2.34, 2020-02
gdb	released in version 10.1, 2020-10
dyninst	released in version 11.0 2021-04
valgrind	released in version 3.17.0, 2021-03
annocheck	released in version 9.03, 2020-01
delve	released in version 1.7.2, 2021-09
llvm	symbolizer merged, server merged, lldb help wanted, see also
bpftrace	released in version 0.21.0, 2021-07
perf	released in linux 5.10, 2021-01
systemd-coredumpd	help wanted
retrace/abrt/faf	in progress amerey@redhat.com
vtune	help wanted
pixie	help wanted
sentry symbolicator	partly released partly help wanted
VS Code	partly automatic via gdb
WinDbg	released in version 1.2104.13002.0, 2021-04
UDB	released in version 6.5
parca	released in version 0.8.0, 2022-02

# Example program: Fibonacci

```
● ● ●

1 func main() {
2     n := fibonacci(os.Args[1])
3     fmt.Println(n)
4 }
5
6 func fibonacci(n int) int {
7     if n <= 1 {
8         return n
9     }
10    return FibonacciRecursion(n-1) + FibonacciRecursion(n-2)
11 }
```

# Captured stack from eBPF

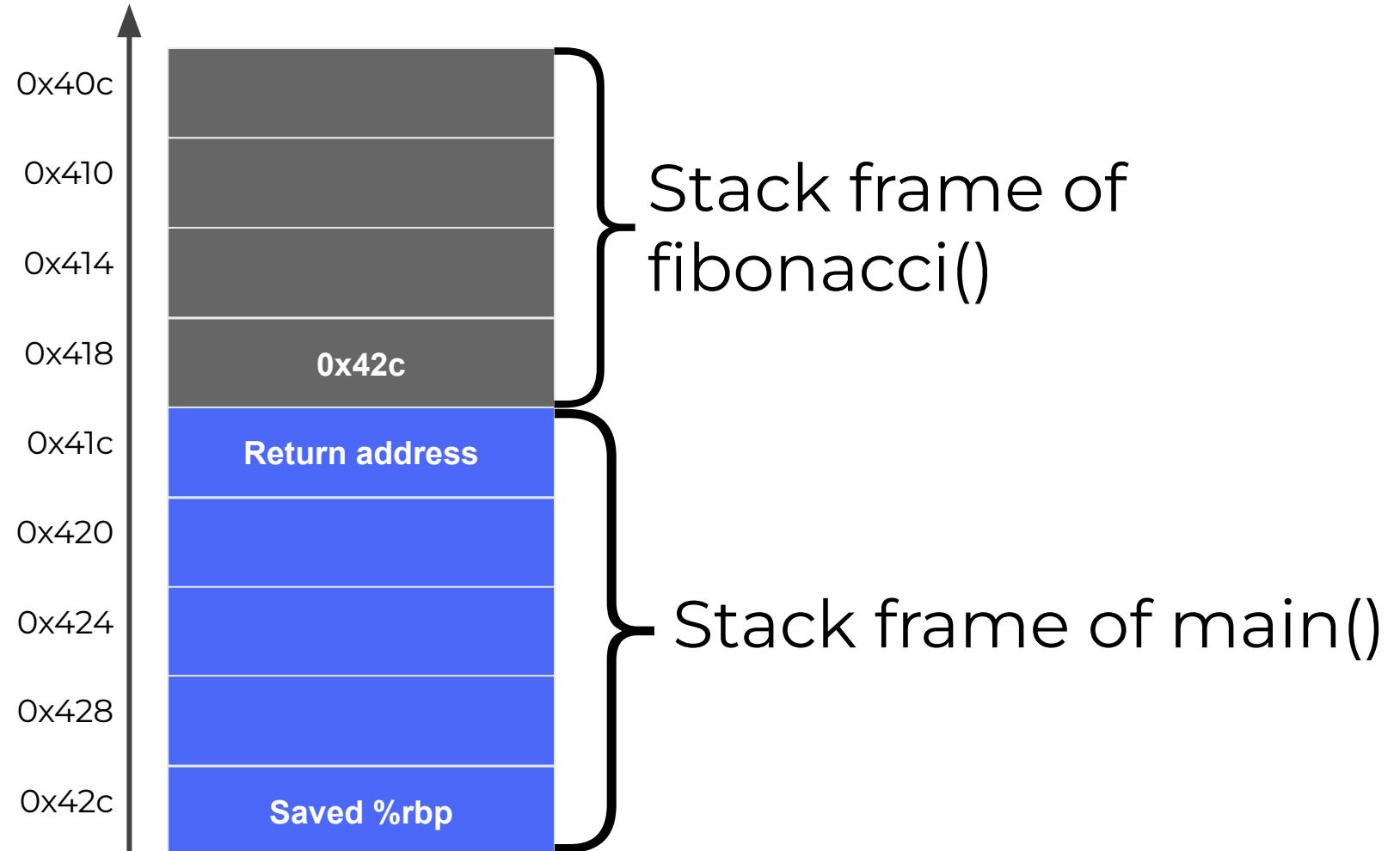


```
1 typedef struct stack_count_key {
2     u32 pid;
3     int user_stack_id;
4     int kernel_stack_id;
5 } stack_count_key_t;
6
7 #define MAX_STACK_ADDRESSES 1024
8 #define MAX_STACK_DEPTH 127
9
10 BPF_HASH(counts, stack_count_key_t, u64);
11 BPF_STACK_TRACE(stack_traces, MAX_STACK_ADDRESSES);
12
13 int profile_cpu(struct bpf_perf_event_data *ctx) {
14     // ... extract PID ...
15
16     stack_count_key_t key = {
17         .pid = tgid,
18     };
19
20     int stack_id = bpf_get_stackid(ctx, &stack_traces, BPF_F_USER_STACK);
21     key.user_stack_id = stack_id;
22
23     int kernel_stack_id = bpf_get_stackid(ctx, &stack_traces, 0);
24     key.kernel_stack_id = kernel_stack_id;
25
26     u64 zero = 0; u64 *count;
27     count = bpf_map_lookup_or_try_init(&counts, &key, &zero);
28     if (!count)
29         return 0;
30
31     __sync_fetch_and_add(count, 1);
32     return 0;
33 }
```

Full code at:  
<https://github.com/parca-dev/parca-agent/blob/main/bpf/cpu/cpu.bpf.c>

# How do we get a stack? Frame pointers!

Register	Value
%rsp	0x40c
%rbp	0x418



# **Frame Pointers**

**==**

# **Linked List**

# Evil compiler optimizations

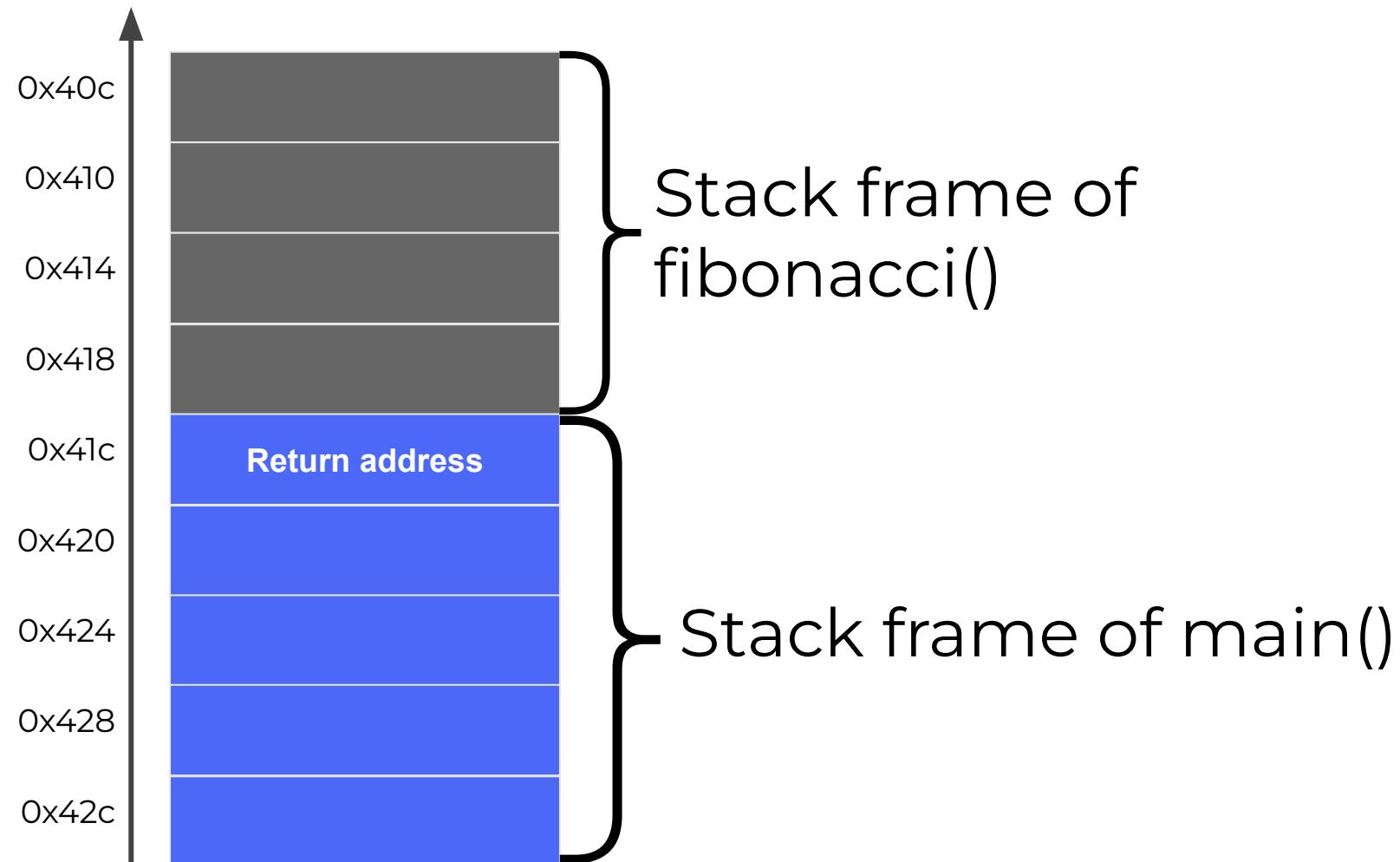
---

```
# Don't do this  
clang -fomit-frame-pointer main.c  
  
# Do this  
clang -fno-omit-frame-pointer main.c
```

# How do we get a stack? Frame pointers!

Register	Value
%rsp	0x40c
%rbp	0x418

Location	CFA	rbp	ra
0x40c	rsp+8	u	c-8
0x410	rsp+16	c-16	c-8
0x414	rbp+16	c-16	c-8
0x418	rsp+8	c-16	c-8



**Linux Plumbers Conference 2022:**

**Developing eBPF profilers for polyglot  
cloud-native applications**

**My point:**

**Stack Walking for any native binaries works!**

**What about interpreters?**

# Interpreters

[libpython3.11.so.1.0] Py\_BytesMain

[libpython3.11.so.1.0] ?

[libpython3.11.so.1.0] ?

[libpython3.11.so.1.0] Py\_RunMain

[libpython3.11.so.1.0] PyRun\_AnyFileExFlags

[libpython3.11.so.1.0] \_PyRun\_AnyFileObject

[libpython3.11.so.1.0] \_PyRun\_InteractiveLoopObject

[libpython3.11.so.1.0] PyRun\_InteractiveOneObjectEx

[libpython3.11.so.1.0] ?

[libpython3.11.so.1.0] run\_mod

[libpython3.11.so.1.0] run\_eval\_code\_obj

[libpython3.11.so.1.0] PyEval\_EvalCode

[libpython3.11.so.1.0] ?

[libpython3.11.so.1.0] \_PyEval\_Vector

[libpython3.11.so.1.0] \_PyEval\_EvalFrameDefault

[libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpython3.11.so.1.0] PyObject\_CallOneArg

[libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpython3.11.so.1.0] cfunction\_vectorcall\_O

[libpython3.11.so.1.0] sys\_displayhook

[libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpython3.11.so.1.0] PyImport\_GetModule

[libpyth] [libp] [libpython3.11.so.1.0] \_PyModuleSpec\_IsInitializing

[lib] [lib] [libpython3.11.so.1.0] PyObject\_GetAttr

[lib] [lib] [libpython3.11.so.1.0] ?

[lib] [lib] [libpython3.11.so.1.0] \_PyObject\_GenericGetA

[lib] [lib] [libpython3.11.so.1.0] ?

[libpython3.11.so.1.0] builtin\_print

[libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpyth] [libpython3.11.so.1.0] PyFile\_WriteString

[libpyth] [libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpyth] [libpython3.11.so.1.0] PyFile\_WriteObject

[libpyth] [libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpyth] [libpython3.11.so.1.0] PyObject\_CallOneArg

[libpyth] [libpython3.11.so.1.0] ?

[libpyth] [libp]

[libpyth] [libpython3.11.so.1.0] cfunction\_vectorcall\_O

[libpyth] [libpython3.11.so.1.0] ?

[libpyth] [libpython3.11.so.1.0] \_io\_Text

[libpyth] [libpython3.11.so.1.0] \_io\_TextIOWrapper\_write\_impl

[libpyth] [libpython3.11.so.1.0] \_io\_Text

Some prior art:

[github.com/javierhonduco/rbperf](https://github.com/javierhonduco/rbperf)

Parca Demo: <https://parca.dev/>



# Future

The future of eBPF

# Where are we today with eBPF?

Demonstrated applicability to a wide set of use cases

- Networking (Load balancing, routing, firewalling...)
- Security (System call filtering, access control...)
- Observability (Application profiling, performance monitoring...)

Breadth and depth of coverage is hardening

- Improving support in programming languages
- Adoption of USDTs for common libraries
- Cloud providers are looking to add eBPF to serverless platforms

Higher-level APIs have been established

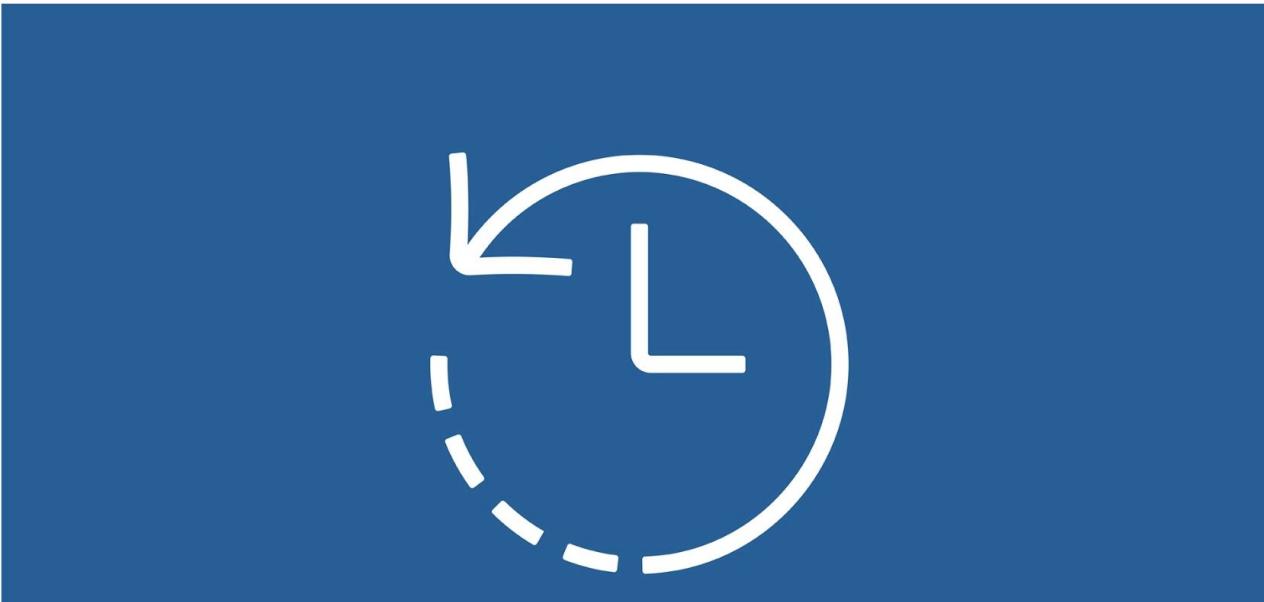
# Time-travel debuggers?

Engineering at Meta

Open Source ▾ Platforms ▾ Infrastructure Systems ▾ Physical Infrastructure ▾ Video Engineering & AR/VR ▾

POSTED ON APRIL 27, 2021 TO DEVINFRA

## Reverse debugging at scale



<https://engineering.fb.com/2021/04/27/developer-tools/reverse-debugging/>

# What frontiers might be addressed next?

- ⌚ *Performance*: Efficiently supporting large numbers of probes reading large amounts of data
- 🎯 *Accessibility*: Identifying what, where, and how to instrument
- 📊 *Analysis*: Interpreting the low-level “firehose” of raw data

# Mini Demo #3

Generating eBPF probes with LLMs



KubeCon



CloudNativeCon

Europe 2023

Session QR Codes will be  
sent via email before the event

Please scan the QR Code above  
to leave feedback on this session