# The Performance Debugging Toolkit

Under the hood of Linux systems

# So you have your project
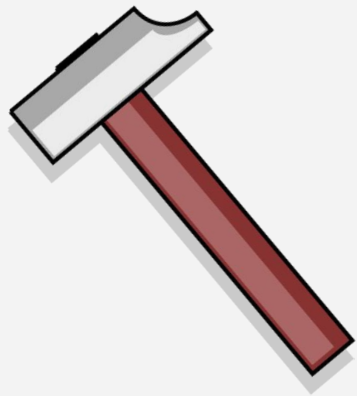
# It should work

# But...

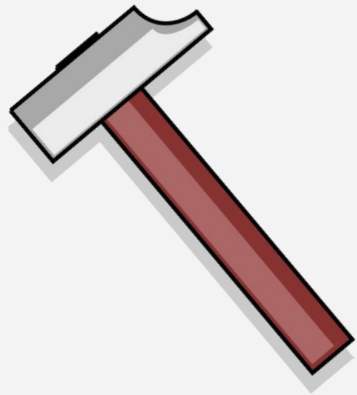# It should work, but…

- **Performance is not as expected**
- **Unpredictable results**
- **Unexpected resource utilization**
- **Weird system behavior**

# What to do???

- Add prints
- Time execution
- Random parameter tuning
- Improvised optimizations
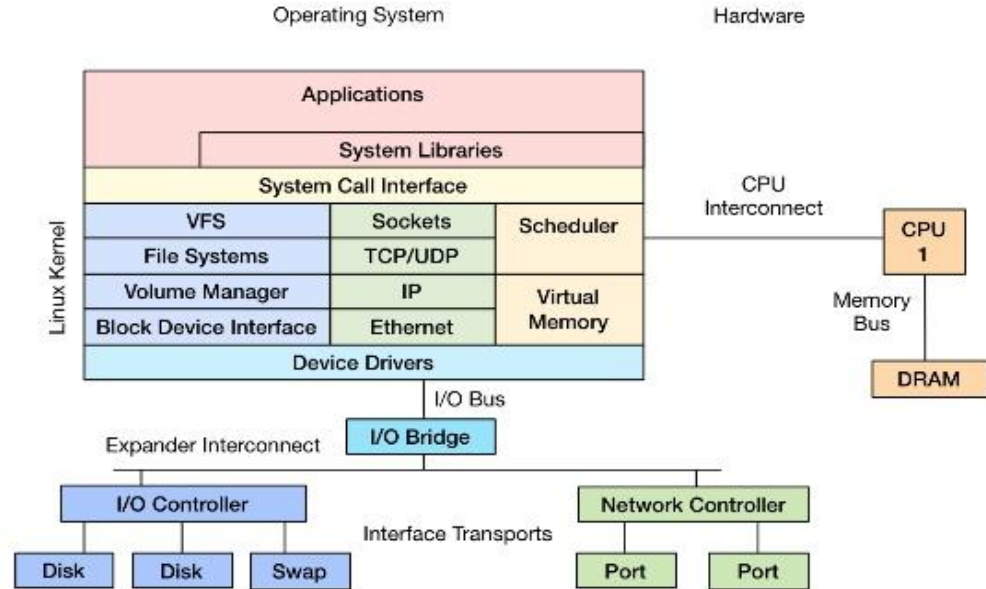- Execute until a good measurement happens?

# What to do???

- Add prints
- Time execution
- Random parameter tuning
- Improvised optimizations
- Execute until a good measurement happens?
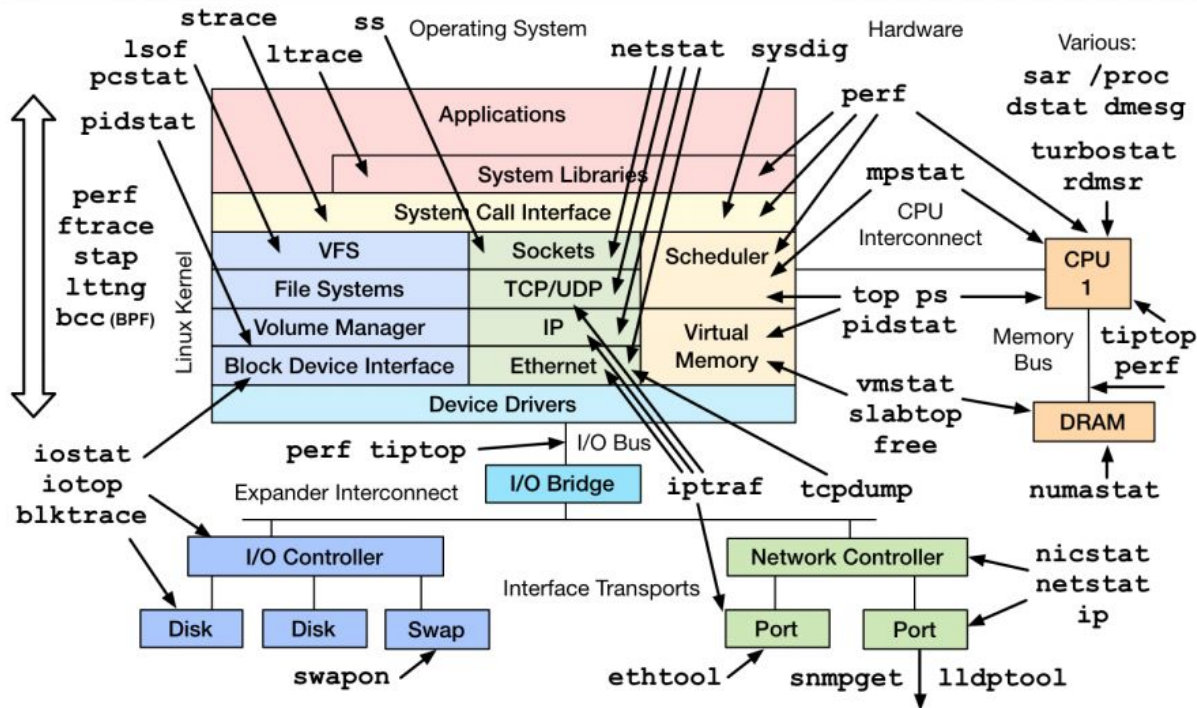
# Where to start?

## Linux Observability Tools

| Operating System | | Hardware |

Brendan Gregg 2014

# Well... get the toolbox!



Linux Performance Observability Tools

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

# CPU - top, htop, ps

# Multicore - mpstat

# Memory - vmstat

# I/O - iostat

# Networking - iftop, netstat

# Execution - ltrace, strace

# Kernel calls - ftrace



```
            |   mutex_unlock() {
            |   smp_irq_work_interrupt() {
            |     irq_enter() {
  0.132 us  |       rcu_irq_enter();
  0.557 us  |     }
            |     __wake_up() {
            |       __wake_up_common_lock() {
  0.046 us  |         _raw_spin_lock_irqsave();
  0.055 us  |         __wake_up_common();
  0.053 us  |         _raw_spin_unlock_irqrestore();
  1.031 us  |       }
  1.349 us  |     }
            |     __wake_up() {
            |       __wake_up_common_lock() {
  0.051 us  |         _raw_spin_lock_irqsave();
            |         __wake_up_common() {
            |           autoremove_wake_function() {
            |             default_wake_function() {
            |               try_to_wake_up() {
  0.164 us  |                 _raw_spin_lock_irqsave();
```

# /proc - all things counters

- A filesystem for system information
- The tools probably read this
- You could write your own!

- /proc/cpuinfo
- /proc/devices
- /proc/ksyms

# The (un)biased favourite

# eBPF - the swiss army knife

- Basically all we just said
- … on steroids!
- Easy to use: BCC
- Not for the faint of heart: native

# One tool to rule them all



Linux bcc/BPF Tracing Tools

https://github.com/iovisor/bcc#tools 2018

# Distributed systems

- **Look into debug tools in frameworks**
- **Look at your systems capability (e.g. hadoop)**
- **Distributed monitoring (e.g. ganglia)**

# Now we have the tools, but we still need the manual

# Remember!

- **Run long enough**
- **Measure one level deeper**
- **Crosscheck with more subsystems**
- **Don't trust your numbers**
- **Back-of-the-envelope estimates**

# Minimize your entropy

- **Execute in controlled environment**
- **Reduce randomness**
- **Cgroups**
- **Cpulimit**
- **Numactl**

# Minimize your entropy

- **Execute in controlled environment**
- **Reduce randomness**
- **Cgroups**
- **Cpulimit**
- **Numactl**

*DIY CONTAINERS!!!*

# You don't have a hammer, don't treat everything as a nail!

# References

- http://www.brendangregg.com/Slides/Velocity2015_LinuxPerfTools.pdf

- http://www.brendangregg.com/blog/2019-01-01/learn-ebpf-tracing.html

- https://cacm.acm.org/magazines/2018/7/229031-always-measure-one-level-deeper/fulltext

- https://github.com/iovisor/bcc

- https://lwn.net/Articles/365835/