- Container
- Docker
- eBPF
- Kubernetes
- Linux
- Python

# Agenda

- How to debug Python code
- How to check system performance
- What is eBPF
- Debug Python container with eBPF

# How to debug Python code

- print
- logging

- profile & timeit
- bdb & pdb
- trace
- ...

# How to debug Python code

- tracemalloc
  - Since Python 3.4
  - PEP 454 – trace Python memory allocations

```
→  ~ cat test.py
#!/usr/bin/env python
# coding=utf-8
import tracemalloc

tracemalloc.start()

a = range(1000)
b = range(1)

snapshot = tracemalloc.take_snapshot()
top_stats = snapshot.statistics('lineno')

print("[ Top 10 ]")
for stat in top_stats[:10]:
    print(stat)
→  ~ python test.py
[ Top 10 ]
test.py:7: size=652 B, count=3, average=217 B
test.py:8: size=48 B, count=1, average=48 B
→  ~
```

# How to debug Python code

- sys.audit
  - Since Python 3.8
  - PEP 578 – Python Runtime Audit Hooks

```
→  ~ cat test.py
#!/usr/bin/env python
# coding=utf-8
import sys


def audit(event, *args):
    print(f'event: {event}; args: {args}')


sys.addaudithook(audit)


with open('test.txt', 'w+') as f:
    f.write('Hello PyCon China 2020')
→  ~ python test.py |grep test.txt
event: open; args: (('test.txt', 'w+', 524866),)
→  ~
```

# How to check system performance

- Physical
  - CPU
  - Memory
  - Network
  - Storage I/O
- Software
  - File descriptors
  - mutex

# How to check system performance

- CPU
  - sar
  - mpstat
  - top
  - ps

# How to check system performance

- Memory
  - sar
  - free
  - vmstat
  - top/htop

# How to check system performance

- Network
  - top/iftop
  - ifconfig
  - iproute2

```
→  ~ sar -n DEV 1 1
Linux 5.9.8-100.fc32.x86_64 (moelove.info)      11/26/20      _x86_64_      (4 CPU)

00:50:18       IFACE   rxpck/s   txpck/s   rxkB/s   txkB/s   rxcmp/s   txcmp/s   rxmcst/s   %ifutil
00:50:19          lo      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19      enp1s0      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19      wlp2s0      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19      virbr0      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19   virbr0-nic     0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19  br-28c4330f1309          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19  br-398149a175d9          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19  br-4539c9c0f42d          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19  br-4814a2909acc          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19  br-a703cf96c002          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19     docker0      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19  br-feb5d9c4f03a          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19  br-19a39f873a23          0.00     0.00     0.00      0.00      0.00      0.00       0.00      0.00
00:50:19    vboxnet0      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19   vethf9d0c9f      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
00:50:19   veth09df558      0.00      0.00     0.00     0.00      0.00      0.00       0.00      0.00
```

# How to check system performance

- Storage
  - sar
  - iostat
  - iotop

http://www.brendangregg.com/linuxperf.html 2017

# What is eBPF

# What is eBPF

- eBPF is a revolutionary technology that can run sandboxed programs in the Linux kernel without changing kernel source code or loading kernel modules.

  - via:https://ebpf.io/what-is-ebpf/

# What is eBPF

- enhanced Berkeley Packet Filter
- Use cases:
  - Networking
  - Tracing
  - Security
  - ...

# How does eBPF work

# Raw BPF

```c
47    struct bpf_insn prog[] = {
48        BPF_MOV64_REG(BPF_REG_6, BPF_REG_1),
49        BPF_LD_ABS(BPF_B, ETH_HLEN + offsetof(struct iphdr, protocol) /* R0 = ip->proto */),
50        BPF_STX_MEM(BPF_W, BPF_REG_10, BPF_REG_0, -4), /* *(u32 *)(fp - 4) = r0 */
51        BPF_MOV64_REG(BPF_REG_2, BPF_REG_10),
52        BPF_ALU64_IMM(BPF_ADD, BPF_REG_2, -4), /* r2 = fp - 4 */
53        BPF_LD_MAP_FD(BPF_REG_1, map_fd),
54        BPF_RAW_INSN(BPF_JMP | BPF_CALL, 0, 0, 0, BPF_FUNC_map_lookup_elem),
55        BPF_JMP_IMM(BPF_JEQ, BPF_REG_0, 0, 2),
56        BPF_MOV64_IMM(BPF_REG_1, 1), /* r1 = 1 */
57        BPF_RAW_INSN(BPF_STX | BPF_XADD | BPF_DW, BPF_REG_0, BPF_REG_1, 0, 0), /* xadd r0 += r1 */
58        BPF_MOV64_IMM(BPF_REG_0, 0), /* r0 = 0 */
59        BPF_EXIT_INSN(),
60    };
61    size_t insns_cnt = sizeof(prog) / sizeof(struct bpf_insn);
```

Linux samples/bpf/sock_example.c

# How to write eBPF programs

## LLVM eBPF compiler

```c
1  SEC("socket1")
2  int bpf_prog1(struct __sk_buff *skb)
3  {
4      int index = load_byte(skb, ETH_HLEN + offsetof(struct iphdr, protocol));
5      long *value;
6
7      if (skb->pkt_type != PACKET_OUTGOING)
8          return 0;
9
10     value = bpf_map_lookup_elem(&my_map, &index);
11     if (value)
12         __sync_fetch_and_add(value, skb->len);
13
14     return 0;
15 }
16 char _license[] SEC("license") = "GPL";
```

Linux samples/bpf/sockex1_kern.c

- BCC proper
- BCC-tools
- Front-ends
  - Python
  - Lua
  - C helper functions

https://github.com/iovisor/bcc

# BPF Compiler Collection

https://ebpf.io/what-is-ebpf/#development-toolchains

# Linux bcc/BPF Tracing Tools

opensnoop statsnoop
syncsnoop

ucalls uflow
uobjnew ustat
uthreads ugc

c* java* node* php*
python* ruby*

mysqld_qslower
dbstat dbslower
bashreadline

gethostlatency
memleak
sslsniff

filetop
filelife fileslower
vfscount vfsstat

syscount
killsnoop

cachestat cachetop
dcstat dcsnoop
mountsnoop

execsnoop
exitsnoop
pidpersec

trace
argdist
funccount
funcslower
funclatency
stackcount
profile

cpudist cpuwalk
runqlat runqlen
runqslower
cpuunclaimed
deadlock

offcputime wakeuptime
offwaketime softirqs

btrfsdist
btrfsslower
ext4dist ext4slower
nfsslower nfsdist
xfsslower xfsdist
zfsslower
zfsdist

slabratetop
oomkill memleak
shmsnoop drsnoop

hardirqs
criticalstat
ttysnoop

mdflush

biotop biosnoop
biolatency bitesize

tcptop tcplife tcptracer
tcpconnect tcpaccept tcpconnlat
tcpretrans tcpsubnet tcpdrop
tcpstates

llcstat

Other:
capable

sofdsnoop

## Applications
### Runtimes
### System Libraries
### System Call Interface
- VFS
- Sockets
- File Systems
- TCP/UDP
- Volume Manager
- IP
- Block Device
- Net Device
- Scheduler
- Virtual Memory
- Device Drivers

CPUs

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

# How to write eBPF programs

## bcc eBPF compiler

- Backends & data structures
  - "restricted C"
- Frontends & loaders
  - Python/Lua

```python
#!/usr/bin/env python
# coding=utf-8
import time
import socket

from bcc import BPF

ebpf_str = """
#include <uapi/linux/if_ether.h>
#include <uapi/linux/ip.h>

BPF_ARRAY(count_map, u64, 256);

int count_packets(struct __sk_buff *skb) {
  int index = load_byte(skb, ETH_HLEN + offsetof(struct iphdr,
  u64 *value = count_map.lookup(&index);
  if (value)
    count_map.increment(index);
  return 0;
}
"""

bpf = BPF(text=ebpf_str)
pfilter = bpf.load_func("count_packets", BPF.SOCKET_FILTER)
BPF.attach_raw_socket(pfilter, "lo")
for i in range(10):
    print("TCP: {0}, UDP: {1}, ICMP: {2}".format(
        bpf["count_map"][socket.IPPROTO_TCP].value,
        bpf["count_map"][socket.IPPROTO_UDP].value,
        bpf["count_map"][socket.IPPROTO_ICMP].value,
    ))
    time.sleep(1)
```
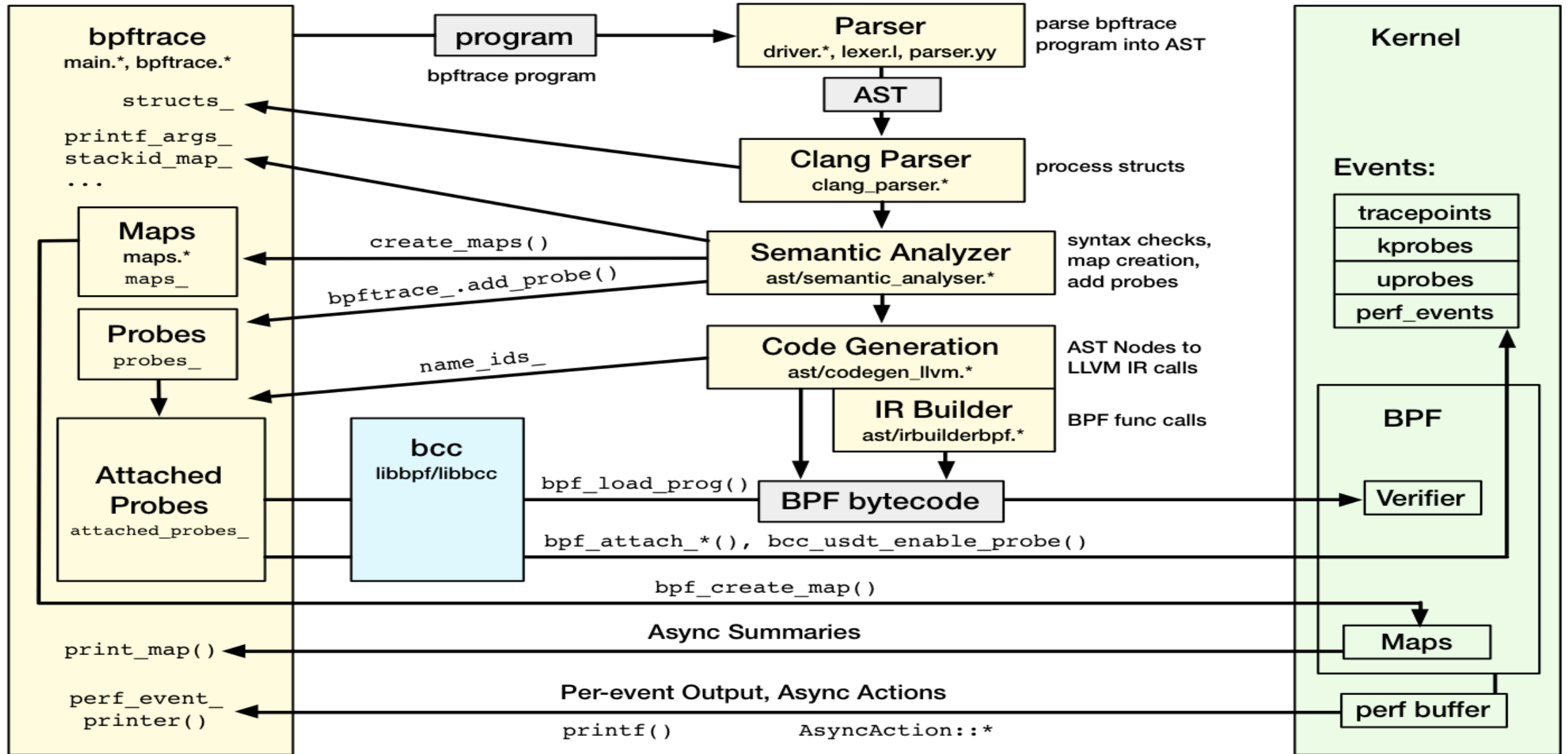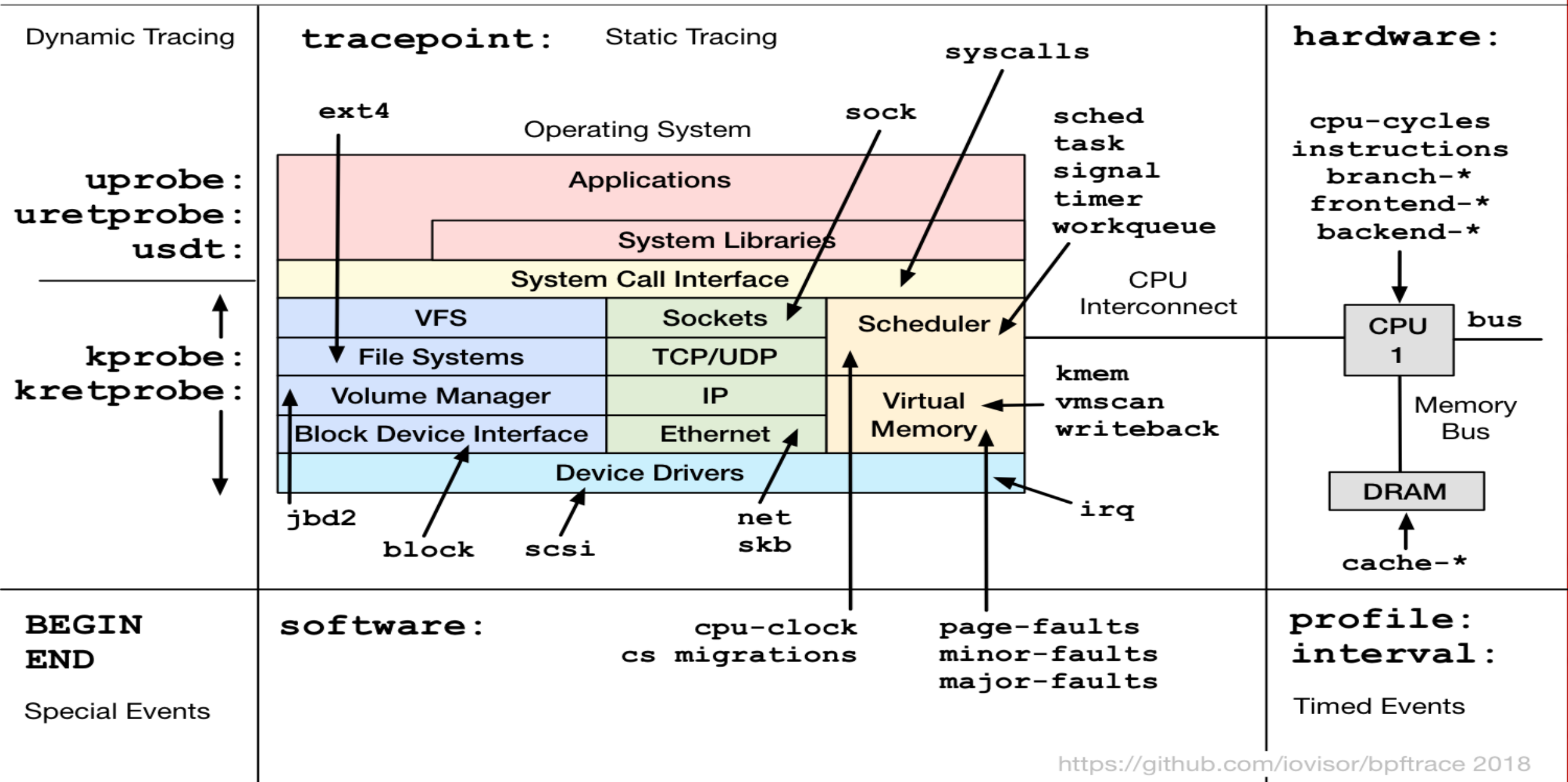
# bpftrace (DTrace 2.0)

- High-level tracing language
- Using LLVM as a backend to compile scripts to BPF-bytecode
- The bpftrace language is inspired by awk and C

# bpftrace Internals



https://github.com/iovisor/bpftrace 2018

# bpftrace Probe Types



Dynamic Tracing

**tracepoint:** Static Tracing

syscalls

**hardware:**

ext4

**uprobe:**
**uretprobe:**
**usdt:**

Operating System

sock

sched
task
signal
timer
workqueue

cpu-cycles
instructions
branch-*
frontend-*
backend-*

Applications

System Libraries

System Call Interface

CPU
Interconnect

**kprobe:**
**kretprobe:**

| VFS | Sockets | Scheduler |
| File Systems | TCP/UDP | |
| Volume Manager | IP | Virtual Memory |
| Block Device Interface | Ethernet | |

Device Drivers

CPU
1

bus

kmem
vmscan
writeback

Memory
Bus

DRAM

jbd2

block    scsi

net
skb

irq

cache-*

**BEGIN**
**END**

Special Events

**software:**

cpu-clock
cs migrations

page-faults
minor-faults
major-faults

**profile:**
**interval:**

Timed Events

https://github.com/iovisor/bpftrace 2018

- Container's Namespace
  - docker run --pid container:xxxx
  - nsenter --target $PID --pid

- USDT (--with-pydebug --with-dtrace)
  - pythoncalls.sh
  - pythonflow.sh
  - pythongc.sh
  - pythonstat.sh