# Diagnosing Performance Bottlenecks in Production Systems

by Julia Iacoviello

# Performance Issue? What to do?

- ✘ **Restart system/services**
- ✘ **Reapply thermal paste**
- ✘ **Reseat components**
- ✘ **Blow out the case with compressed air**
- ✘ **Kill random processes**
- ✘ **Percussive maintenance**
- ✘ **Upgrade hardware**

LIN·BIT

# Julia Iacoviello
# Systems Engineer
# LINBIT

- **Implements and supports Highly Available (HA) production environments**
- **High availability: Groups of servers (clusters) where one is configured to automatically switch (fail over) to a working node if there is an issue**
- **Many LINBIT clients have tight SLAs for service uptime**



LIN·BIT

# Are These *Good* Numbers?



???

???

# Talk Outline

- **Introduction** *(you are here)*

- **Basic Concepts & Methodologies**

- **Demos**

LIN·BIT

# Overhead

- **The impact that *gathering the data itself* has on the system**
- **Different methods have more or less**

```
440036 julia      20   0   13480    4560    3464 R   0.3   0.0   0:02.95 top
440057 julia      20   0 2527284  167616  104808 S   0.3   0.3   0:06.24 spectacle
440222 root       20   0       0       0       0 I   0.3   0.0   0:01.49 kworker/12:1-mm_percpu_wq
     1 root       20   0  167156   12776    8420 S   0.0   0.0   0:05.80 systemd
     2 root       20   0       0       0       0 S   0.0   0.0   0:00.24 kthreadd
     3 root        0 -20       0       0       0 I   0.0   0.0   0:00.00 rcu_gp
     4 root        0 -20       0       0       0 I   0.0   0.0   0:00.00 rcu_par_gp
     5 root        0 -20       0       0       0 I   0.0   0.0   0:00.00 slub_flushwq
```

LIN:BIT

# Observability

- **The ability to measure the system state based on what it is *already doing***
- **Can also refer to static configuration of the system**
- **Preferred (in most cases) for systems already deployed to production**

```
julia@julia-XPS-15-9510:/$ cat /proc/pressure/cpu
some avg10=0.00 avg60=0.00 avg300=0.00 total=390403417
full avg10=0.00 avg60=0.00 avg300=0.00 total=0
julia@julia-XPS-15-9510:/$ cat /proc/pressure/memory
some avg10=0.00 avg60=0.00 avg300=0.00 total=1248900
full avg10=0.00 avg60=0.00 avg300=0.00 total=1189250
julia@julia-XPS-15-9510:/$ ifconfig wlp0s20f3
wlp0s20f3: flags=4163<UP,BROADCAST,RUNNING,MULTICAST>  mtu 1500
        inet 192.168.1.15  netmask 255.255.255.0  broadcast 192.168.1.255
        inet6 fe80::4971:74b0:db73:8939  prefixlen 64  scopeid 0x20<link>
        ether 4c:79:6e:d3:93:a2  txqueuelen 1000  (Ethernet)
        RX packets 14660848  bytes 18705361239 (18.7 GB)
```
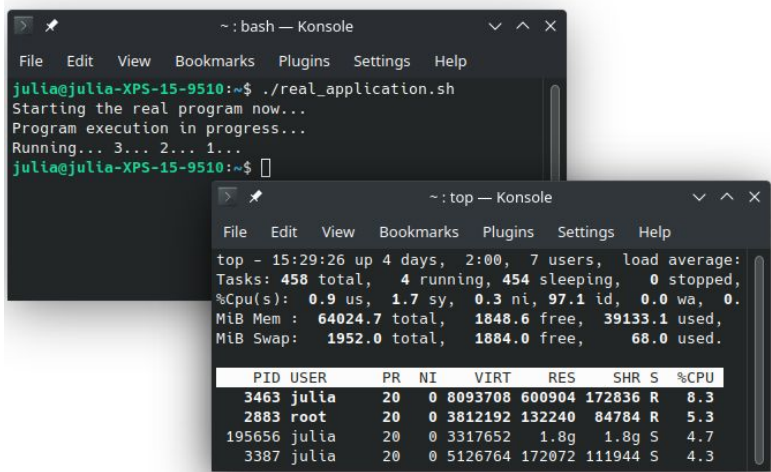
LIN·BIT

# Micro-Benchmarking

- **Metrics derived from a *simulated workload* applied to one component or one subset of system components**
- **Less *overhead* than Macro-Benchmarking**

```
julia@julia-XPS-15-9510:/$ ping 8.8.8.8
PING 8.8.8.8 (8.8.8.8) 56(84) bytes of data.
64 bytes from 8.8.8.8: icmp_seq=1 ttl=120 time=7.25 ms
64 bytes from 8.8.8.8: icmp_seq=4 ttl=120 time=7.00 ms
^C
--- 8.8.8.8 ping statistics ---
8 packets transmitted, 2 received, 75% packet loss, time 7098ms
rtt min/avg/max/mdev = 7.000/7.123/7.247/0.123 ms
```

LIN·BIT

# Macro-Benchmarking

- **Applying the application workload fully as it would flow through the data path**
- **Run explicitly as a test to observe system metrics while it is applied**



LIN:BIT

# Workload Characterization

- **Determining quantitative aspects of the production workload to better simulate and observe how it performs on the system**

```
julia@julia-XPS-15-9510:~$ ./real_application.sh
Starting the real program now...
Program execution in progress...
(I am writing 512 sequential bytes to the disk, by the way...
Running... 3... 2... 1...
julia@julia-XPS-15-9510:~$
```

```
julia@julia-XPS-15-9510:~$ dd if=/dev/zero of=/dev/null obs=bytes bs=512 count=1
```

# Performance Tuning

- **Changing aspects of the environment and configuration with the intent to improve performance**
- **Especially disruptive/risky for production systems, so prior diagnostics via other methods are critical**

```
root@sar-1:/home/vagrant# ps -eo pid,ppid,ni,comm | grep backup.sh
1223286 1217870  -2 backup.sh
root@sar-1:/home/vagrant# renice -n 5 1223286
1223286 (process ID) old priority -2, new priority 5
root@sar-1:/home/vagrant#
```

# The USE Method

- **Developed by Brendan Gregg**
- **Focus on *Utilization*, *Saturation*, and *Errors* of system resources to quickly diagnose performance issues**
- **Learn more at www.brendangregg.com/usemethod.html**
- **Or his book, *Systems Performance*, 2nd ed.**

LIN·BIT

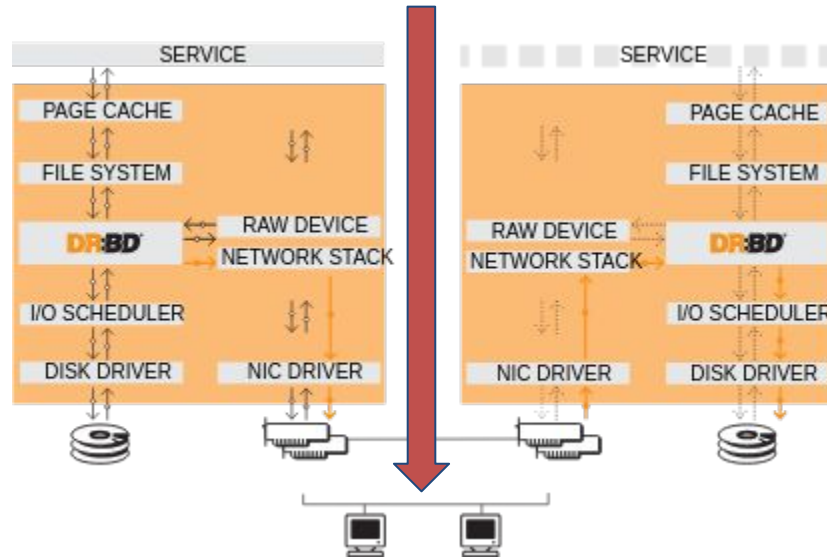# Functional Block Diagram

# The USE Method, Cont'd

- **Utilization**: The percentage of time the resource is doing work
- **Saturation**: The degree to which a resource has more work than it can process
- **Errors**: An discrete incident of a system not working as intended. In this case, refers to logged/loggable errors.

LIN·BIT

# Demo #1

# Demo #2

# Demo #3

# Demo #4

# Demo #5

# Demo #6

# Considerations for Disk I/O Benchmarking

- **Random vs sequential**
- **Ratio of reads/writes**
- **Size of individual writes performed**
- ***Working Set Size*: How much memory is needed by the application to perform the work**
- **Flash SSDs vs rotational HDDs**
- **RAID configuration (striped, parity?)**

LIN:BIT

# Disk Performance Tools & Metrics

- **iostat** - used to measure disk I/O (or determine if disks are performing I/O at all). Use with the -x flag
- **iowait** - a measure of the time CPUs are waiting for disk I/O to complete. Can be misleading!
- **smartctl** - can be used to report health metrics from disks, if supported by the disks used

LIN·BIT

# Where to go from here?

- **perf:  The Linux CPU profiler. Lightweight, powerful**
- **eBPF: Extended Berkeley Packet Filter. Kernel technology (available since 4.4) to run programs in kernel space and has numerous use cases for observability, tracing and profiling**

LIN:BIT

# Thanks for listening!

**Julia Iacoviello**
**julia@linbit.com**