# HIGH PERFORMANCE $Co_{ding}^{mputing}$ (HPC)

ALBERTO DASSATTI - 2023

# You need to measure!

TO MEASURE IS TO KNOW

"In God we trust. All others must bring data".

W. Edwards Deming

Every measure has its own tool. CS is not different.

And measuring correctly is even harder.

# Measuring time!

# THE BASICS: TIME

```bash
[al@lap ~]# time your_command
[al@lap ~]# <your_command output omitted>
[al@lap ~]# real    0m3.168s
[al@lap ~]# user    0m2.952s
[al@lap ~]# sys     0m0.180s
```

Let's test hyperfine.

# THE BASICS: TIME

Can we have a little more visibility?

```bash
[al@lap ~]# strace -tt
```

```bash
[al@lap ~]# ltrace -tt
```

# THE BASICS: TIME

If you have access at the source code and the toolchain, some more options are available.

```bash
[al@lap ~]# gcc -pg -O3 -Wall source.c -o source
[al@lap ~]# ./source
[al@lap ~]# gprof source gmon.out > analysis.txt
```

With the help of -pg there is also the option of uftrace

## MANUAL INSTRUMENTATION

Gcc can help.

### rdtsc

```
static __inline__ unsigned long long rdtsc(void)
{
  unsigned long long int x;
  __asm__ volatile (".byte 0x0f, 0x31" : "=A" (x));
  return x;
}
```
source

cpucycles is a library that improve a lot on this raw technique.

# MANUAL INSTRUMENTATION ON STEROIDS

You are not forced to trace all the code, but if you know your critical path you can use an instrumentation library (eg. easy_profiler or Likwid marker API).

a good tutorial

# WARNING AND LIES

Understand the numbers provided by the tools is not always easy. Please read this article and this discussion to be sure you get the most accurate info from your tool. I strongly reccomend you this fun and instructive video.

# What our system is doing?

For a first look we may use

```bash
[al@lap ~]# top
```

htop is a feature packed alternative

# Linux Performance Tools

**strace**
**perf trace**
**sysdig**

**opensnoop statsnoop**
**syncsnoop**

**tcptop tcplife**
**tcpconnect tcpaccept**
**tcpconnlat tcpretrans**

**mysqld_qslower, ...**

Various, observability:
**sar /proc**
**dstat dmesg**

Various, static:
**sysctl /sys**
**dmesg lshw**
**journalctl**
**lsmod**

Various, tracing:
**capable**

**cachestat dcsnoop**
**filetop fileslower**
**mountsnoop**

App Config

ugc ucalls

Operating System

ldd

gethostlatency

ltrace

**cpudist execsnoop**
**runlat cpudist**
**offcputime**

Hardware

**intel_gpu_top**
**intel_gpu_time**

**intel_gpu_\**
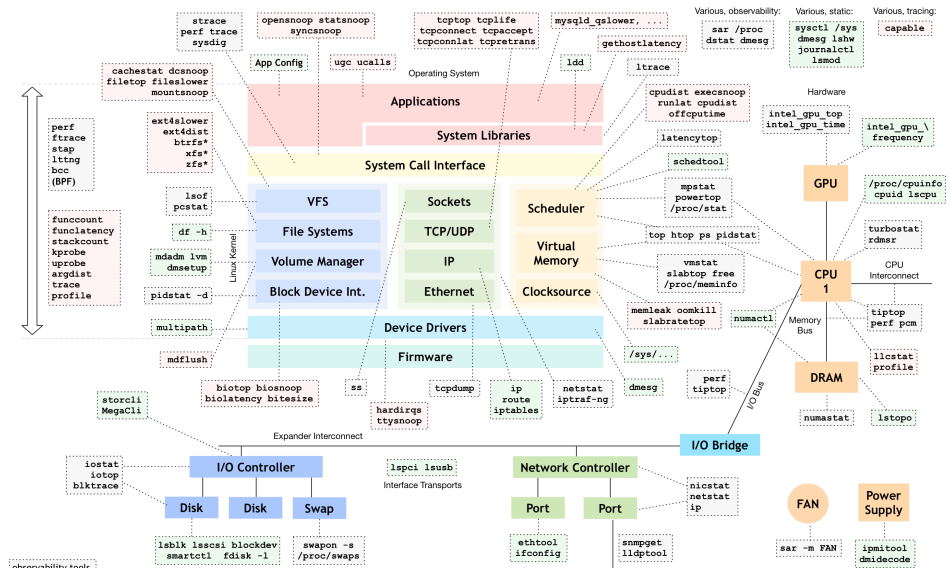**frequency**

**perf**
**ftrace**
**stap**
**lttng**
**bcc**
**(BPF)**

**ext4slower**
**ext4dist**
**btrfs***
**xfs***
**zfs***

## Applications

### System Libraries

### System Call Interface

latencytop

schedtool

GPU

**/proc/cpuinfo**
**cpuid lscpu**

**lsof**
**pcstat**

VFS

Sockets

Scheduler

**mpstat**
**powertop**
**/proc/stat**

**df -h**

File Systems

TCP/UDP

Virtual
Memory

**top htop ps pidstat**

**turbostat**
**rdmsr**

**funccount**
**funclatency**
**stackcount**
kprobe
**uprobe**
**argdist**
**trace**
**profile**

**mdadm lvm**
**dmesetup**

Volume Manager

IP

**vmstat**
**slabtop free**
**/proc/meminfo**

CPU
1

CPU
Interconnect

**pidstat -d**

Block Device Int.

Ethernet

Clocksource

**memleak oomkill**
**slabratetop**

numactl

Memory
Bus

**tiptop**
**perf pcm**

**multipath**

### Device Drivers

**perf**
**tiptop**

**mdflush**

### Firmware

**/sys/...**

dmesg

**llcstat**
**profile**

DRAM

**biotop biosnoop**
**biolatency bitesize**

ss

tcpdump

**ip**
**route**
**iptables**

**netstat**
**iptraf-ng**

**numastat**

**lstopo**

**storcli**
**MegaCli**

**hardirqs**
**ttysnoop**

Expander Interconnect

I/O Bridge

**iostat**
**iotop**
**blktrace**

I/O Controller

**lspci lsusb**

Network Controller

Interface Transports

**nicstat**
**netstat**
**ip**

Disk

Disk

Swap

Port

Port

FAN

Power
Supply

**lsblk lsscsi blockdev**
**smartctl fdisk -l**

**swapon -s**
**/proc/swaps**

**ethtool**
**ifconfig**

**snmpget**
**lldptool**

**sar -m FAN**

**ipmitool**
**dmidecode**

Linux Kernel

I/O Bus

observability tools
static performance tools          these can observe the state of the system at rest, without load
perf-tools/bcc tracing tools      https://github.com/brendangregg/perf-tools  https://github.com/iovisor/bcc

style inspired by reddit.com/u/redct
http://www.brendangregg.com/linuxperf.html 2017

Up to date tools catalog can be found here and here.

OK, BUT WHY IT IS TAKING SO LONG?

# Workload Characterization!

## micro Benchmarking

Isolate a small kernel or a function and analyze it alone

vs

## System Benchmarking

Looking at the performance of a system fully in service

### Statistical

Periodically stop the system and look at its state. If you collect enough data you can have a good picture of the situation.

- ▶ Advantage: do not need any code modification
- ▶ Warning: can be missleading, use long runs
- ▶ Example: `perf`

vs

### Instrumentation

Modify your code manually or automatically to emit `events` and then process them

- ▶ Advantage: tracing can be very precise
- ▶ Warning: you cannot see what's outside of your code or what your input is not stimulating
- ▶ Example: `gprof`

### Time-based sampling (TBS)

Program a timer. When the timer ends look at the `pc` and at the *stack* to get useful information.

- ▶ It's useful to have the debug symbols (-g)
- ▶ Trade sampling frequency for accuracy (limit the overhead)
- ▶ Some more ideas (`-fno-exceptions`, `-fno-rtti`, `-fno-omit-frame-pointer`)

vs

### Event-based sampling (EBS)

Count some events (HW or SW generated) and try to get what's going on.

- ▶ Low overhead (Intel $< 2\%$ using performance monitor counters)
- ▶ Use event ratios (cache miss/ cache access)

# PERFORMANCE COUNTERS

## oprofile, perf, VTune

- ► Total instruction count and mix
- ► Branch events
- ► Load/store events
- ► L1/L2 cache events
- ► Prefetching events
- ► TLB events
- ► Multicore events

Use event ratios

```
sudo perf list
```

# Perf Tutorial

blog, blog, pdf, Cern

# PERF GOODIES

- Pmc-cloud-tools
- HotSpot
- Flamegraph

but, what if what you are interested in is not already an event or you want more visibility?

Dynamic instrumentation is the answer!

Since long time (3.15 at least), Linux has the infrastructure to dynamically add trace point in both kernel space (kprobes) and user space (uprobe).

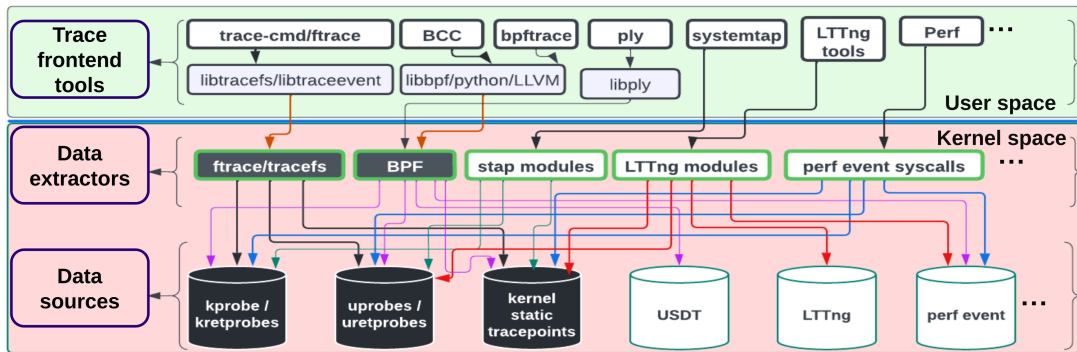They are not usually used directly. `perf probe` can help.

# Tracing

Bootlin Debugging-slides from pg 144

# Linux trace system



From Here.

# TRACING



ftrace    perf_events    eBPF    SystemTap

LTTng    ktap    dtrace4linux    OEL DTrace    sysdig

[ read this ]

To know more:

- ▶ have a look here
- ▶ read carefully and understand these Methods
- ▶ dig this mine of information

# QUESTIONS