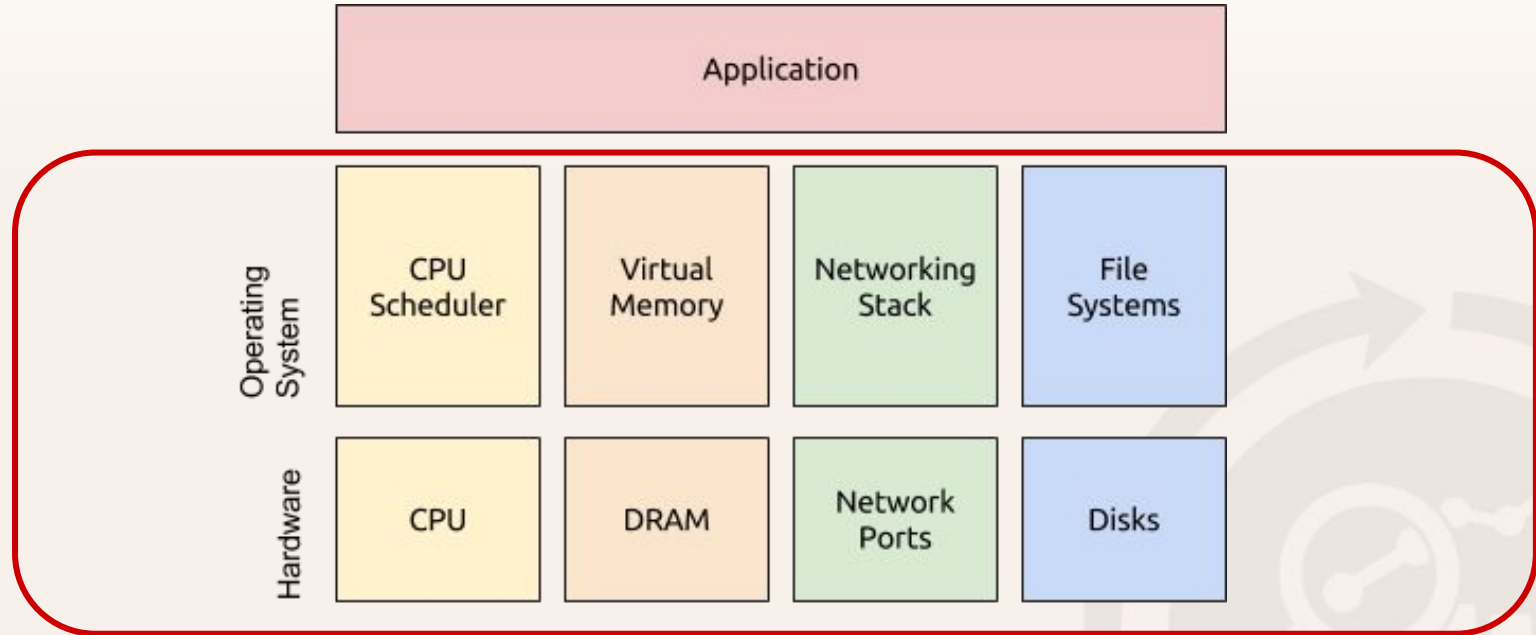# Linux System Monitoring with eBPF

*DevOpsDays Zurich, 2018-05-03*

*Heinrich Hartmann*

# System Monitoring is about Kernel & Hardware

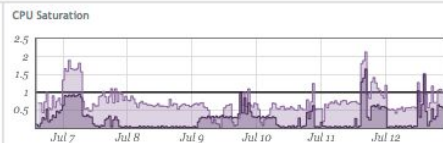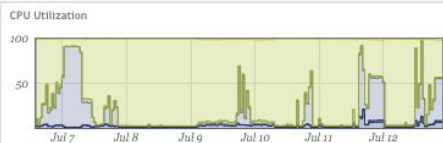# Best Practice: The USE Method

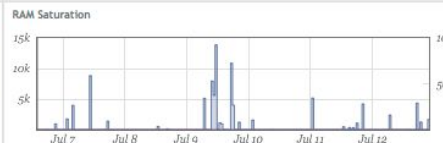https://www.circonus.com/2017/08/system-monitoring-with-the-use-dashboard

# Best Practice: The USE Method

https://www.circonus.com/2017/08/system-monitoring-with-the-use-dashboard

# Lot's of Unknowns remaining

https://www.circonus.com/2017/08/system-monitoring-with-the-use-dashboard



Heinrich.Hartmann@Circonus.com

# eBPF allows unparalleled insights

https://github.com/iovisor/bcc



Credits:
- Brendan Gregg @ Netflix (Sun)
- Sasha Goldshtein @ Sela, Microsoft
- Brenden Blanco @ VMWare
- Linus Torvalds, et. al.

# eBPF allows unparalleled insights

https://github.com/iovisor/bcc



Linux bcc/BPF Tracing Tools

Credits:
- Brendan Gregg @ Netflix (Sun)
- Sasha Goldshtein @ Sela, Microsoft
- Brenden Blanco @ VMWare
- Linus Torvalds, et. al.

# CPU: Scheduling Latency



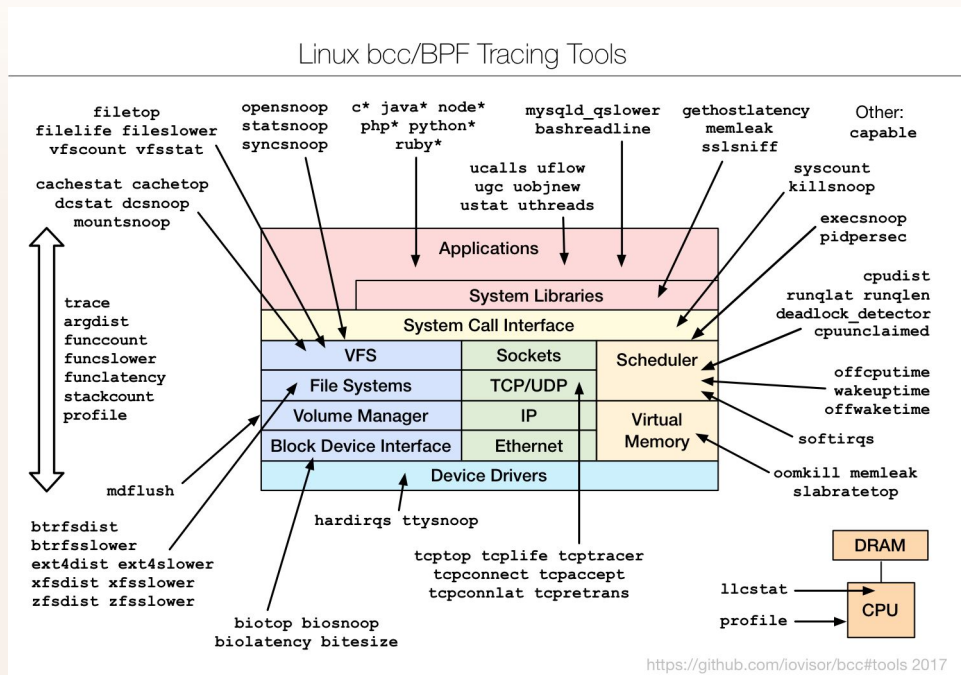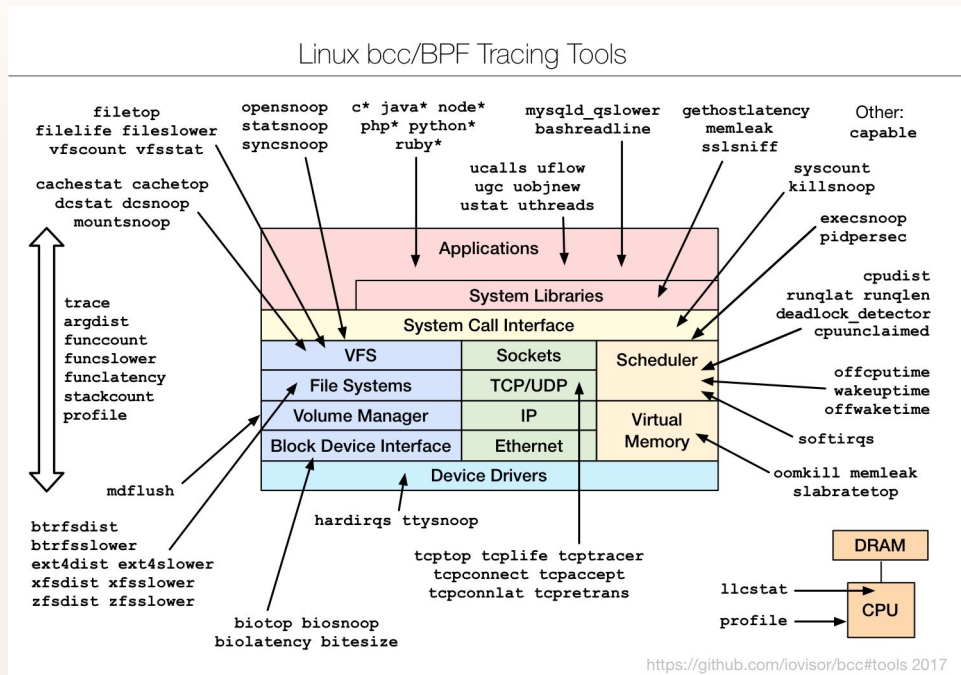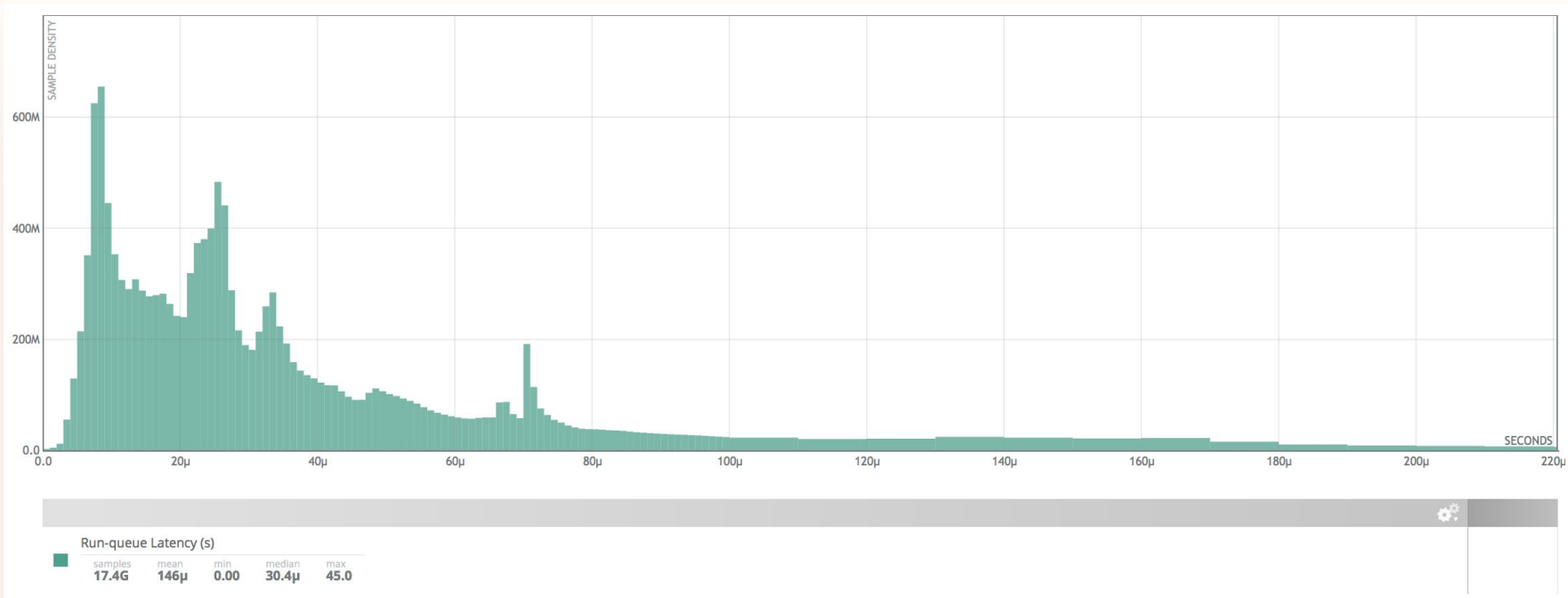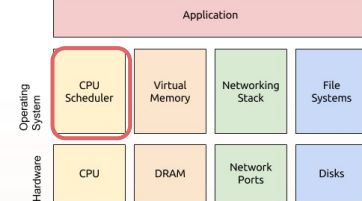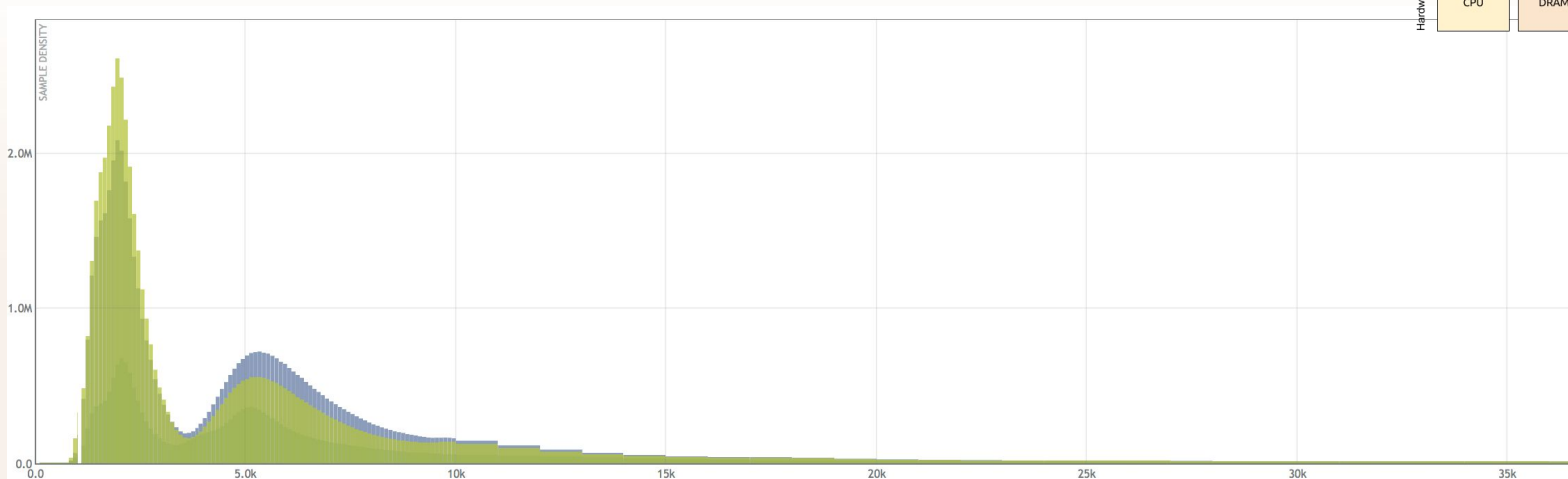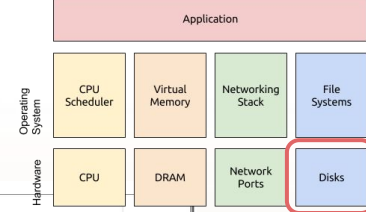Heinrich.Hartmann@Circonus.com

# Disk: Block-I/O Latency



excite-compute cosi/system: iolatency `vda (on excite-compute, from Chicago, IL, US)

| | samples | mean | min | median | max |
|---|---|---|---|---|---|
| | 24.3M | 32.9k | 120 | 5.23k | 55.0M |

9.83M **40.4%** | 256k **1.1%** | 14.2M **58.5%**

excite-compute cosi/system: iolatency `vdc (on excite-compute, from Chicago, IL, US)

| | samples | mean | min | median | max |
|---|---|---|---|---|---|
| | 61.4M | 13.0k | 450 | 4.87k | 39.0M |

28.6M **46.5%** | 517k **0.8%** | 32.3M **52.6%**

excite-compute cosi/system: iolatency `vdb (on excite-compute, from Chicago, IL, US)

| | samples | mean | min | median | max |
|---|---|---|---|---|---|
| | 59.7M | 11.9k | 400 | 3.14k | 44.0M |

33.0M **55.4%** | 416k **0.7%** | 26.2M **43.9%**

# Disk: Block-I/O Latency

# Disk: Block-I/O Latency over time

# Disk: Block-I/O Latency over time

# Don't shout in the Datacenter

Brendan Gregg (2008) https://www.youtube.com/watch?v=tDacjrSCeq4



Heinrich.Hartmann@Circonus.com

# System Calls: The Kernel API



Application

System Call API

Monitor

**R**ate
**E**rrors
**D**uration

Operating System
- CPU Scheduler
- Virtual Memory
- Networking Stack
- File Systems

Hardware
- CPU
- DRAM
- Network Ports
- Disks

IRONdb
POWERED BY CIRCONUS

Heinrich.Hartmann@Circonus.com

CIRCONUS

# Syscalls: Rate / Count



Heinrich.Hartmann@Circonus.com

# Syscalls: Duration



Heinrich.Hartmann@Circonus.com

# Syscall durations span >8 orders of magnitude



10 us

100 ms

1.5 tn events total

1s

IRONdb
POWERED BY CIRCONUS

CIRCONUS

# File System: Latency

# Memory: Allocation Latency



excite-compute cosi/system: bpf `syscall `latency `sys_mremap (on excite-compute, from Chicago, IL, US)
samples 402k | mean 29.3µ | min 0.00 | median 7.28µ | max 38.0m

excite-compute cosi/system: bpf `syscall `latency `sys_brk (on excite-compute, from Chicago, IL, US)
samples 12.5M | mean 12.2µ | min 0.00 | median 6.58µ | max 86.0m

excite-compute cosi/system: bpf `syscall `latency `sys_mmap (on excite-compute, from Chicago, IL, US)
samples 65.3M | mean 19.5µ | min 0.00 | median 8.64µ | max 7.40

excite-compute cosi/system: bpf `syscall `latency `sys_munmap (on excite-compute, from Chicago, IL, US)
samples 11.9M | mean 121µ | min 0.00 | median 26.5µ | max 7.80

Heinrich.Hartmann@Circonus.com

# Further Reading

Slides: [@HeinrichHartman](#) / [#DevOpsDaysZH](#)

Code: [https://github.com/circonus-labs/nad/.../bccbpf](#)

Blog: [http://www.circonus.com/2018/05/linux-system-monitoring-with-ebpf/](#)

IRONdb
POWERED BY CIRCONUS

CIRCONUS