



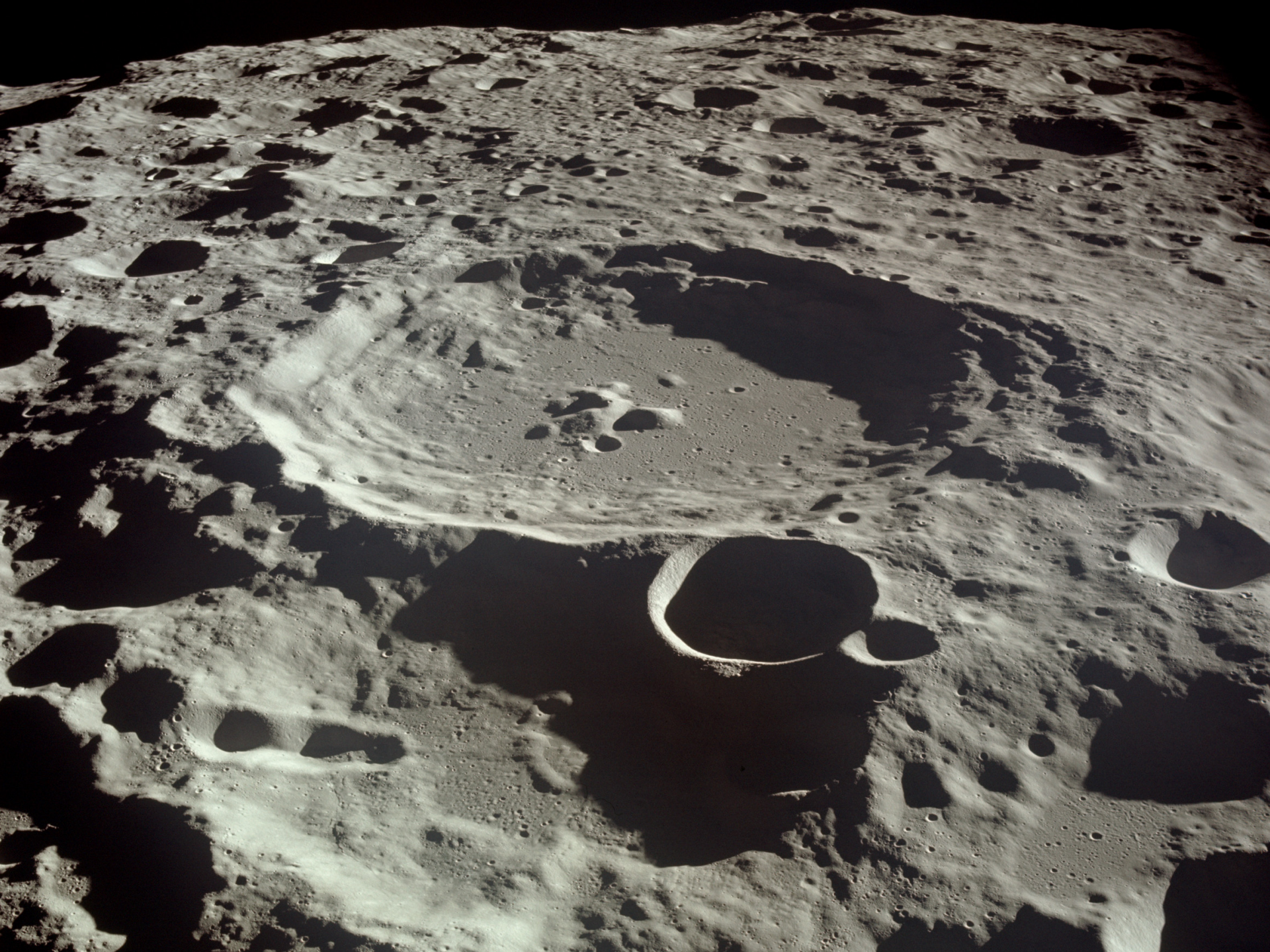
# System Methodology

## Holistic Performance Analysis on Modern Systems

Brendan Gregg

*Senior Performance Architect*

**NETFLIX**



# Apollo LMGC performance analysis

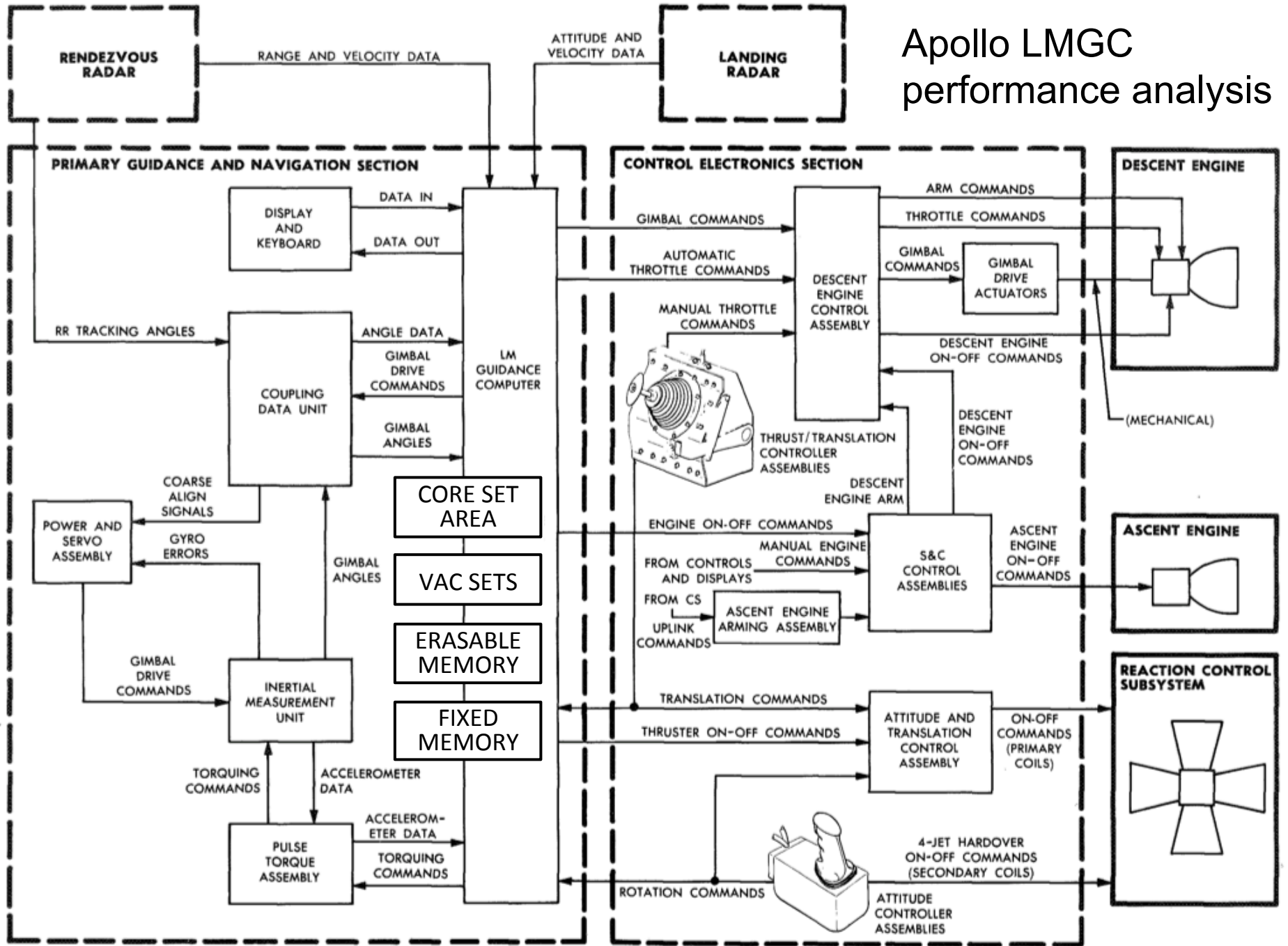
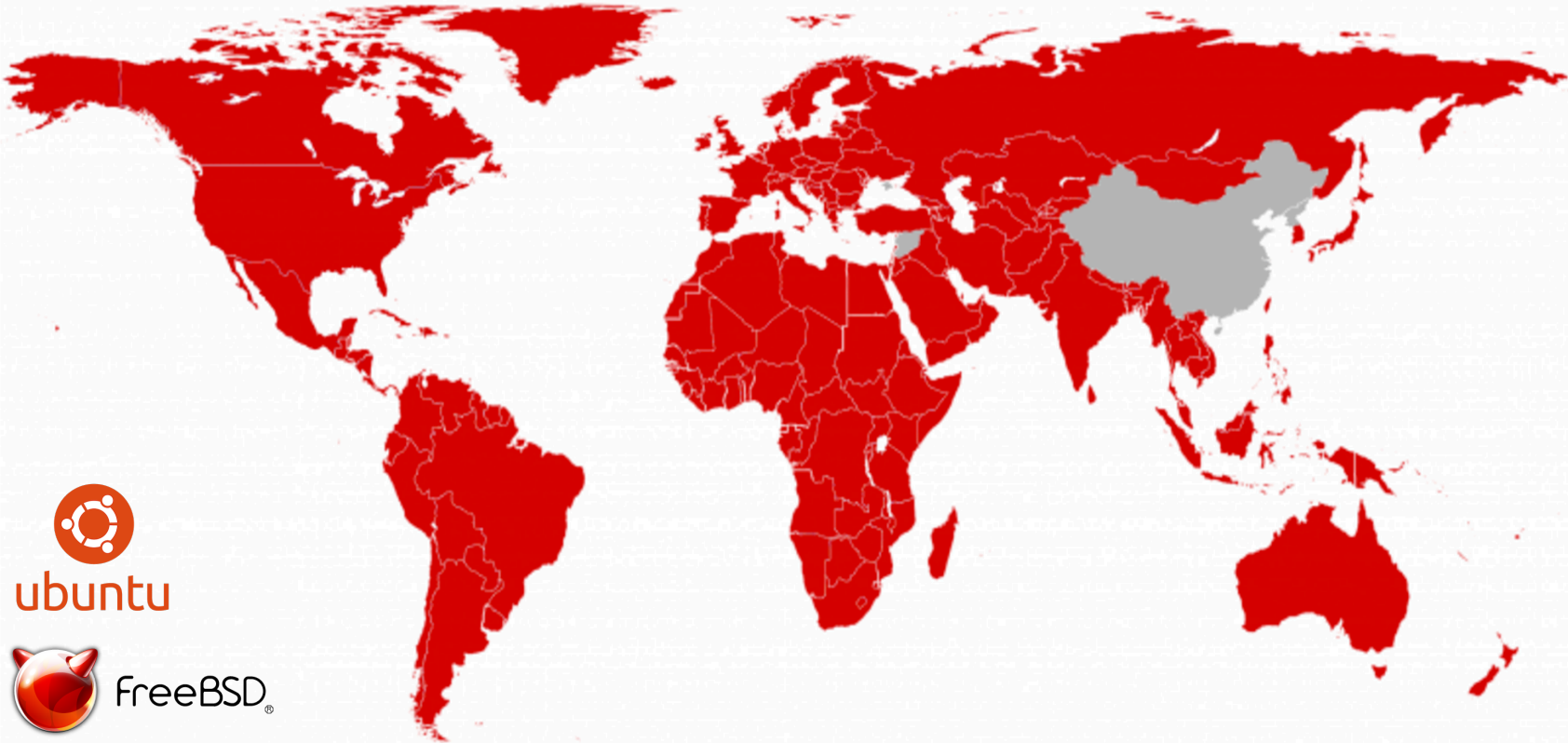


Figure 3-2.4. Primary Guidance Path - Simplified Block Diagram

# NETFLIX

REGIONS WHERE NETFLIX IS AVAILABLE



Background

# History

- System Performance Analysis up to the '90s:
  - Closed source UNIXes and applications
  - Vendor-created metrics and performance tools
  - Users interpret given metrics
- Problems
  - Vendors may not provide the best metrics
  - Often had to *infer*, rather than *measure*
  - Given metrics, what do we do with them?

```
# ps alx
F S UID    PID   PPID  CPU  PRI  NICE  ADDR  SZ  WCHAN  TTY  TIME  CMD
3 S  0      0     0    0    0    20   2253  2   4412  ?   186:14  swapper
1 S  0      1     0    0   30   20   2423  8   46520  ?    0:00  /etc/init
1 S  0     16     1    0   30   20   2273 11   46554  co   0:00  -sh
[...]
```

# Today

## 1. Open source

- Operating systems: Linux, BSDs, illumos, etc.
- Applications: source online (Github)

## 2. Custom metrics

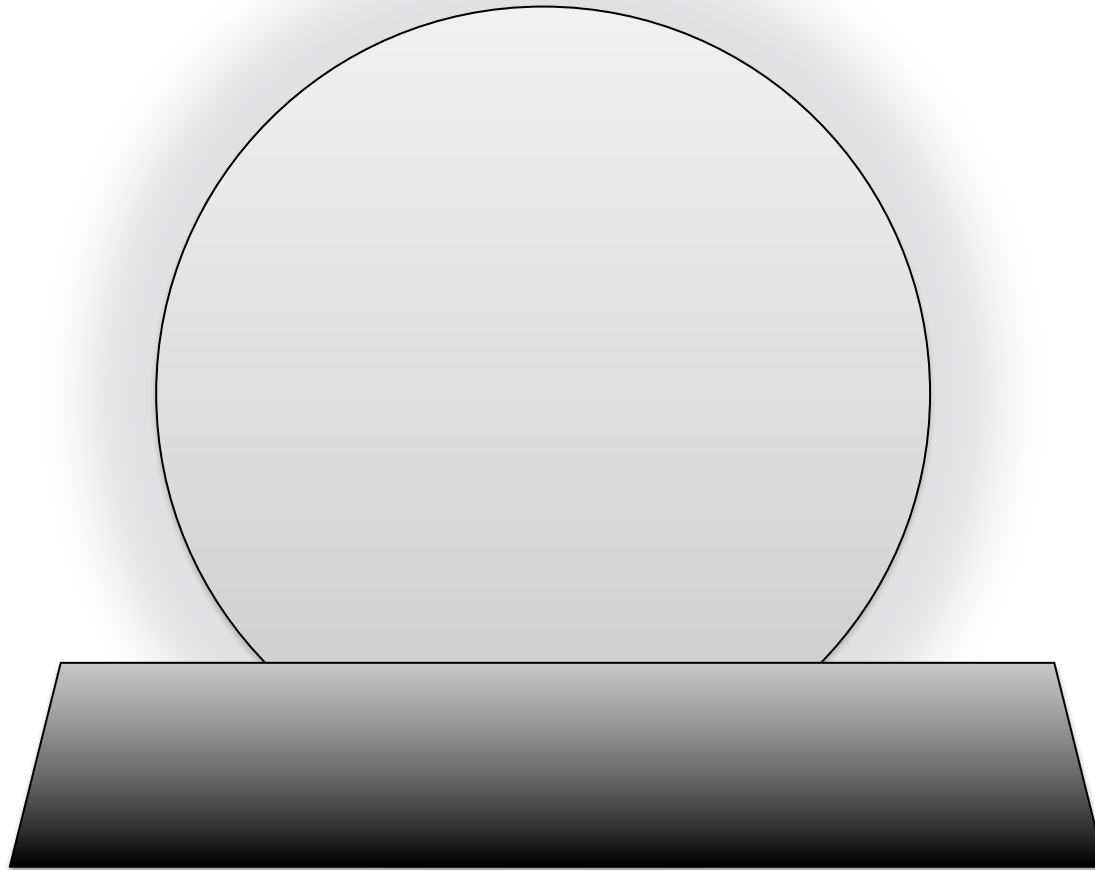
- Can patch the open source, or,
- Use dynamic tracing (open source helps)

## 3. Methodologies

- Start with the questions, then make metrics to answer them
- Methodologies can pose the questions

Biggest problem with dynamic tracing has been what to do with it.  
Methodologies guide your usage.

# Crystal Ball Thinking





# *Anti-Methodologies*

# Street Light *Anti*-Method

1. Pick observability tools that are
  - Familiar
  - Found on the Internet
  - Found at random
2. Run tools
3. Look for obvious issues



# Drunk Man *Anti*-Method

- ~~Drink~~ Tune things at random until the problem goes away



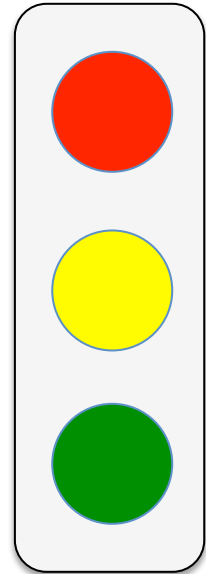
# Blame Someone Else *Anti-Method*

1. Find a system or environment component you are not responsible for
2. Hypothesize that the issue is with that component
3. Redirect the issue to the responsible team
4. When proven wrong, go to 1



# Traffic Light *Anti*-Method

1. Turn all metrics into traffic lights
  2. Open dashboard
  3. Everything green? No worries, mate.
- Type I errors: red instead of green
    - team wastes time
  - Type II errors: green instead of red
    - performance issues undiagnosed
    - team wastes more time looking elsewhere



Traffic lights are suitable for *objective* metrics (eg, errors), not *subjective* metrics (eg, IOPS, latency).

# Methodologies

# Performance Methodologies

- For system engineers:
  - ways to analyze unfamiliar systems and applications
- For app developers:
  - guidance for metric and dashboard design



Collect your own toolbox of methodologies

## System Methodologies:

- Problem statement method
- Functional diagram method
- Workload analysis
- Workload characterization
- Resource analysis
- USE method
- Thread State Analysis
- On-CPU analysis
- CPU flame graph analysis
- Off-CPU analysis
- Latency correlations
- Checklists
- Static performance tuning
- Tools-based methods

...

# Problem Statement Method

1. What makes you **think** there is a performance problem?
2. Has this system **ever** performed well?
3. What has **changed** recently?
  - software? hardware? load?
4. Can the problem be described in terms of **latency**?
  - or run time. not IOPS or throughput.
5. Does the problem affect **other** people or applications?
6. What is the **environment**?
  - software, hardware, instance types?  
versions? config?



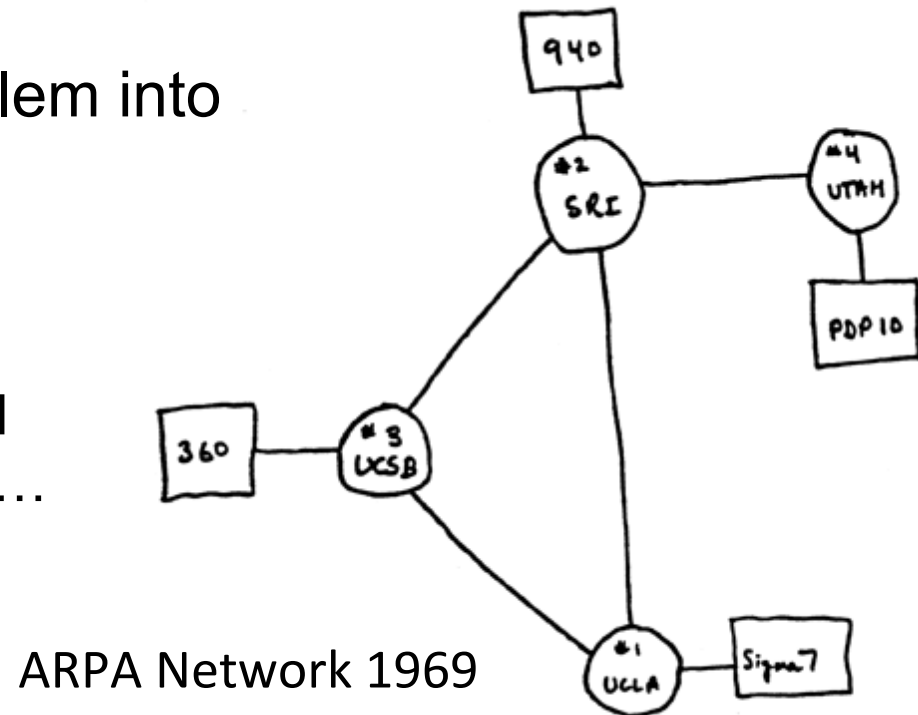


# Functional Diagram Method

1. Draw the functional diagram
2. Trace all components in the data path
3. For each component, check performance

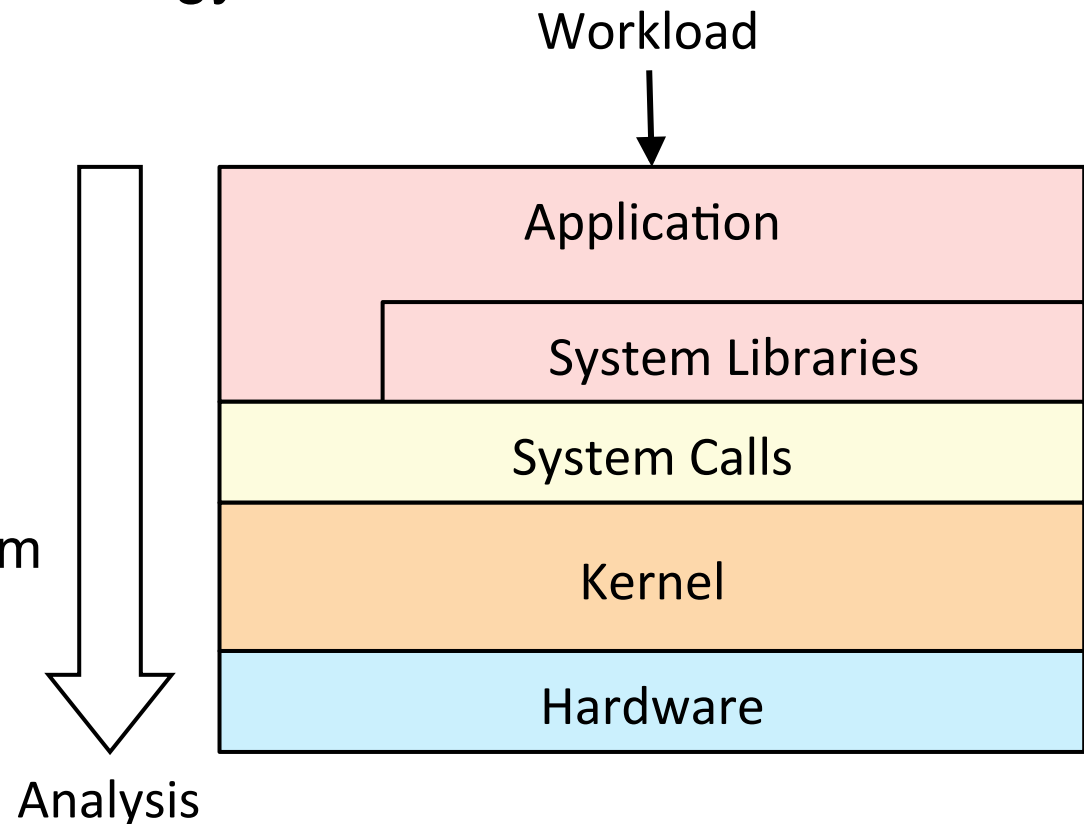
Breaks up a bigger problem into smaller, relevant parts

Eg, imagine throughput between the UCSB 360 and the UTAH PDP10 was slow...



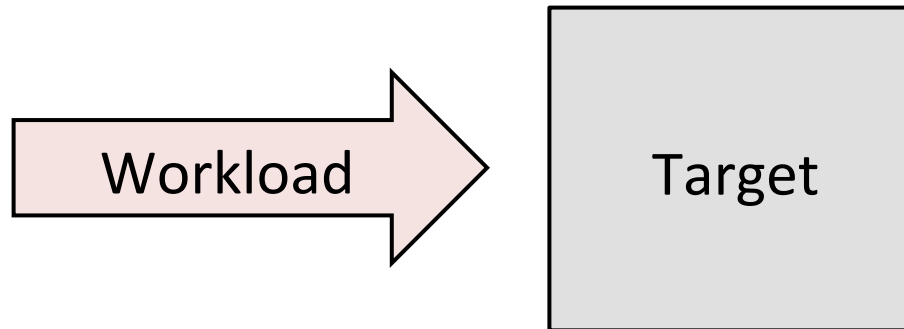
# Workload Analysis

- Begin with application metrics & context
- A **drill-down** methodology
- Pros:
  - Proportional, accurate metrics
  - App context
- Cons:
  - App specific
  - Difficult to dig from app to resource



# Workload Characterization

- Check the workload: **who, why, what, how**
  - not resulting performance



- Eg, for CPUs:
  1. Who: which PIDs, programs, users
  2. Why: code paths, context
  3. What: CPU instructions, cycles
  4. How: changing over time

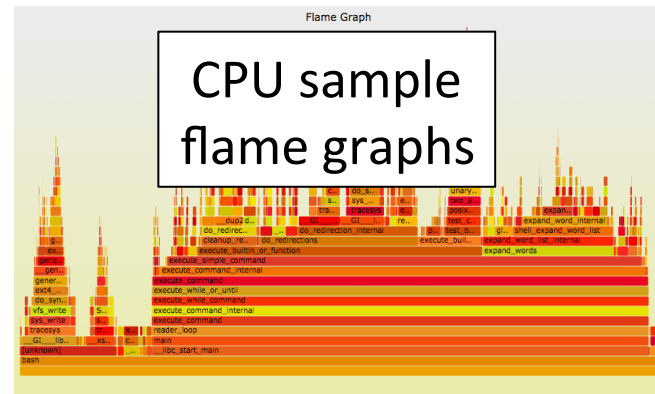
# Workload Characterization: CPUs

## Who

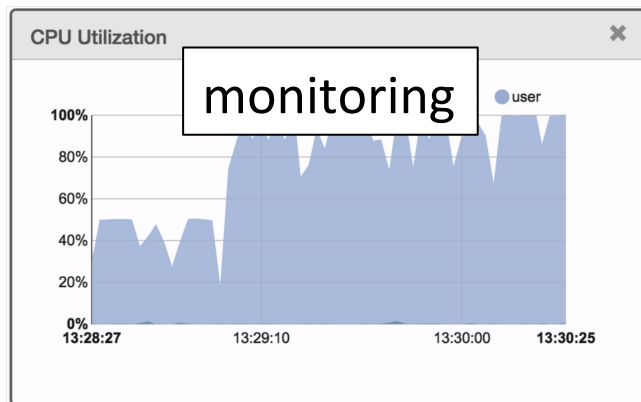
PID	USER	VIRT	RES	CPU%	MEM%	TIME+	Command
27983	root	3233M	20M	0.0	0.0	0:02:10.50	/usr/lib/jvm/java
28004	root	3233M	20M	0.0	0.0	0:02:02.60	/usr/lib/jvm/java
28173	root	63488	49M	0.0	0.0	0:02:02.68	ab -k -c 100 -n 1
28170	root	24660	21M	0.0	0.0	0:00:00.62	htop
2730	root	202M	58600	0.0	0.0	2h31:25	/apps/epic/perl/b
2752	root	151M	10308	0.0	0.1	1h48:36	postgres: bgregg-
28000	root	3233M	204M	0.0	2.7	0:00:00.26	/usr/lib/jvm/java
1	root	24320	2256	0.0	0.0	0:01:29	/sbin/init
341	root	17236	632	0.0	0.0	0:00:00.04	upstart-udev-brid
346	root	21600	1304	0.0	0.0	0:00:00.06	/sbin/udev --dae
357	root	23944	1164	0.0	0.0	0:00:00.21	dbus-daemon --sys
408	root	21464	792	0.0	0.0	0:00:00.00	/sbin/udev --dae
549	root	15192	392	0.0	0.0	0:00:00.00	upstart-socket-br
612	root	7268	1028	0.0	0.0	0:00:00.24	dhclient3 -e IF_M
644	root	50036	2920	0.0	0.0	0:00:00.06	/usr/sbin/sshd -D
772	root	14508	956	0.0	0.0	0:00:00.00	/sbin/getty -8 38
777	root	14508	952	0.0	0.0	0:00:00.00	/sbin/getty -8 38
785	root	14508	952	0.0	0.0	0:00:00.00	/sbin/getty -8 38

top

## Why



## How



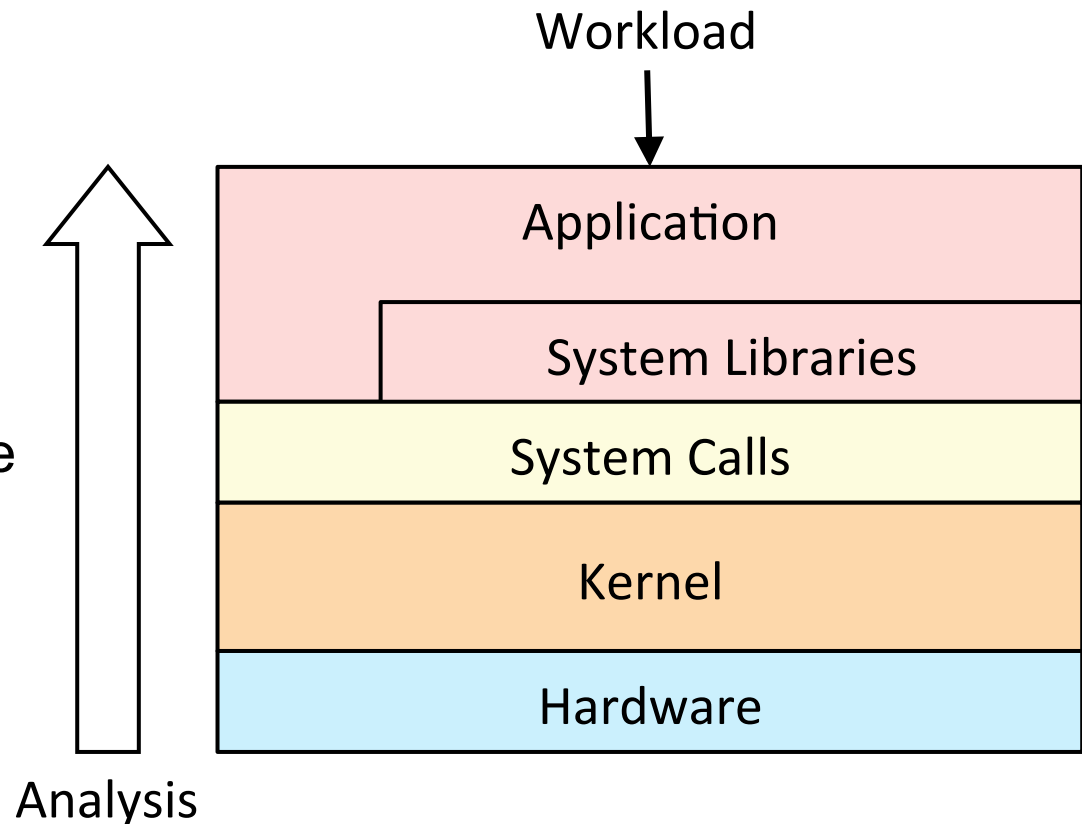
## What

```
root@lgud-bgregg:~# perf stat -a -d sleep 10
```

Performance counter	stat	side'
39996.388668 task-clock	PMCs	# 3.999 CPUs ut
1,026,540 context-switches		# 0.026 M/sec
193,563 cpu-migrations		# 0.005 M/sec
4,835 page-faults		# 0.121 K/sec
83,859,543,001 cycles		# 2.097 GHz
61,028,919,136 stalled-cycles-frontend		# 72.78% frontend
50,812,852,642 stalled-cycles-backend		# 60.59% backend
52,969,864,055 instructions		# 0.63 insns p
		# 1.15 stalled
10,223,584,755 branches		# 255.613 M/sec
376,529,869 branch-misses		# 3.68% of all
0 L1-dcache-loads		# 0.000 K/sec
1,339,950,792 L1-dcache-load-misses		# 0.00% of all
762,761,193 LLC-loads		# 19.071 M/sec

# Resource Analysis

- Typical approach for system performance analysis: begin with system tools & metrics
- Pros:
  - Generic
  - Aids resource perf tuning
- Cons:
  - Uneven coverage
  - False positives

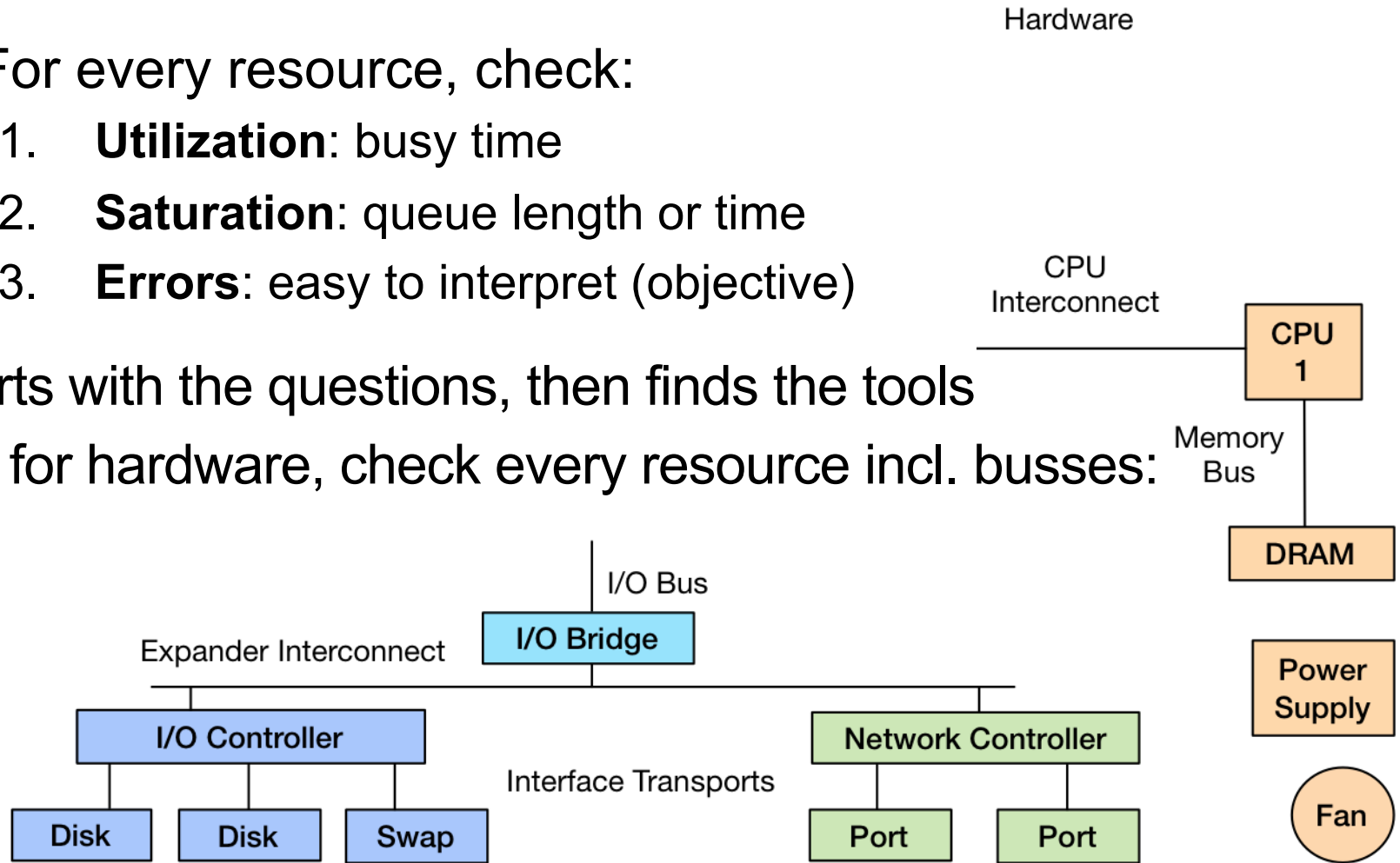


# The USE Method

- For every resource, check:
  1. **Utilization**: busy time
  2. **Saturation**: queue length or time
  3. **Errors**: easy to interpret (objective)

Starts with the questions, then finds the tools

Eg, for hardware, check every resource incl. busses:



# USE Method: Rosetta Stone of Performance Checklists

The following [USE Method](#) example checklists are automatically generated from the individual pages for: [Linux](#), [Solaris](#), [Mac OS X](#), and [FreeBSD](#). These analyze the performance of the physical host. You can customize this table using the checkboxes on the right.

- Linux
  - Solaris
  - FreeBSD
  - Mac OS X
- 

There are some additional USE Method example checklists not included in this table: the [SmartOS](#) checklist, which is for use within an OS virtualized guest, and the [Unix 7th Edition](#) checklist for historical interest.

For general purpose operating system differences, see the [Rosetta Stone for Unix](#), which was the inspiration for this page.

## Hardware Resources

<http://www.brendangregg.com/USEmethod/use-rosetta.html>

Resource	Metric	Linux	Solaris	FreeBSD
CPU	errors	<code>perf</code> (LPE) if processor specific error events (CPC) are available; eg, AMD64's "04Ah Single-bit ECC Errors Recorded by Scrubber" [4]	<code>fmadm faulty</code> ; <code>cpustat</code> (CPC) for whatever error counters are supported (eg, thermal throttling)	<code>dmesg</code> ; <code>/var/log/messages</code> ; <code>pmcstat</code> for PMC and whatever error counters are supported (eg, thermal throttling)
CPU	saturation	system-wide: <code>vmstat 1, "r" &gt; CPU count</code> [2]; <code>sar -q, "runq-sz" &gt; CPU count</code> ; <code>dstat -p, "run" &gt; CPU count</code> ; per-process: <code>/proc/PID/schedstat</code> 2nd field ( <code>sched_info.run_delay</code> ); <code>perf sched latency</code> (shows "Average" and "Maximum" delay per-schedule); dynamic tracing, eg, SystemTap <code>schedtimes.stp "queued(us)"</code> [3]	system-wide: <code>uptime</code> , <a href="#">load averages</a> ; <code>vmstat 1, "r"</code> ; DTrace <code>dispqlen.d (DTI)</code> for a better "vmstat r"; per-process: <code>prstat -mLc 1, "LAT"</code>	system-wide: <code>uptime</code> , "load averages" <code>&gt; CPU count</code> ; <code>vmstat 1, "procs:r" &gt; CPU count</code> ; per-cpu: DTrace to profile CPU run queue lengths [1]; per-process: DTrace of scheduler events [2]
CPU	utilization	system-wide: <code>vmstat 1, "us" + "sy" + "st"</code> ; <code>sar -u</code> , sum fields except "%idle" and "%iowait"; <code>dstat -c</code> , sum fields except "idl" and "wai"; per-cpu: <code>mpstat -P ALL 1</code> , sum fields except "%idle" and "%iowait"; <code>sar -P ALL</code> , same as <code>mpstat</code> ; per-process: <code>top, "%CPU"</code> ; <code>htop, "CPU%"</code> ; <code>ps -o pcpu</code> ; <code>pidstat 1, "%CPU"</code> ; per-kernel-thread: <code>top/htop ("K" to toggle)</code> , where <code>VIRT == 0</code> (heuristic). [1]	per-cpu: <code>mpstat 1, "usr" + "sys"</code> ; system-wide: <code>vmstat 1, "us" + "sy"</code> ; per-process: <code>prstat -c 1 ("CPU" == recent)</code> , <code>prstat -mLc 1 ("USR" + "SYS")</code> ; per-kernel-thread: <code>lockstat -Ii rate</code> , DTrace <code>profile stack()</code>	system-wide: <code>vmstat 1, "us" + "sy"</code> ; per-cpu: <code>vmstat -P</code> ; per-process: <code>top, "WCPU"</code> for weighted and recent usage; per-kernel-process: <code>top -s, "WCPU"</code>
CPU interconnect	errors	LPE (CPC) for whatever is available	<code>cpustat</code> (CPC) for whatever is available	<code>pmcstat</code> and relevant PMCs for whatever is available
CPU	saturation	LPE (CPC) for stall	<code>cpustat</code> (CPC) for stall	<code>pmcstat</code> and relevant PMCs for

# USE Method: Unix 7th Edition Performance Checklist

Out of curiosity, I've developed a [USE Method](#)-based performance checklist for [Unix 7th Edition](#) on a [PDP-11/45](#), which I've been running via a PDP [simulator](#). 7th Edition is from 1979, and was the first Unix with `iostat(1M)` and `pstat(1M)`, enabling more serious performance analysis from shipped tools. Were I to write a checklist for earlier Unixes, it would contain many more "unknowns".

I often work on the [illumos](#) kernel, a direct descendant of Unix which contains some original AT&T code. It's been interesting to study this earlier version, and see familiar code that has survived over 30 years of development.

Example screenshots from various tools are shown at the end of this page.



*PDP 11/70 front panel (similar to the 11/45)*

## Physical Resources

component	type	metric
CPU	utilization	system-wide: <code>iostat 1</code> , utilization is "user" + "nice" + "system"; per-process: <code>ps alx</code> , "CPU" shows recent CPU usage (max 255), and "TIME" shows cumulative minutes:seconds of CPU time
CPU	saturation	<code>ps alx   awk '\$2 == "R" { r++ } END { print r - 1 }'</code> , shows the number of runnable processes
CPU	errors	console message if lucky, otherwise panic
Memory capacity	utilization	system-wide: unknown [1]; per-type: unknown [2]; per-process: <code>ps alx</code> , "SZ" is the in-core (main memory) in blocks (512 bytes); <code>pstat -p</code> , "SIZE" is in-core size, in units of core clicks (64 bytes) and printed in octal!
Memory capacity	saturation	system-wide: <code>iostat 1</code> , sustained "tpm" may be caused by swapping to disk; significant delays as processes wait for space to swap in
Memory capacity	errors	<code>malloc()</code> returns 0; ENOMEM
Disk I/O	utilization	system-wide: <code>iostat -i 1</code> , "IO active" plus "IO wait" percents; per-disk-controller: <code>iostat -i 1, RF, RK, RP</code> "active" percents; rough estimate using <code>iostat 1</code> , and "tpm" for transactions per minute on expected max; per-disk: listen to each rattle; unknown from Unix, unless only 1 disk per controller; per-process: unknown
Disk I/O	saturation	unknown [3]
Disk I/O	errors	might get a console message, eg, "err on dev", "ECC on dev" or "no space on dev", otherwise unknown [4]
Tape I/O	utilization	look at tape drives and watch them spin [5]



# Apollo Guidance Computer

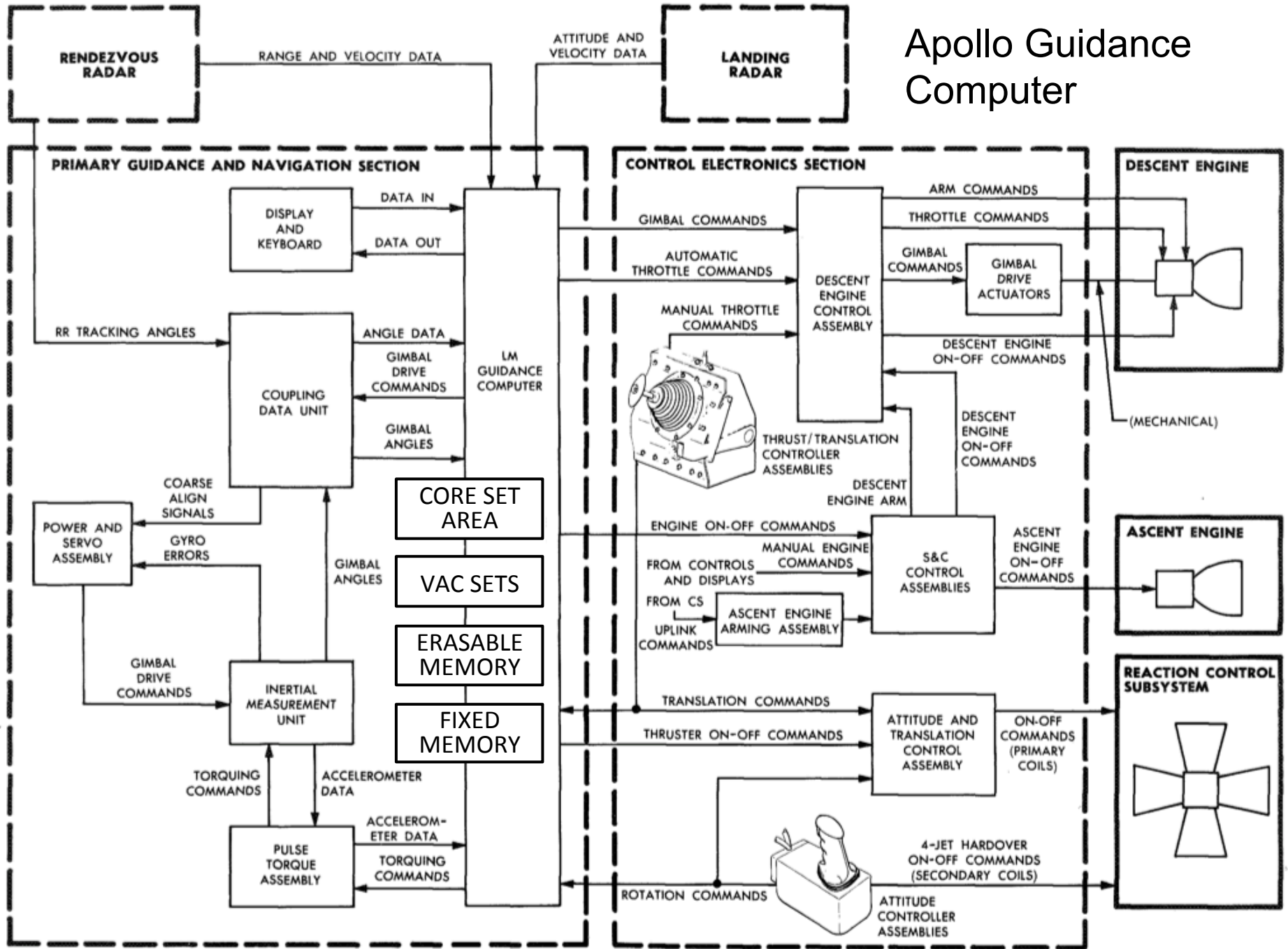
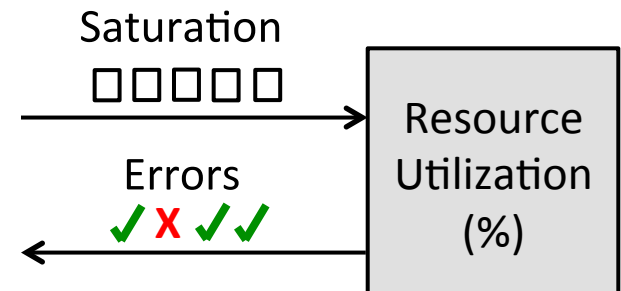


Figure 3-2.4. Primary Guidance Path - Simplified Block Diagram

# USE Method: Software

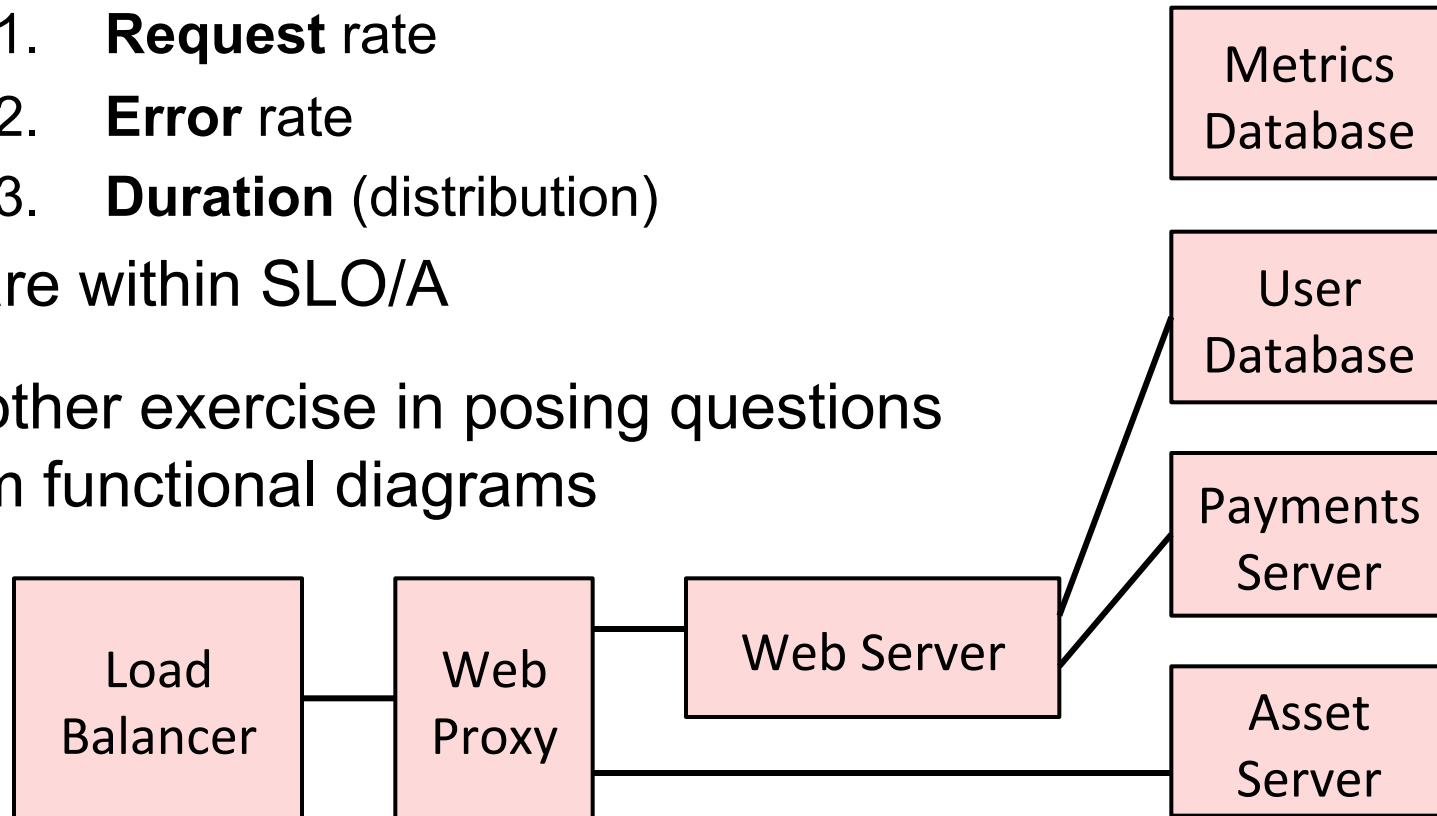
- USE method can also work for software resources
  - kernel or app internals, cloud environments
  - small scale (eg, locks) to large scale (apps). Eg:
- Mutex locks:
  - utilization → lock hold time
  - saturation → lock contention
  - errors → any errors
- Entire application:
  - utilization → percentage of worker threads busy
  - saturation → length of queued work
  - errors → request errors



# RED Method

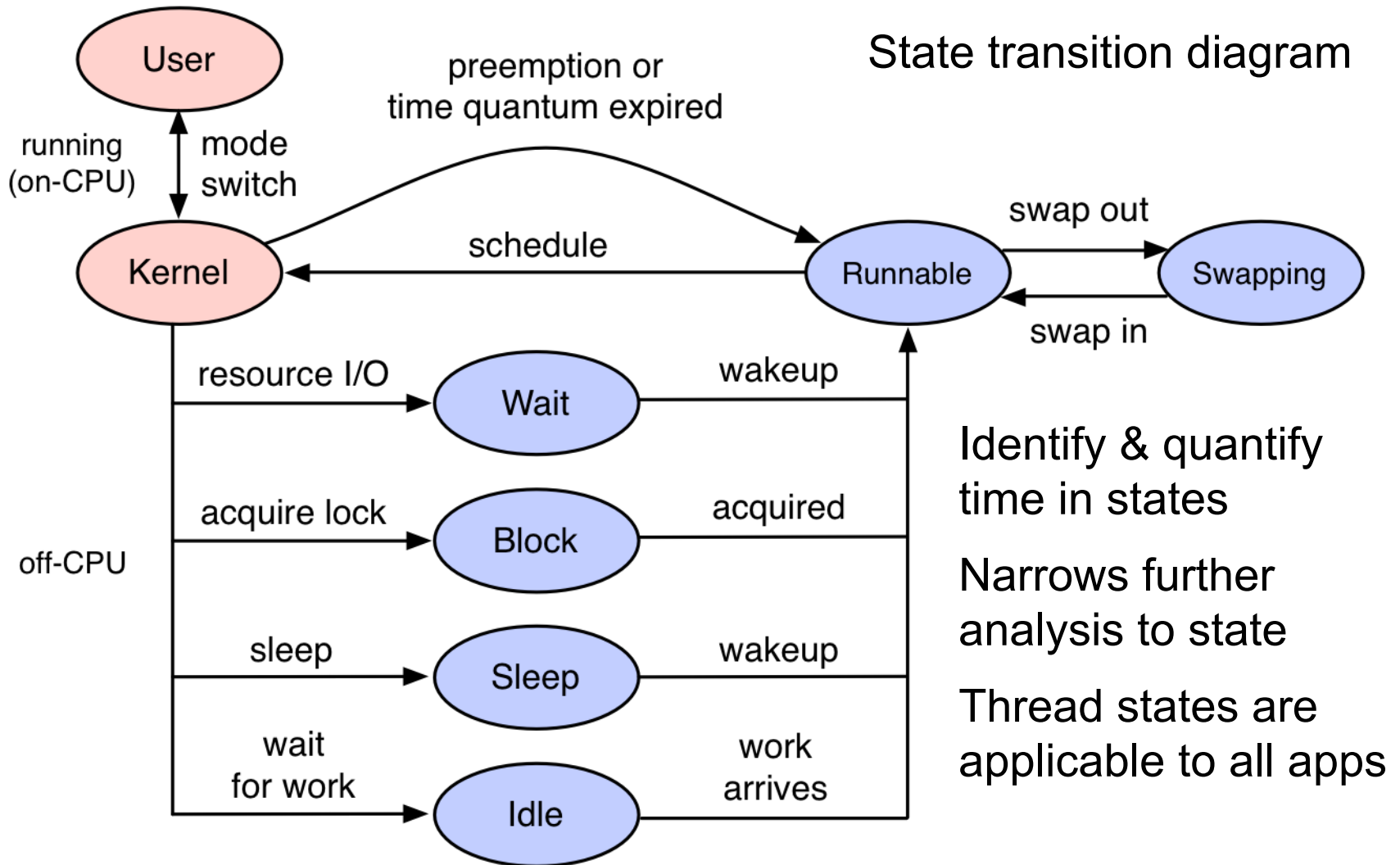
- For every service, check that:
  1. **Request** rate
  2. **Error** rate
  3. **Duration** (distribution)are within SLO/A

Another exercise in posing questions from functional diagrams



# Thread State Analysis

State transition diagram



Identify & quantify time in states

Narrows further analysis to state

Thread states are applicable to all apps

# TSA: eg, Solaris

1) \$ `prstat -mLc 1`

```

PID USERNAME  USR  SYS  TRP  TFL  DFL  LCK  SLP  LAT  VCX  ICX  SCL  SIG  PROCESS/LWPID
45747 1000      35  28  0.0  0.0  22  0.0  16  0.1  216  93  38K  0  beam.smp/192
[...]
```

Executing

Anon.  
Paging

Lock  
+Idle

Sleep  
+Idle

Runnable

2)



Fields	States	Analysis	Actions
USR+SYS	Executing	<ol style="list-style-type: none"> <li>Profile stacks using DTrace; Flame Graphs</li> <li>Check CPU stall cycles: <code>cpustat</code>, DTrace</li> <li>If SYS time, analyze syscalls using DTrace</li> </ol>	<ul style="list-style-type: none"> <li>Look for inefficiencies</li> <li>Look for tunables/config in active code</li> </ul>
DFL	Anon. Paging	<ol style="list-style-type: none"> <li>Confirm using: <code>vmstat -p 1, "api"</code></li> <li>Check system-wide memory free: <code>vmstat 1</code></li> <li>Check any resource controls; eg: <code>zonememstat</code></li> </ol>	<ul style="list-style-type: none"> <li>Upgrade memory</li> <li>Increase memory limits</li> <li>Look for leaks/growth</li> </ul>
LCK	Lock + Idle	<ol style="list-style-type: none"> <li>Coarse: profile CPU stacks and look for spins</li> <li>Analyze using DTrace <code>[p]lockstat</code> providers</li> <li>Separate locks and the Idle state using DTrace <code>sched:::off-cpu</code> with <code>ustack()</code></li> </ol>	<ul style="list-style-type: none"> <li>Check config</li> </ul>
SLP	Sleep + Idle	<ol style="list-style-type: none"> <li>Quick resource check: <code>iostat -xnz 1, nicstat 1</code></li> <li>Identify both sleep reason and separate from Idle: DTrace <code>sched:::off-cpu</code> with <code>ustack()</code> and <code>stack()</code></li> </ol>	<ul style="list-style-type: none"> <li>Tune or upgrade resource</li> </ul>
LAT	Runnable	<ol style="list-style-type: none"> <li>Check system CPU usage: <code>mpstat 1</code></li> <li>Check any resource controls; eg, <code>prctl, kstat -p caps:::cpucaps_zone*</code></li> <li>Check for <code>pbind/psets</code> limiting migrations</li> </ol>	<ul style="list-style-type: none"> <li>Upgrade CPUs</li> <li>Increase CPU limits</li> <li>Move/tune other load</li> <li>Unbind apps</li> </ul>

# TSA: eg, RSTS/E

RSTS: DEC OS  
from the 1970's

TENEX (1969-72)  
also had Control-T  
for job states

State Column (Job Status)		
RN	Run	Job is running or waiting to run.
RS	Residency	Job is waiting for residency. (The job has been swapped out of memory and is waiting to be swapped back in.)
BF	Buffers	Job is waiting for buffers (no space is available for I/O buffers).
SL	Sleep	Job is sleeping (SLEEP statement).
SR	Send/Receive	Job is sleeping and is a message receiver.
FP	File Processor	Job is waiting for file processing by the system (opening or closing a file, file search).
TT	Terminal	Job is waiting to perform output to a terminal.
HB	Hibernating	Job is detached and waiting to perform I/O to or from a terminal. (Someone must attach to the job before it can resume execution.)
KB	Keyboard	Job is waiting for input from a terminal.
^C	CTRL/C	Job is at command level, awaiting a command. (In other words, the keyboard monitor has displayed its prompt and is waiting for input.)
CR	Card Reader	Job is waiting for input from a card reader.
MT,MM, or MS	Magnetic Tape	Job is waiting for magnetic tape I/O.
LP	Line Printer	Job is waiting to perform line printer output.
DT	DEctape	Job is waiting for DEctape I/O.
DK,DM,DB, DP,DL,DR	Disk	Job is waiting to perform disk I/O.

# TSA: eg, OS X

## Instruments: Thread States

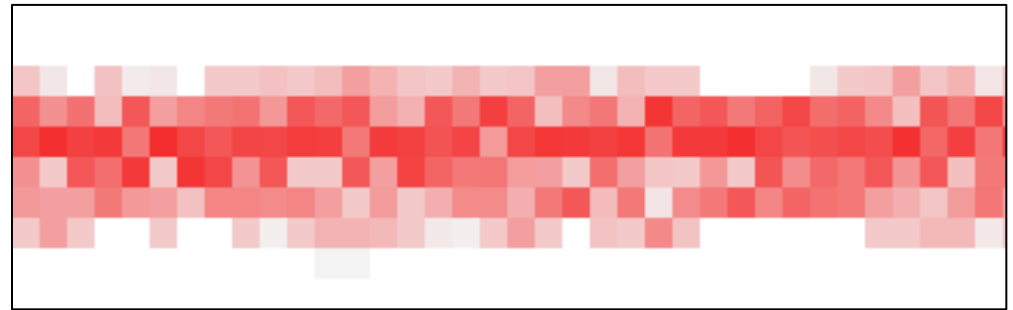
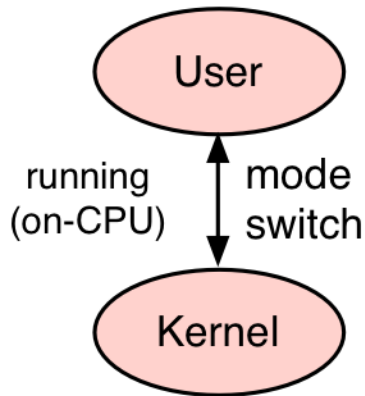
The screenshot shows the Instruments application interface. At the top, there's a toolbar with buttons for Record, Target (firefox (1027)), Inspection Range (00:00:42), View, Library, and Filter (Recorded Data). The main area displays a timeline with a Thread States track. A settings panel for Thread States is open, showing the following options:

- Target: firefox (1027)
- Track Display: Style: Thread States, Type: Stacked, Zoom: 4x
- Thread States:
  - Unknown
  - Waiting
  - Suspended
  - Requested to suspend
  - Running
  - On run queue
  - Waiting and uninterruptible
  - At termination
  - Idling processor
- Track Behavior:
  - Size track by thread count

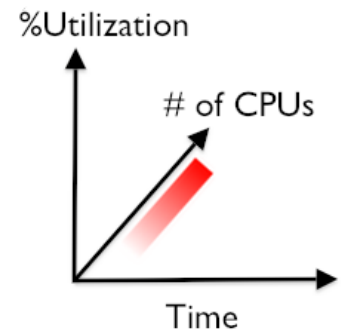
Below the Thread States track, a table provides summary statistics for the threads:

Alive	ms On CPU	Switches	Children	% Living Children
•	856	966	1	100%
	1	16	1	100%
•		29		
•	2,428,558	7,528	17	100%

# On-CPU Analysis



CPU Utilization  
Heat Map



1. Split into user/kernel states
  - /proc, vmstat(1)
2. Check CPU balance
  - mpstat(1), CPU utilization heat map
3. Profile software
  - User & kernel stack sampling (as a **CPU flame graph**)
4. Profile cycles, caches, busses
  - PMCs, CPI flame graph



# CPU Flame Graph Analysis

1. Take a CPU profile
2. Render it as a flame graph
3. Understand all software that is in  $>1\%$  of samples

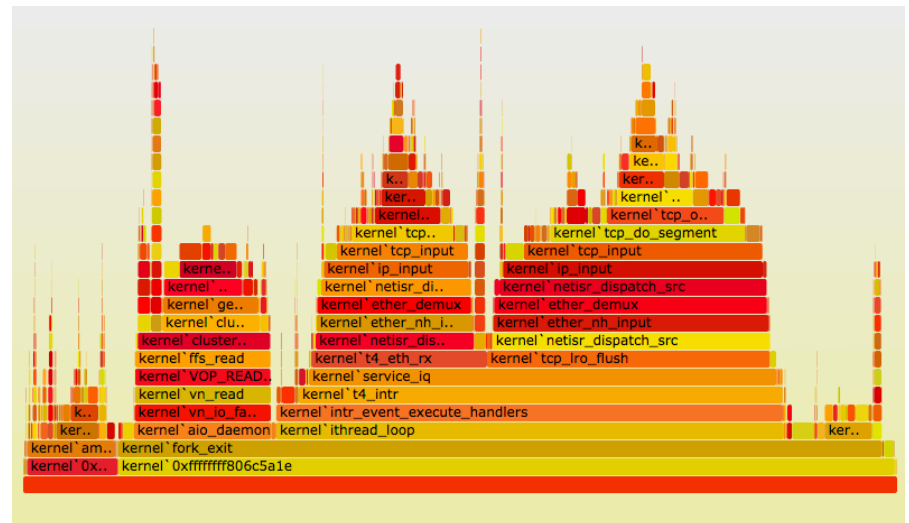
Discovers issues by their CPU usage

- Directly: CPU consumers
- Indirectly: initialization of I/O, locks, times, ...

Narrows target of study to only running code

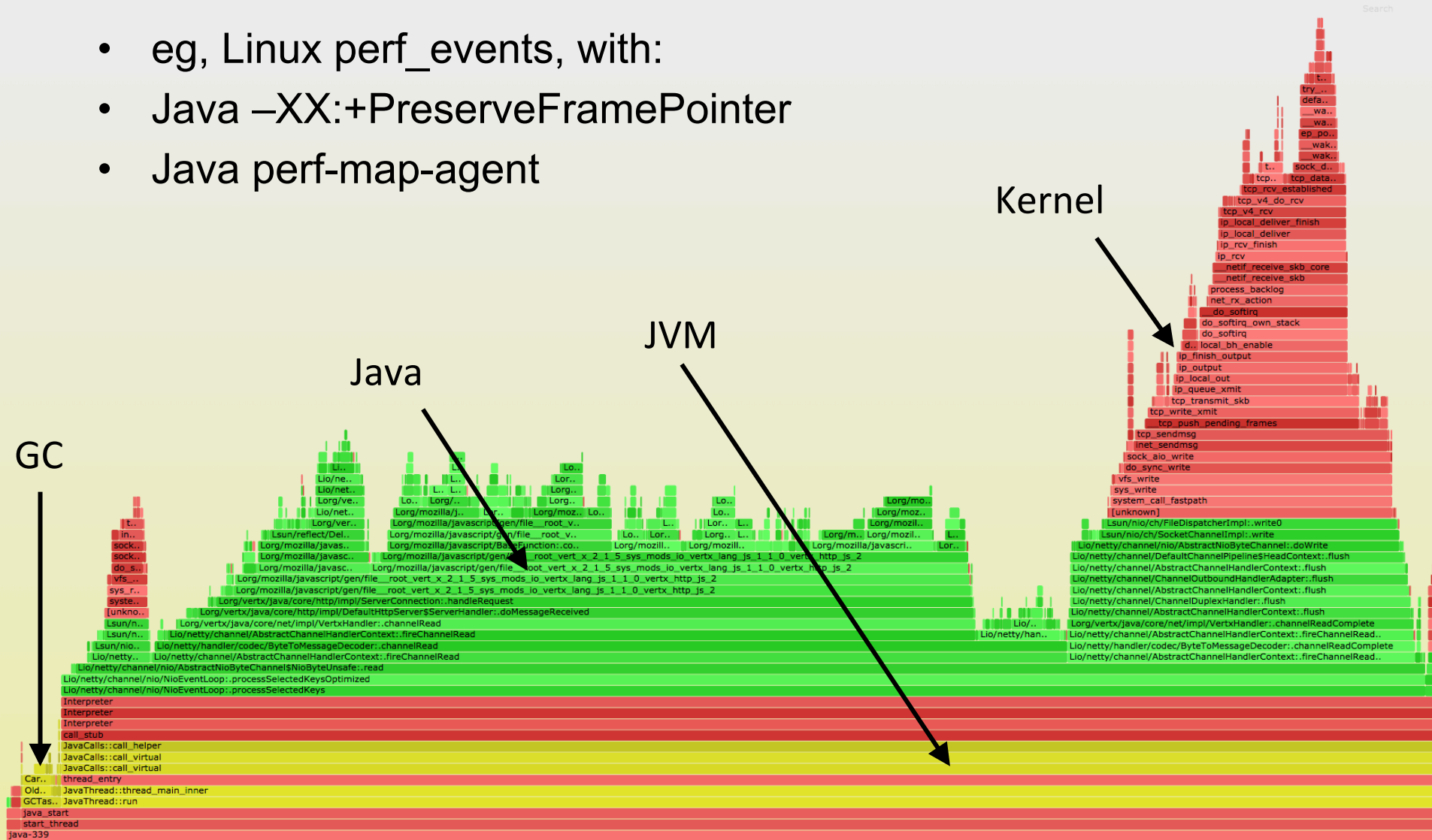
- See: "The Flame Graph", CACM, June 2016

Flame Graph



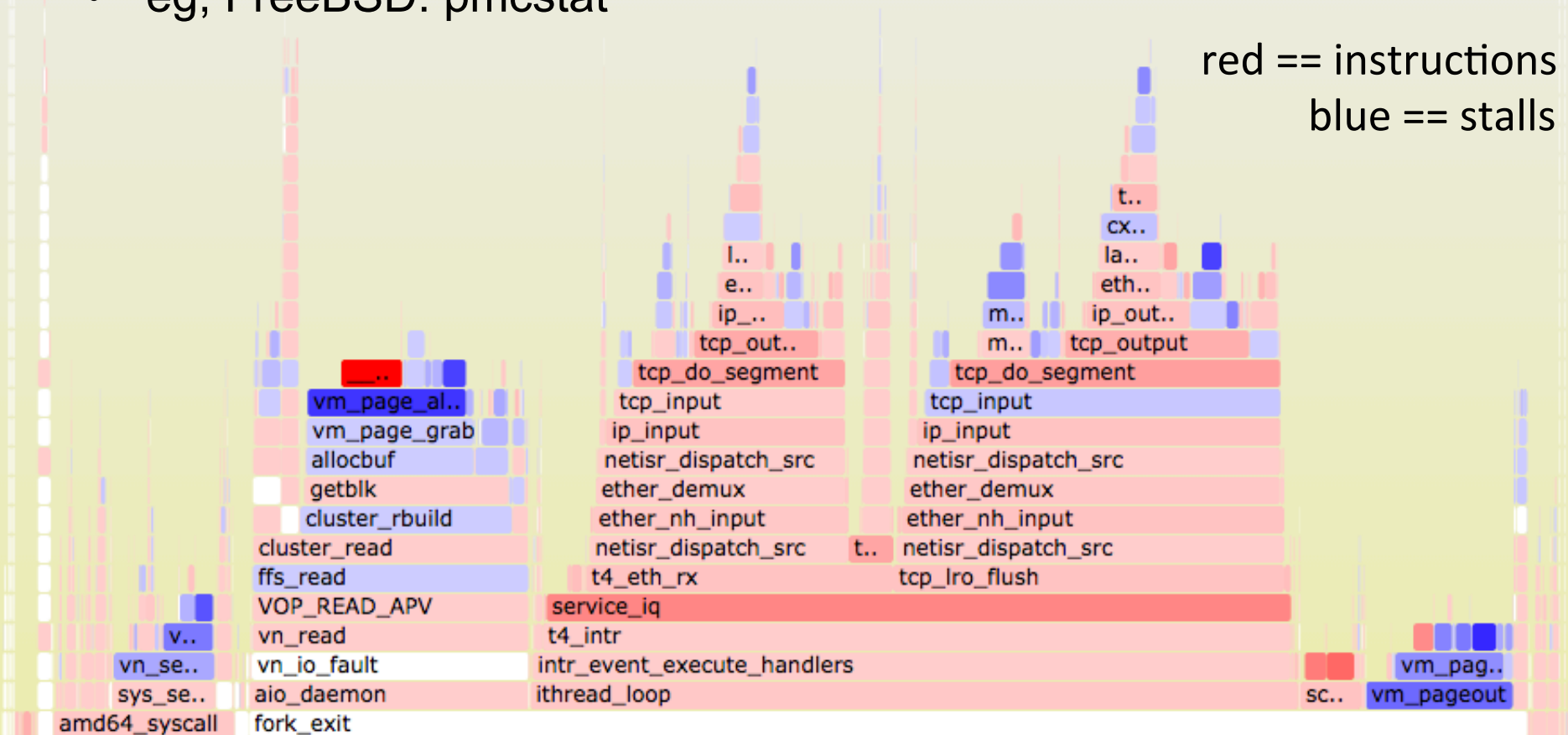
# Java Mixed-Mode CPU Flame Graph

- eg, Linux perf\_events, with:
- Java -XX:+PreserveFramePointer
- Java perf-map-agent

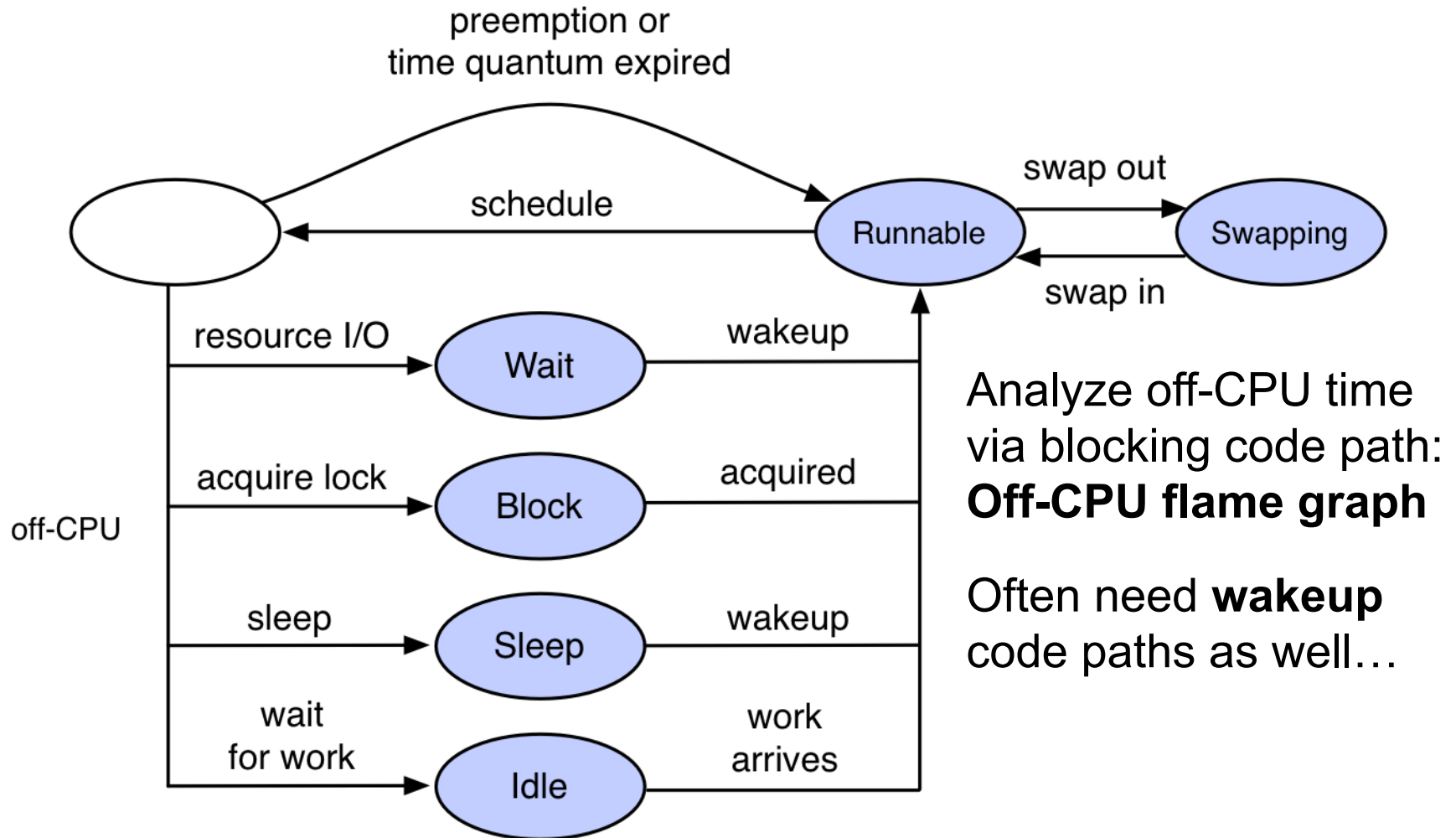


# CPI Flame Graph

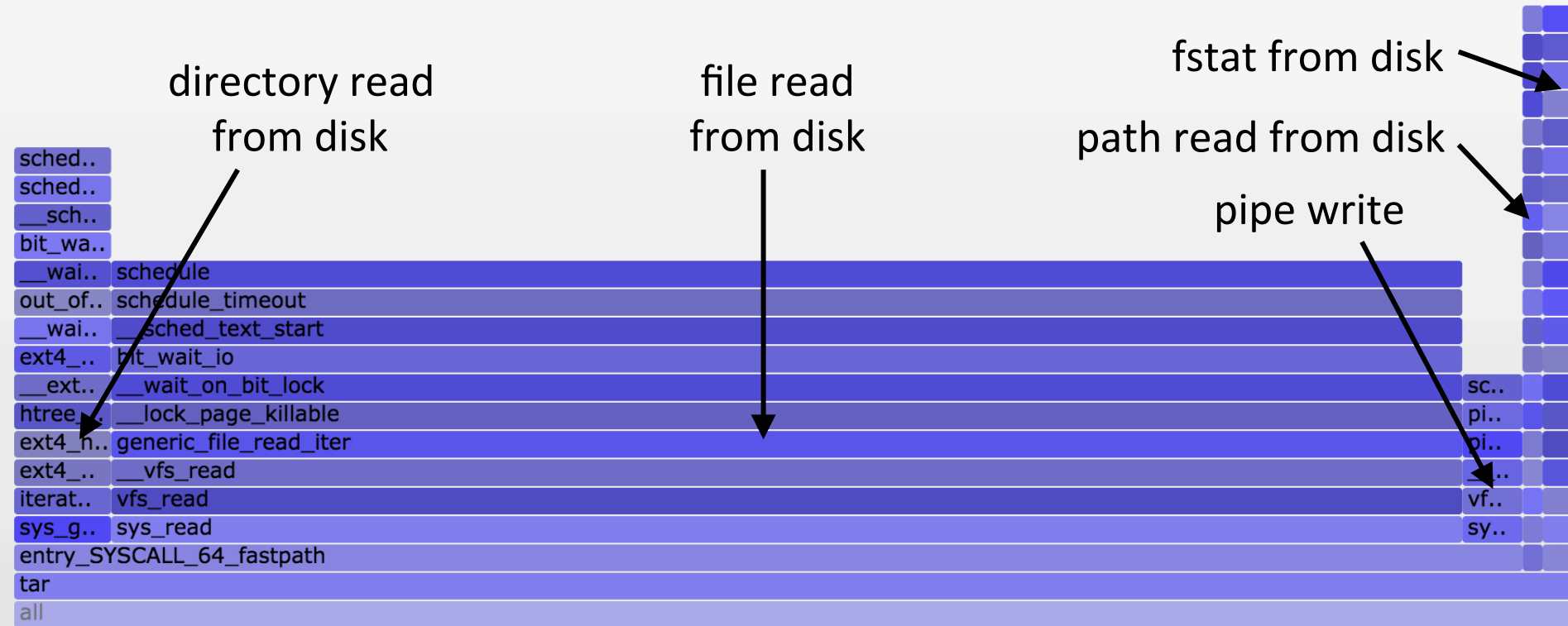
- Profile cycle stack traces and instructions or stalls separately
- Generate CPU flame graph (cycles) and color using other profile
- eg, FreeBSD: pmcstat



# Off-CPU Analysis



# Off-CPU Time Flame Graph



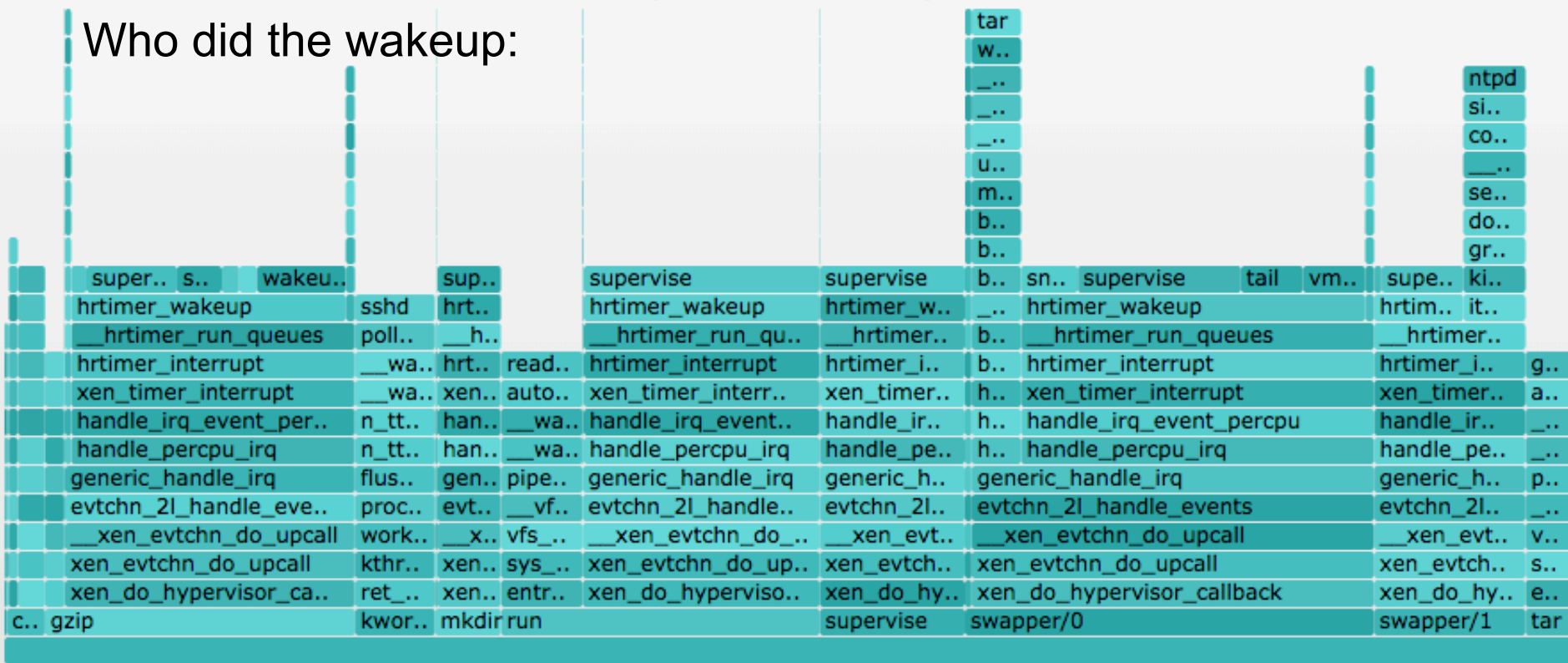
Trace blocking events with kernel stacks & time blocked (eg, using Linux BPF)

← Off-CPU time →

Stack depth ↑

# Wakeup Time Flame Graph

Who did the wakeup:



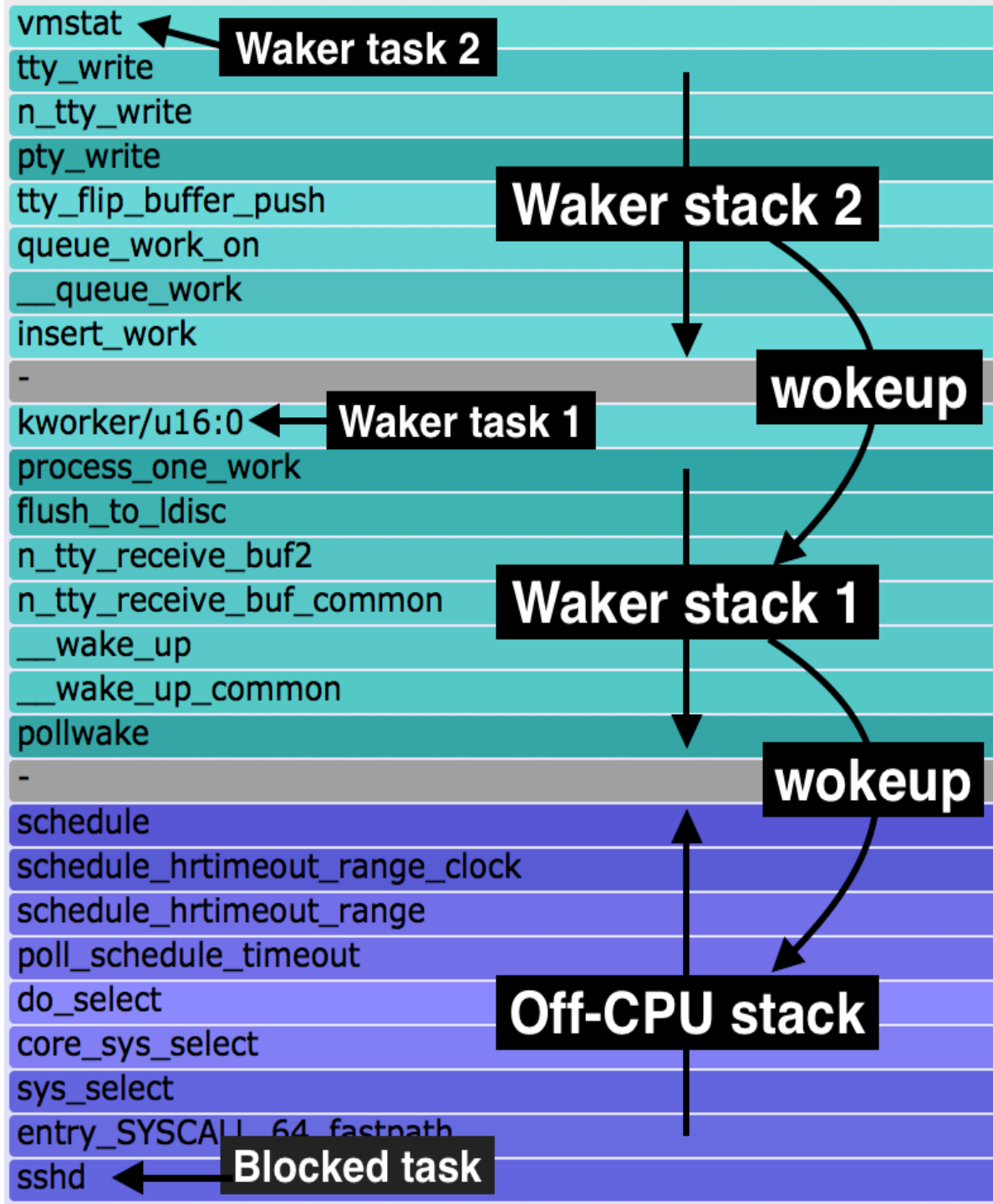
... can also associate wake-up stacks *with* off-CPU stacks (eg, Linux 4.6: `samples/bpf/offwaketime*`)

# Chain Graphs

Associate more than one waker: the full chain of wakeups

With enough stacks, all paths lead to metal

An approach for analyzing *all* off-CPU issues

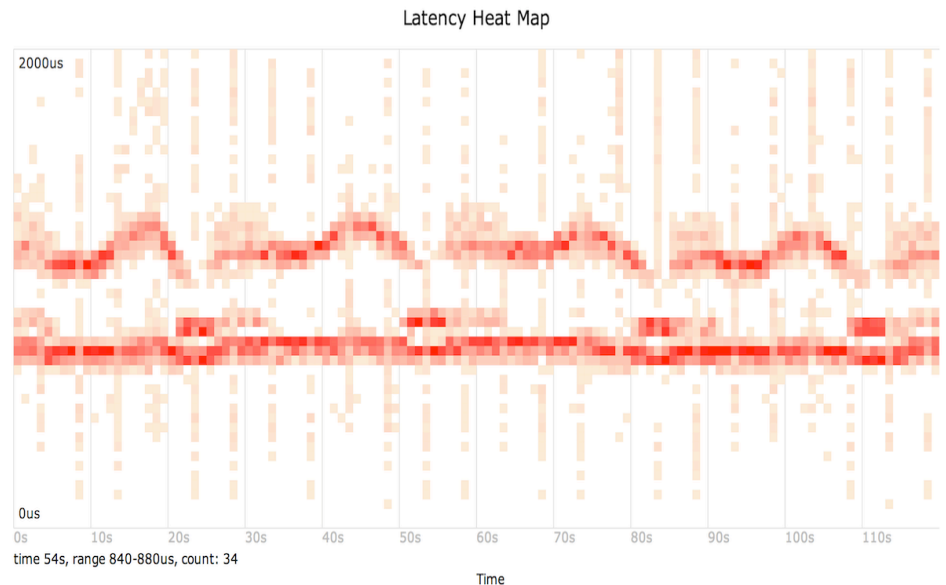
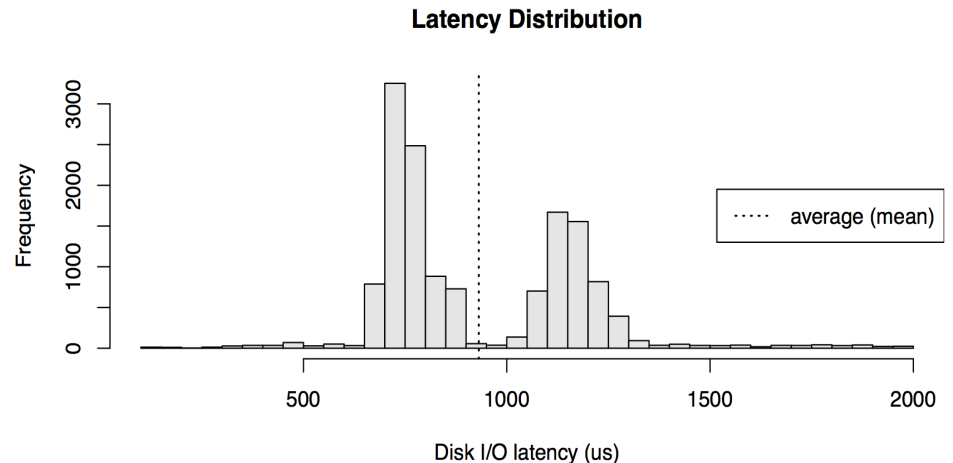


# Latency Correlations

1. Measure latency histograms at different stack layers
2. Compare histograms to find latency origin

Even better, use latency heat maps

- Match outliers based on both latency and time



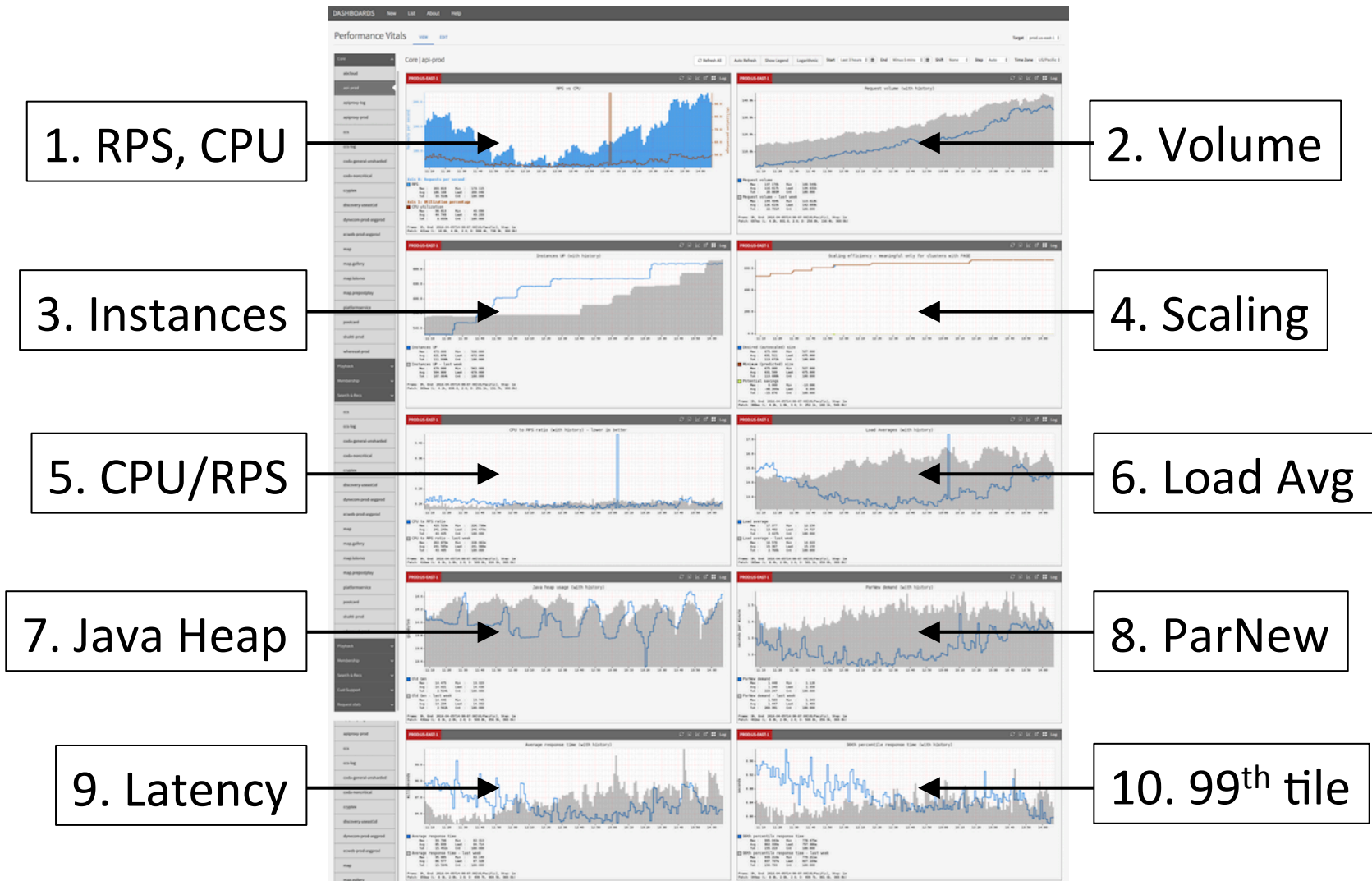


# Checklists: eg, Linux Perf Analysis in 60s

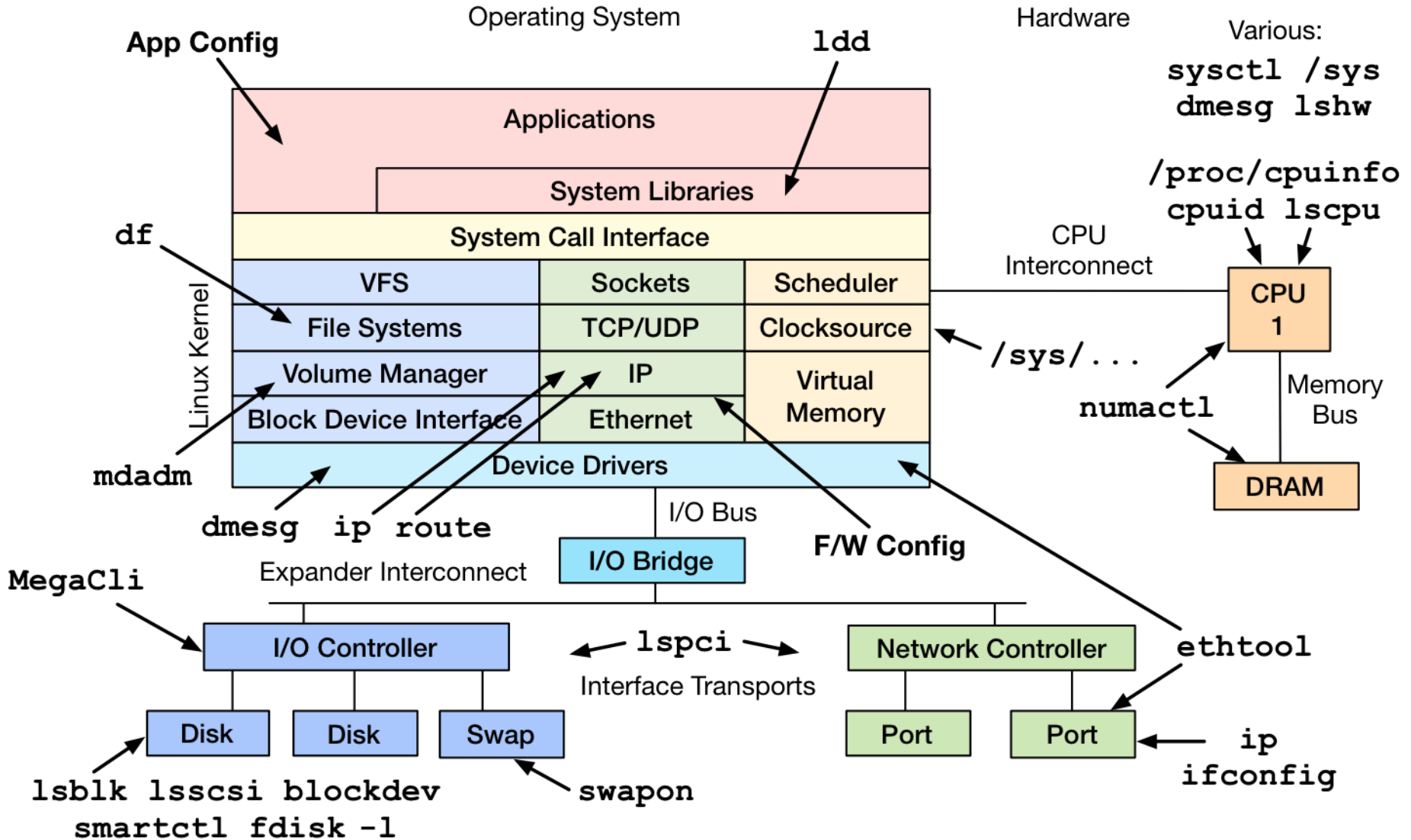
1. `uptime` -----▶ load averages
2. `dmesg | tail` -----▶ kernel errors
3. `vmstat 1` -----▶ overall stats by time
4. `mpstat -P ALL 1` -----▶ CPU balance
5. `pidstat 1` -----▶ process usage
6. `iostat -xz 1` -----▶ disk I/O
7. `free -m` -----▶ memory usage
8. `sar -n DEV 1` -----▶ network I/O
9. `sar -n TCP,ETCP 1` -----▶ TCP stats
10. `top` -----▶ check overview

<http://techblog.netflix.com/2015/11/linux-performance-analysis-in-60s.html>

# Checklists: eg, Netflix perfvitals Dashboard

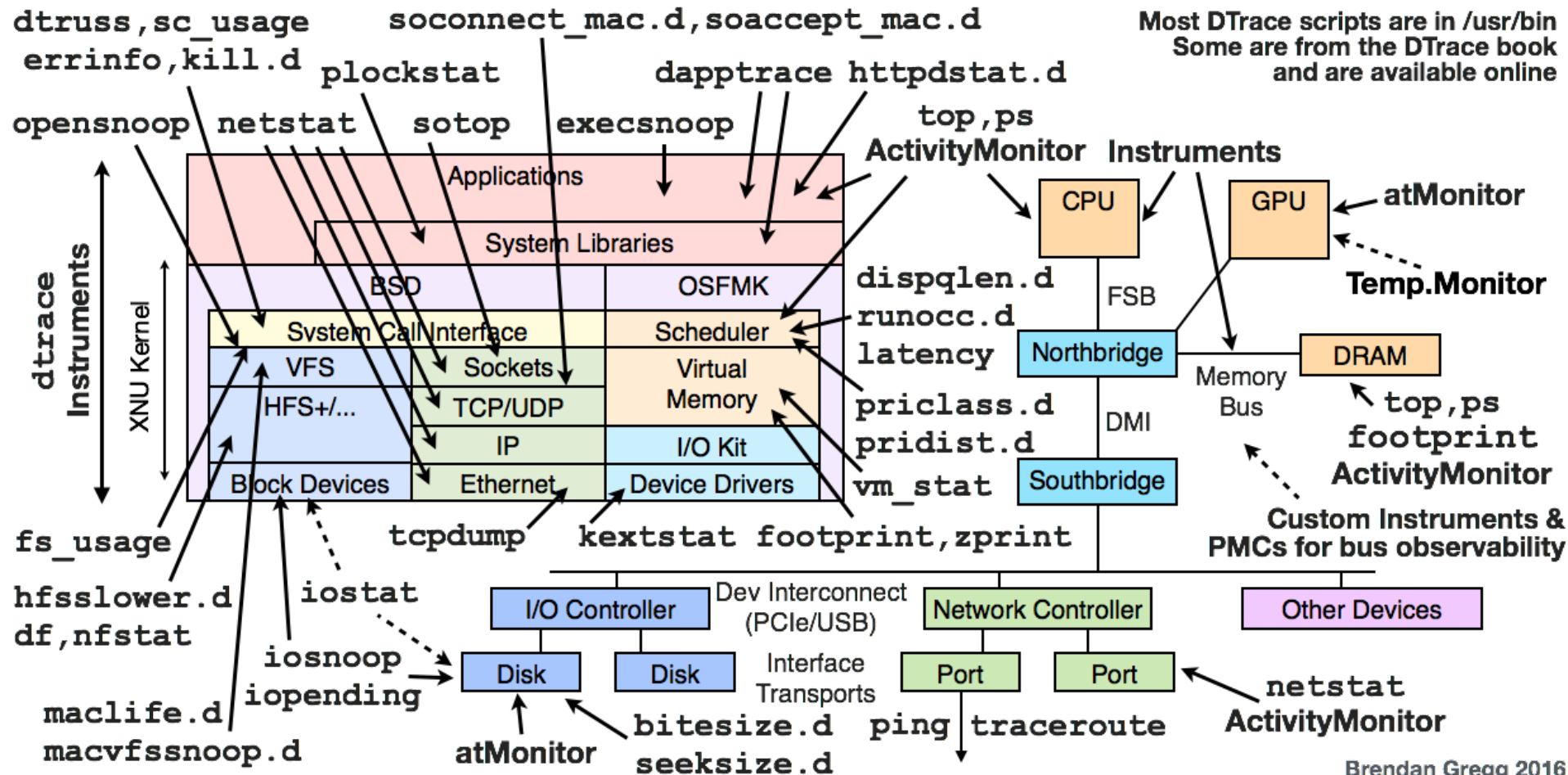


# Static Performance Tuning: eg, Linux



# Tools-Based Method

1. Try all the tools! May be an anti-pattern. Eg, OS X:



# Other Methodologies

- Scientific method
- 5 Why's
- Process of elimination
- Intel's Top-Down Methodology
- Method R

**What You Can Do**

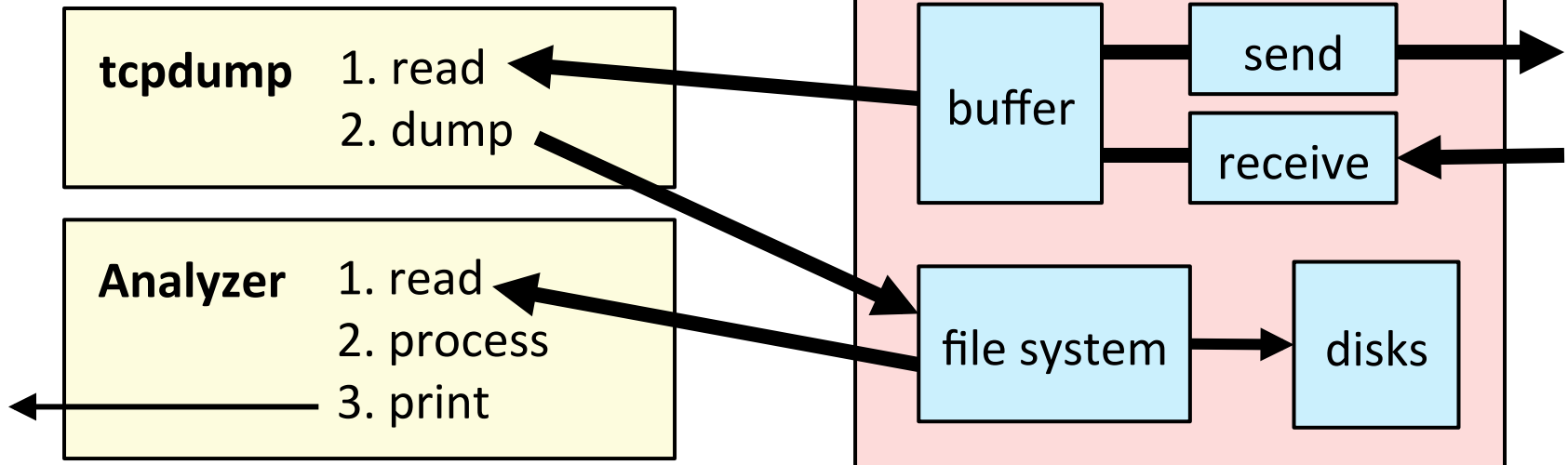
# What you can do

1. Know what's now possible on modern systems
  - Dynamic tracing: efficiently instrument any software
  - CPU facilities: PMCs, MSRs (model specific registers)
  - Visualizations: flame graphs, latency heat maps, ...
2. Ask questions first: use methodologies to ask them
3. Then find/build the metrics
4. Build or buy dashboards to support methodologies

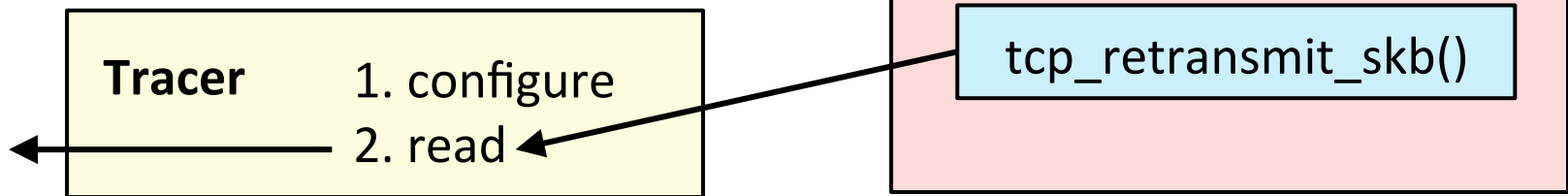
# Dynamic Tracing: Efficient Metrics

Eg, tracing TCP retransmits

**Old way:** packet capture

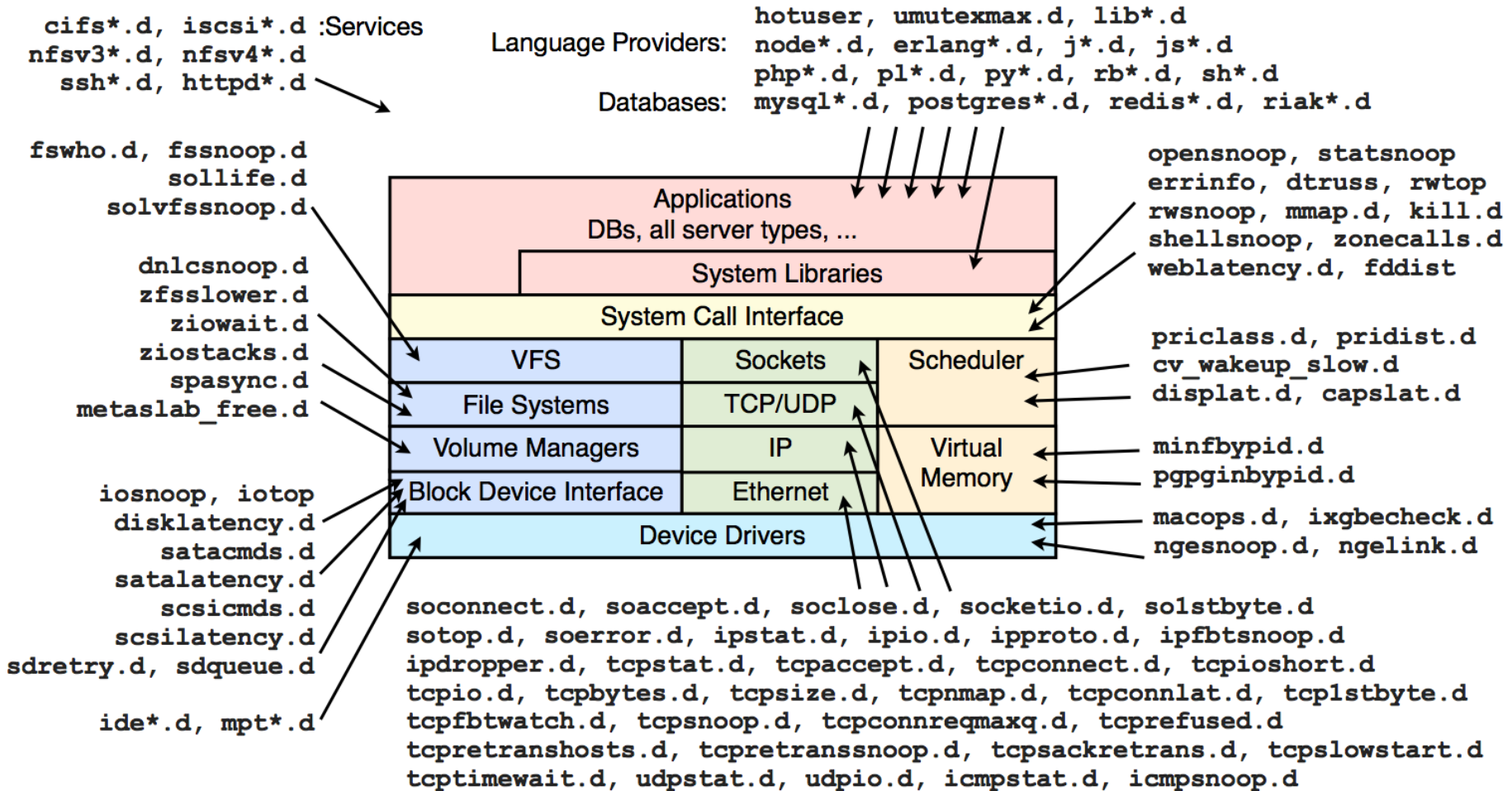


**New way:** dynamic tracing





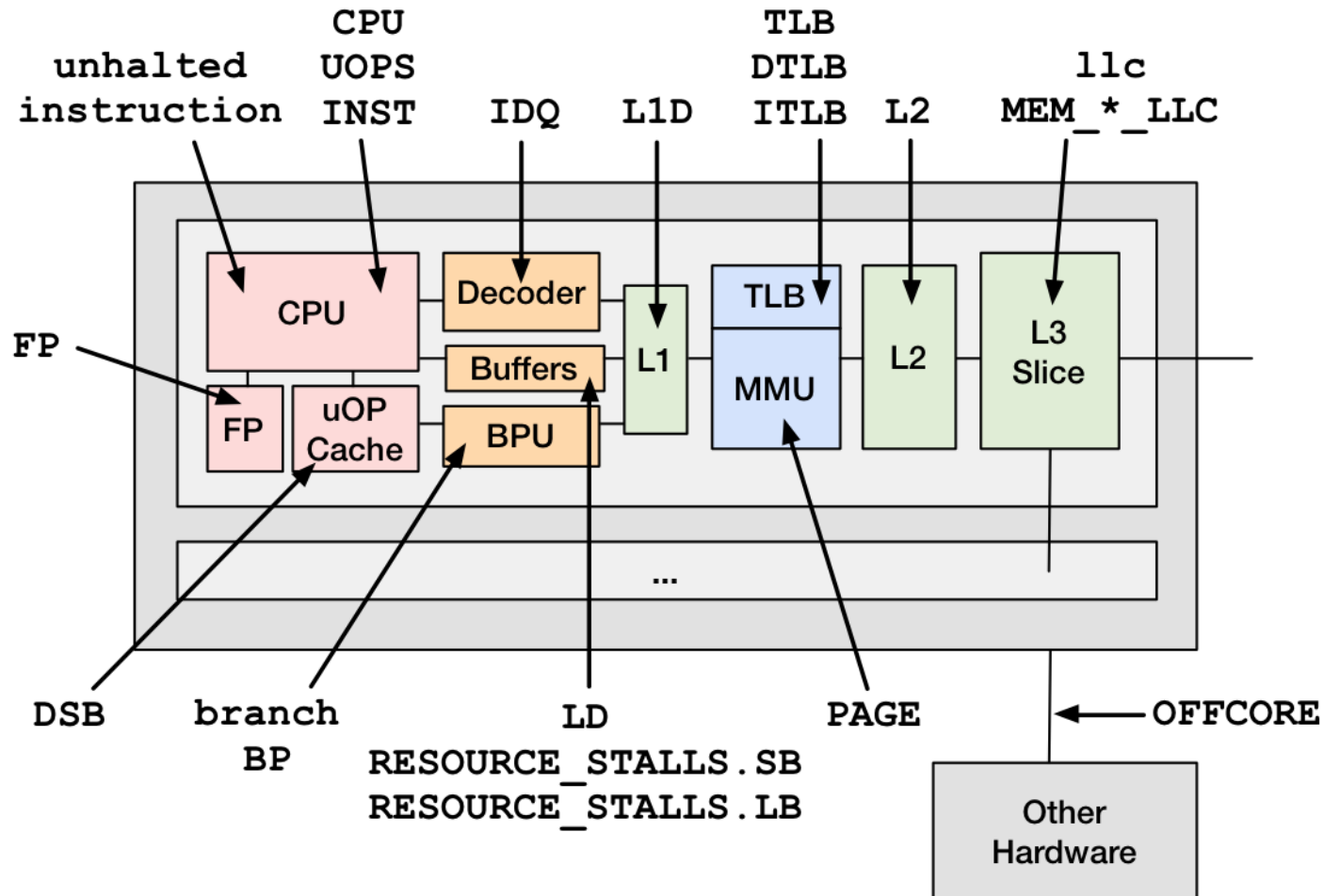
# Dynamic Tracing: Measure Anything



Those are Solaris/DTrace tools. Now becoming possible on all OSes:  
 FreeBSD & OS X DTrace, Linux BPF, Windows ETW

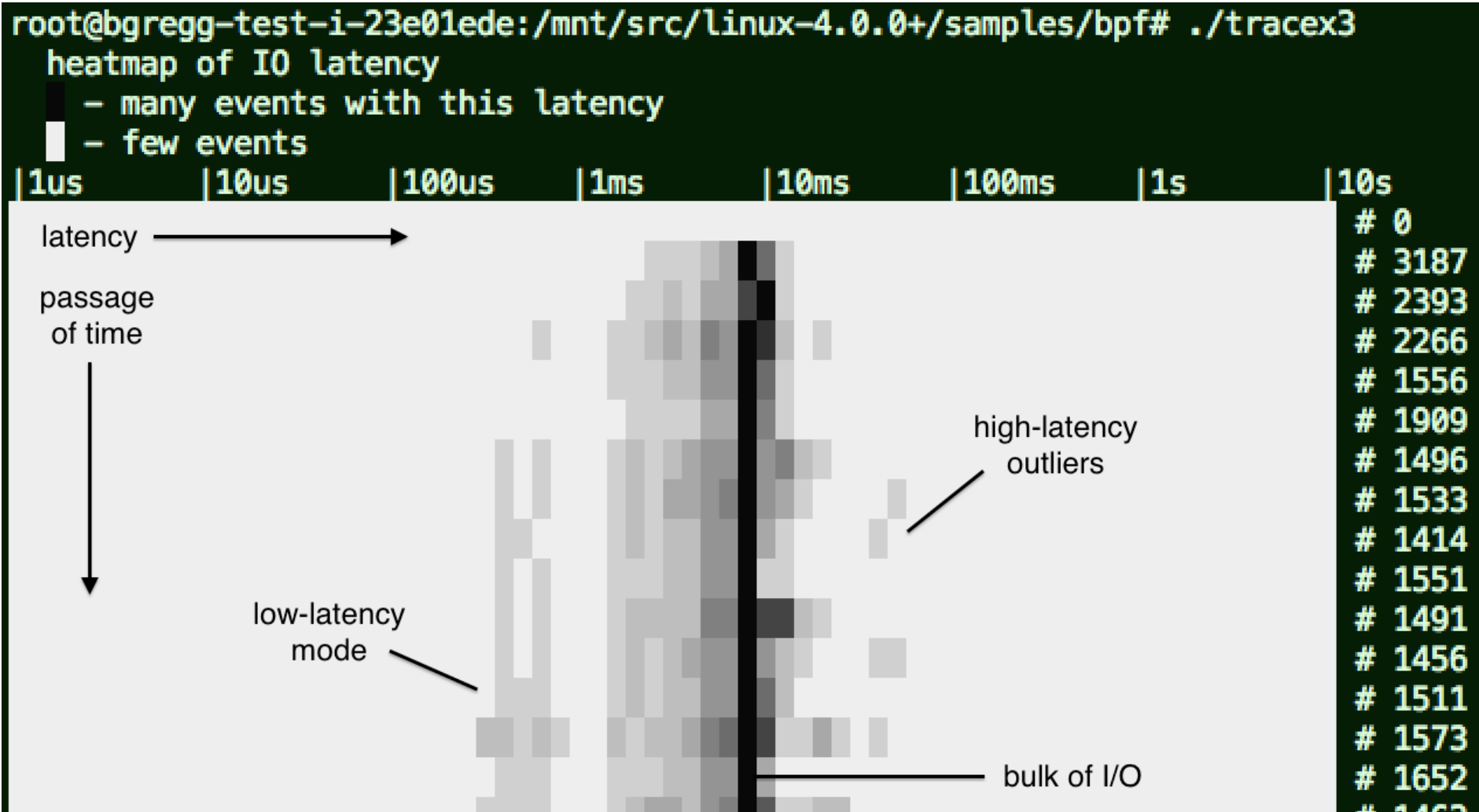
# Performance Monitoring Counters

Eg, FreeBSD PMC groups for Intel Sandy Bridge:

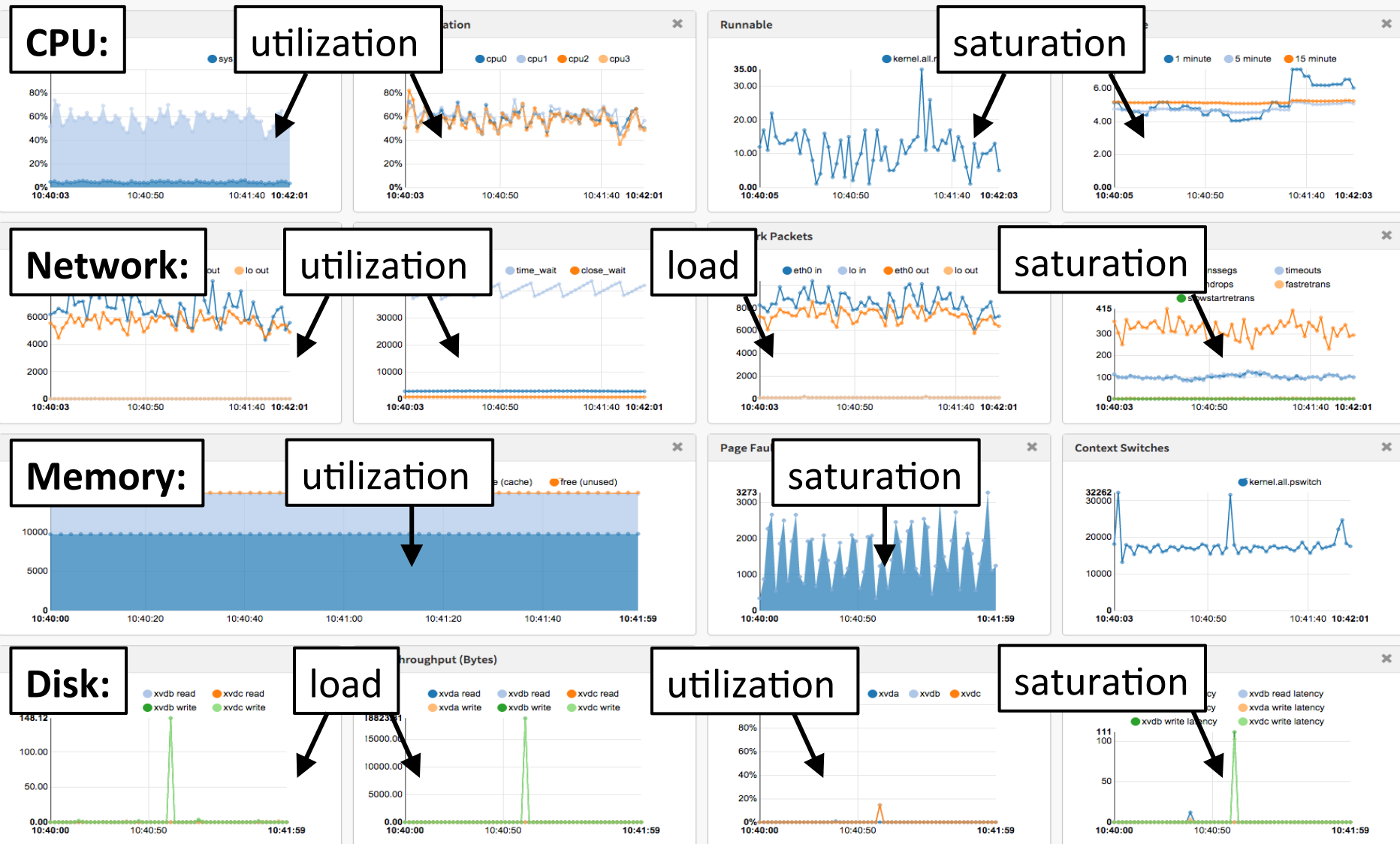


# Visualizations

Eg, Disk I/O latency as a heat map, quantized in kernel:



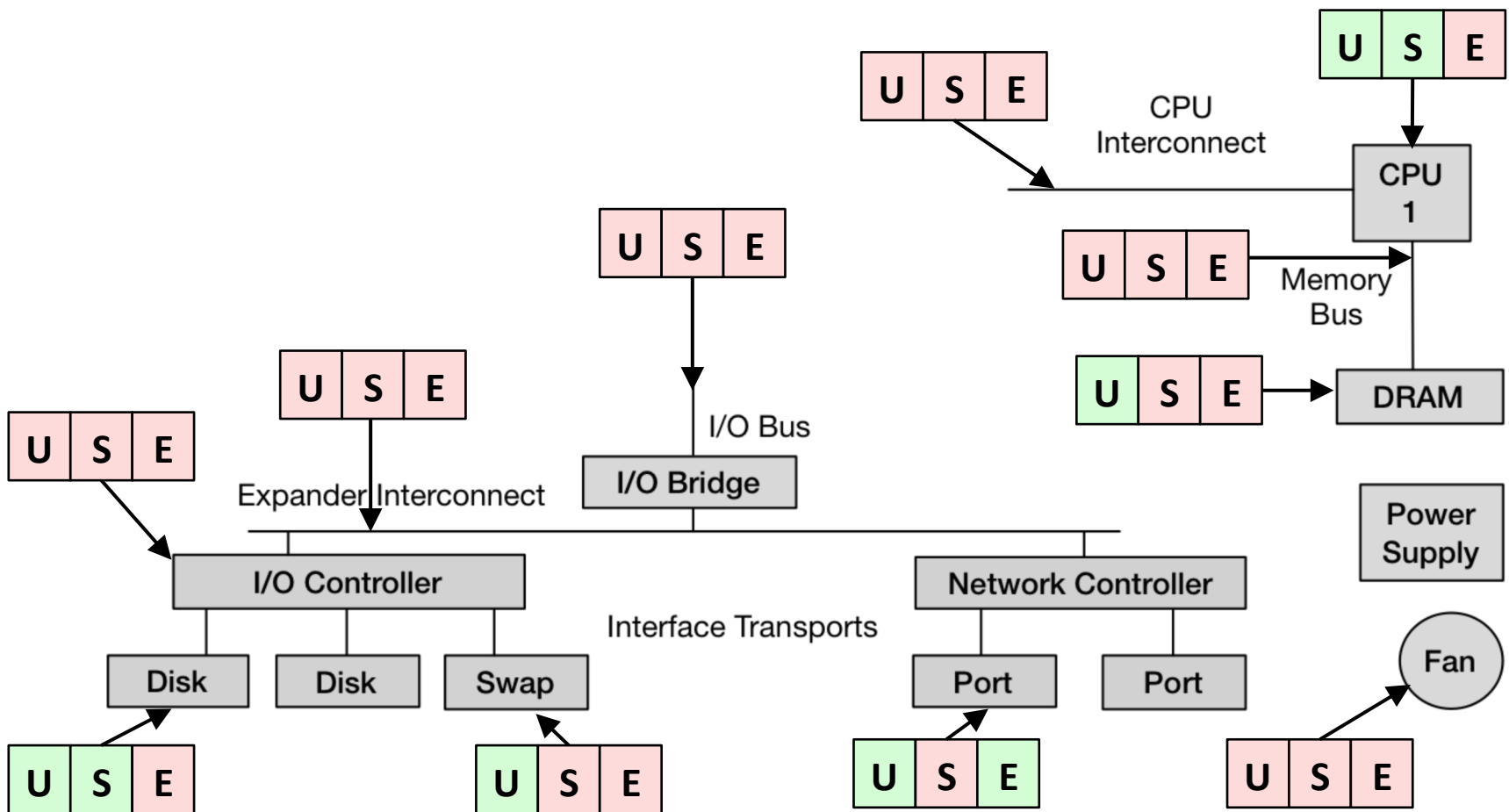
# USE Method: eg, Netflix Vector



# USE Method: To Do

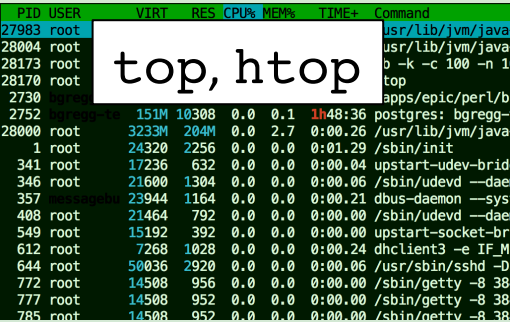
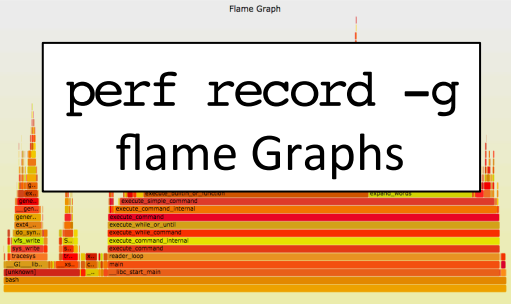
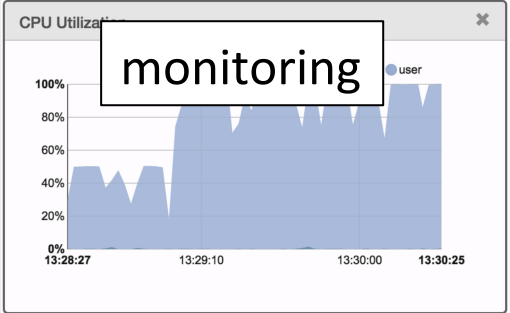
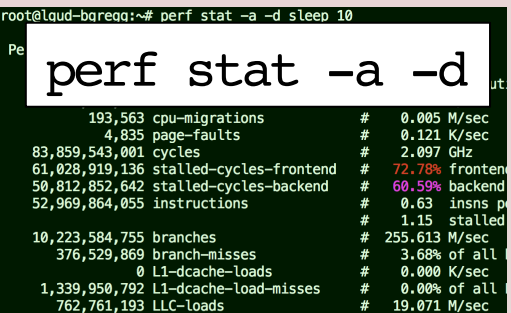
Hardware

Showing what **is** and **is not** commonly measured



# CPU Workload Characterization: To Do

Showing what **is** and **is not** commonly measured

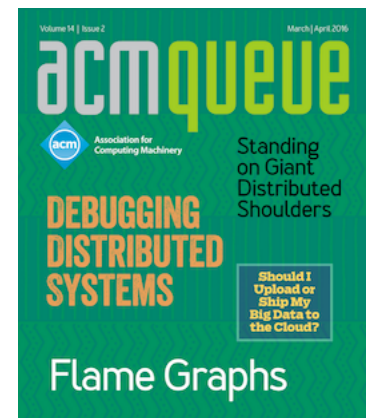
Who	Why
 <pre>PID USER VIRT RES CPU% MEM% TIME+ Command 27983 root 0 0 0.0 0.0 0:00.00 /usr/lib/jvm/java 28004 root 0 0 0.0 0.0 0:00.00 /usr/lib/jvm/java 28173 root 0 0 0.0 0.0 0:00.00 b -k -c 100 -n 1 28170 root 0 0 0.0 0.0 0:00.00 top 2730 0 0 0.0 0.0 0:00.00 apps/epic/perl/b 2752 151M 10308 0.0 0.1 1h48:36 postgres: bgregg- 28000 root 3233M 204M 0.0 2.7 0:00.26 /usr/lib/jvm/java 1 root 24320 2256 0.0 0.0 0:01.29 /sbin/init 341 root 17236 632 0.0 0.0 0:00.04 upstart-udev-brid 346 root 21600 1304 0.0 0.0 0:00.06 /sbin/udevmd --dae 357 23944 1164 0.0 0.0 0:00.21 dbus-daemon --sys 408 root 21464 792 0.0 0.0 0:00.00 /sbin/udevmd --dae 549 root 15192 392 0.0 0.0 0:00.00 upstart-socket-br 612 root 7268 1028 0.0 0.0 0:00.24 dhclient3 -e IF M 644 root 50036 2920 0.0 0.0 0:00.06 /usr/sbin/sshd -D 772 root 14508 956 0.0 0.0 0:00.00 /sbin/getty -8 38 777 root 14508 952 0.0 0.0 0:00.00 /sbin/getty -8 38 785 root 14508 952 0.0 0.0 0:00.00 /sbin/getty -8 38</pre>	 <pre>Flame Graph perf record -g flame Graphs</pre>
How	What
 <pre>CPU Utilization monitoring</pre>	 <pre>root@laud-bqregq:~# perf stat -a -d sleep 10 Perf perf stat -a -d 193,563 cpu-migrations # 0.005 M/sec 4,835 page-faults # 0.121 K/sec 83,859,543,001 cycles # 2.097 GHz 61,028,919,136 stalled-cycles-frontend # 72.78% frontend 50,812,852,642 stalled-cycles-backend # 60.59% backend 52,969,864,055 instructions # 0.63 insns p 10,223,584,755 branches # 255.613 M/sec 376,529,869 branch-misses # 3.68% of all 0 L1-dcache-loads # 0.000 K/sec 1,339,950,792 L1-dcache-load-misses # 0.00% of all 762,761,193 LLC-loads # 19.071 M/sec</pre>

# Summary

- It is the crystal ball age of performance observability
- What matters is the questions you want answered
- Methodologies are a great way to pose questions

# References & Resources

- USE Method
  - <http://queue.acm.org/detail.cfm?id=2413037>
  - <http://www.brendangregg.com/usemethod.html>
- TSA Method
  - <http://www.brendangregg.com/tsamethod.html>
- Off-CPU Analysis
  - <http://www.brendangregg.com/offcpuanalysis.html>
  - <http://www.brendangregg.com/blog/2016-01-20/ebpf-offcpu-flame-graph.html>
  - <http://www.brendangregg.com/blog/2016-02-05/ebpf-chaingraph-prototype.html>
- Static Performance Tuning, Richard Elling, Sun blueprint, May 2000
- RED Method: <http://www.slideshare.net/weaveworks/monitoring-microservices>
- Other system methodologies
  - Systems Performance: Enterprise and the Cloud, Prentice Hall 2013
  - <http://www.brendangregg.com/methodology.html>
  - The Art of Computer Systems Performance Analysis, Jain, R., 1991
- Flame Graphs
  - <http://queue.acm.org/detail.cfm?id=2927301>
  - <http://www.brendangregg.com/flamegraphs.html>
  - <http://techblog.netflix.com/2015/07/java-in-flames.html>
- Latency Heat Maps
  - <http://queue.acm.org/detail.cfm?id=1809426>
  - <http://www.brendangregg.com/HeatMaps/latency.html>
- ARPA Network: <http://www.computerhistory.org/internethistory/1960s>
- RSTS/E System User's Guide, 1985, page 4-5
- DTrace: Dynamic Tracing in Oracle Solaris, Mac OS X, and FreeBSD, Prentice Hall 2011
- Apollo: <http://www.hq.nasa.gov/office/pao/History/alsj/a11> <http://www.hq.nasa.gov/alsj/alsj-LMdocs.html>



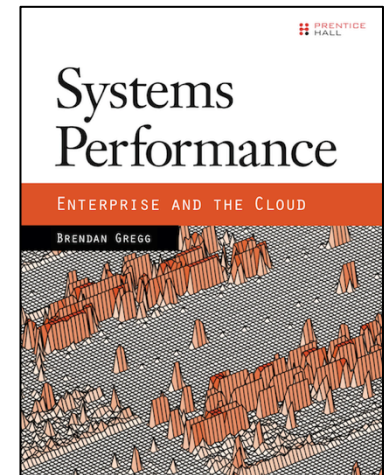




# ACM Applicative 2016

Jun, 2016

- Questions?
- <http://slideshare.net/brendangregg>
- <http://www.brendangregg.com>
- [bgregg@netflix.com](mailto:bgregg@netflix.com)
- @brendangregg



**NETFLIX**