



Camp-Con-World-Fest-Summit

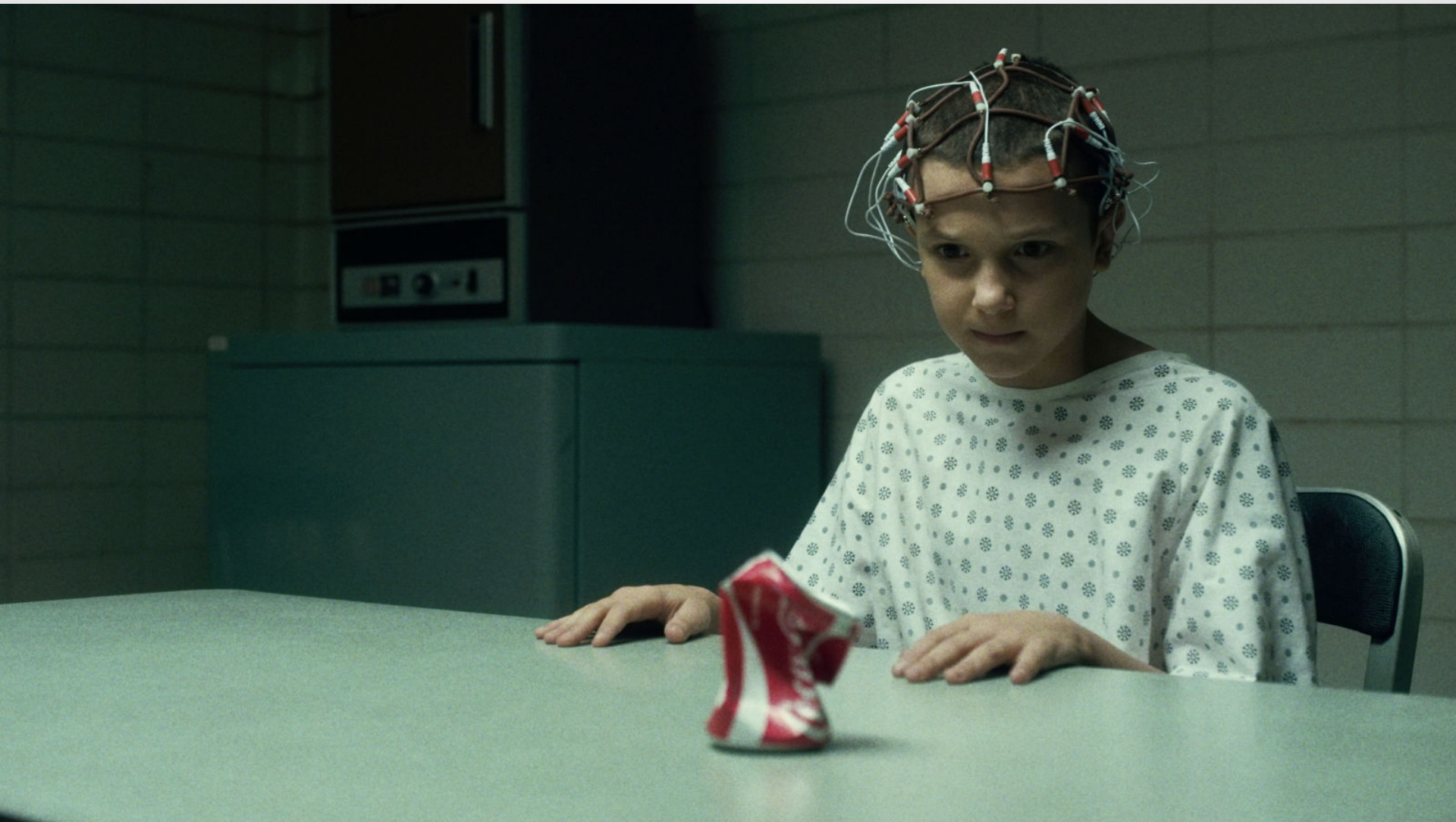
Designing Tracing Tools



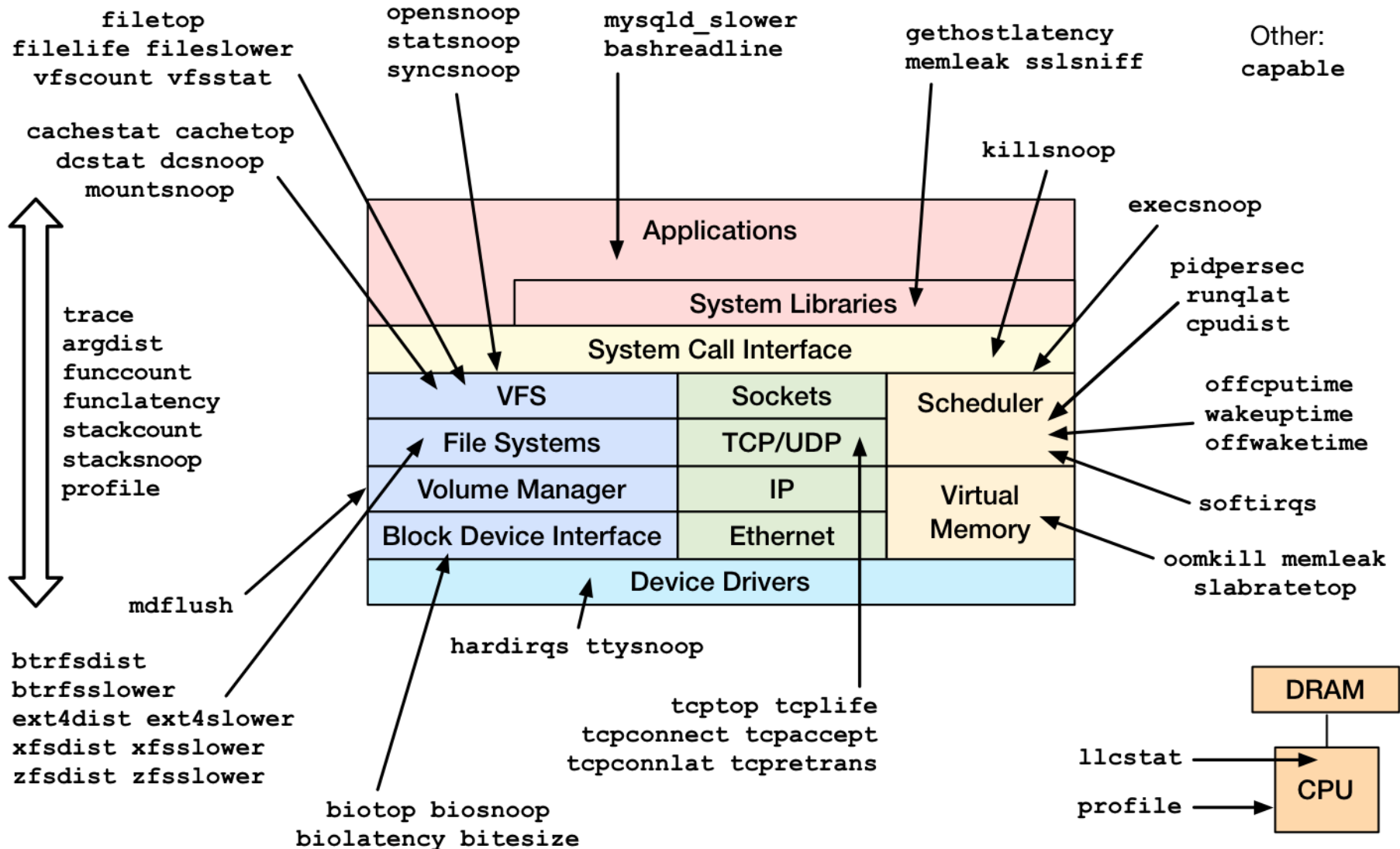
NETFLIX

Brendan Gregg, *Senior Performance Architect*

Wielding Superpowers



I'm currently developing more tracing tools (bcc/BPF)



Tool Design

- For tool developers
- For everyone else: what you can ask for
 - Tool templates
 - GUI visualizations
- The following is applicable to all tracers
 - sysdig, bcc/BPF, DTrace, SystemTap, LTTng, ...

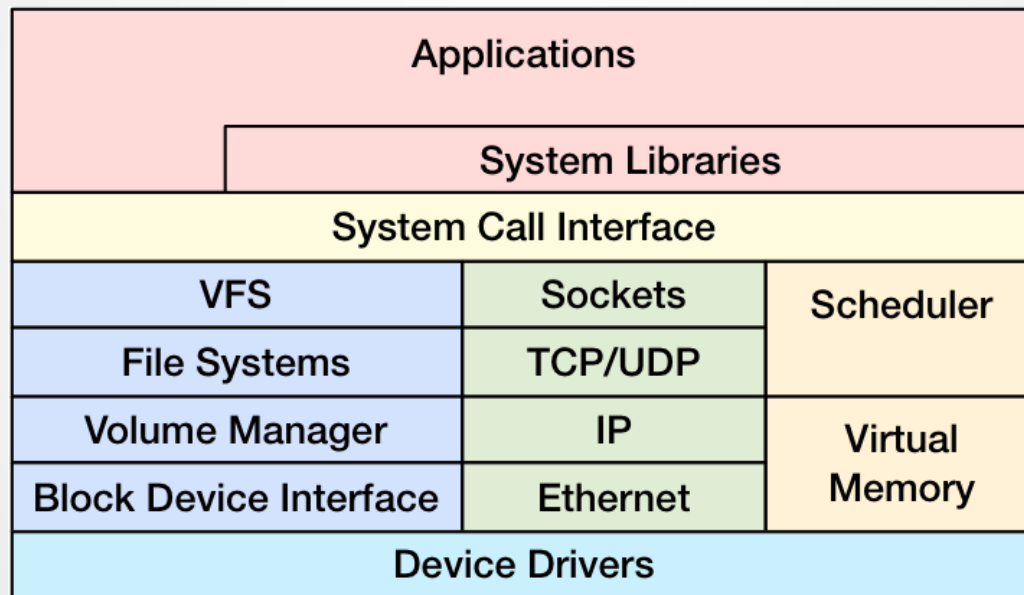
Methodologies



Methodology-driven Design

- Methodologies provide ideas for purposeful tools
- Find/draw a functional diagram, apply methods

Operating Systems



See: <http://www.brendangregg.com/methodology.html>

Methodology Examples

Eg, at the syscall layer (well known & documented):

- Workload Characterization
 - `exec()` or `open()` per-event trace (execsnoop, opensnoop)
 - `connect()/accept()` per-event trace (tcpconnect, tcpaccept)
 - `read()/write()` size histogram (one-liners)
- Latency Analysis
 - `read()/write()` latency histogram (funclatency, ...)
- USE Method
 - network utilization by thread (not done yet)
 - syscall errors (fserrors, soerrors)

CLI Tracing Tools



CLI Templates

1. Per event output

- `*snoop, *slower 0, ...`

2. Filtered event output

- `*slower`

3. Interval summary

- `*stat, *top`

4. Count summary

- `*count`

5. Histogram summary

- `*dist, *latency`

6. Heatmap summary

- `spectrogram.lua, subsecoffset.lua, ...`

Template 1: Per Event Output

Examples: *snoop, *slower 0, ...

```
# opensnoop
PID      COMM          FD  ERR  PATH
10085    sshd          3   0    /lib/x86_64-linux-gnu/libkeyutils.so.1
10085    sshd          3   0    /lib/x86_64-linux-gnu/libresolv.so.2
10085    sshd          3   0    /lib/x86_64-linux-gnu/libgpg-error.so.0
10085    sshd          3   0    /dev/urandom
10085    sshd         -1   2    /lib/x86_64-linux-gnu/.libcrypto.so.1.0.0.hmac
10085    sshd         -1   2    /proc/sys/crypto/fips_enabled
10085    sshd          3   0    /proc/filesystems
10085    sshd          3   0    /dev/null
10085    sshd          3   0    /proc/10085/fd
10085    sshd          3   0    /usr/lib/ssl/openssl.cnf
10085    sshd          3   0    /etc/gai.conf
10085    sshd          3   0    /etc/nsswitch.conf
10085    sshd          3   0    /etc/ld.so.cache
10085    sshd          3   0    /lib/x86_64-linux-gnu/libnss_compat.so.2
10085    sshd          3   0    /etc/ld.so.cache
10085    sshd          3   0    /lib/x86_64-linux-gnu/libnss_nis.so.2
[...]
```

Template 2: Filtered Event Output

Examples: *slower

```
# sysdig -c fileslower 1
TIME                PROCESS           TYPE           LAT(ms) FILE
2014-04-13 20:40:43.973 cksum             read            2 /mnt/partial.0.0
2014-04-13 20:40:44.187 cksum             read            1 /mnt/partial.0.0
2014-04-13 20:40:44.689 cksum             read            2 /mnt/partial.0.0
2014-04-13 20:40:45.005 cksum             read            2 /mnt/partial.0.0
2014-04-13 20:40:45.193 cksum             read            1 /mnt/partial.0.0
[...]
```

Tools like this can also do all event output:

```
# sysdig -c fileslower 0
TIME                PROCESS           TYPE           LAT(ms) FILE
2014-04-13 20:59:04.414 ls                 read            0 /lib/x86_64-linux-gnu/librt.so.1
2014-04-13 20:59:04.414 ls                 read            0 /lib/x86_64-linux-gnu/libacl.so.1
2014-04-13 20:59:04.414 ls                 read            0 /lib/x86_64-linux-gnu/libc.so.6
2014-04-13 20:59:04.414 ls                 read            0 /lib/x86_64-linux-gnu/libdl.so.2
2014-04-13 20:59:04.414 ls                 read            0 /lib/x86_64-linux-gnu/libattr.so.1
2014-04-13 20:59:04.415 ls                 read            0 /proc/filesystems
2014-04-13 20:59:04.415 ls                 read            0 /proc/filesystems
[...]
```

Template 3: Interval Summary

Examples: *stat, *top

```
# dcstat
TIME          REFS/s    SLOW/s    MISS/s    HIT%
08:11:47:    2059      141       97        95.29
08:11:48:    79974     151       106       99.87
08:11:49:    192874    146       102       99.95
08:11:50:    2051      144       100       95.12
08:11:51:    73373     17239     17194     76.57
08:11:52:    54685     25431     25387     53.58
08:11:53:    18127     8182      8137     55.12
08:11:54:    22517     10345     10301     54.25
08:11:55:    7524      2881      2836     62.31
08:11:56:    2067      141       97        95.31
08:11:57:    2115      145       101       95.22
[...]
```

Template 4: Count Summary

Examples: *count

```
# funccount 'vfs_*'
Tracing... Ctrl-C to end.
^C
ADDR                FUNC                COUNT
ffffffff811efe81    vfs_create          1
ffffffff811f24a1    vfs_rename          1
ffffffff81215191    vfs_fsync_range    2
ffffffff81231df1    vfs_lock_file      30
ffffffff811e8dd1    vfs_fstata         152
ffffffff811e8d71    vfs_fstat          154
ffffffff811e4381    vfs_write          166
ffffffff811e8c71    vfs_getattr_nosec  262
ffffffff811e8d41    vfs_getattr        262
ffffffff811e3221    vfs_open           264
ffffffff811e4251    vfs_read           470
Detaching...
```

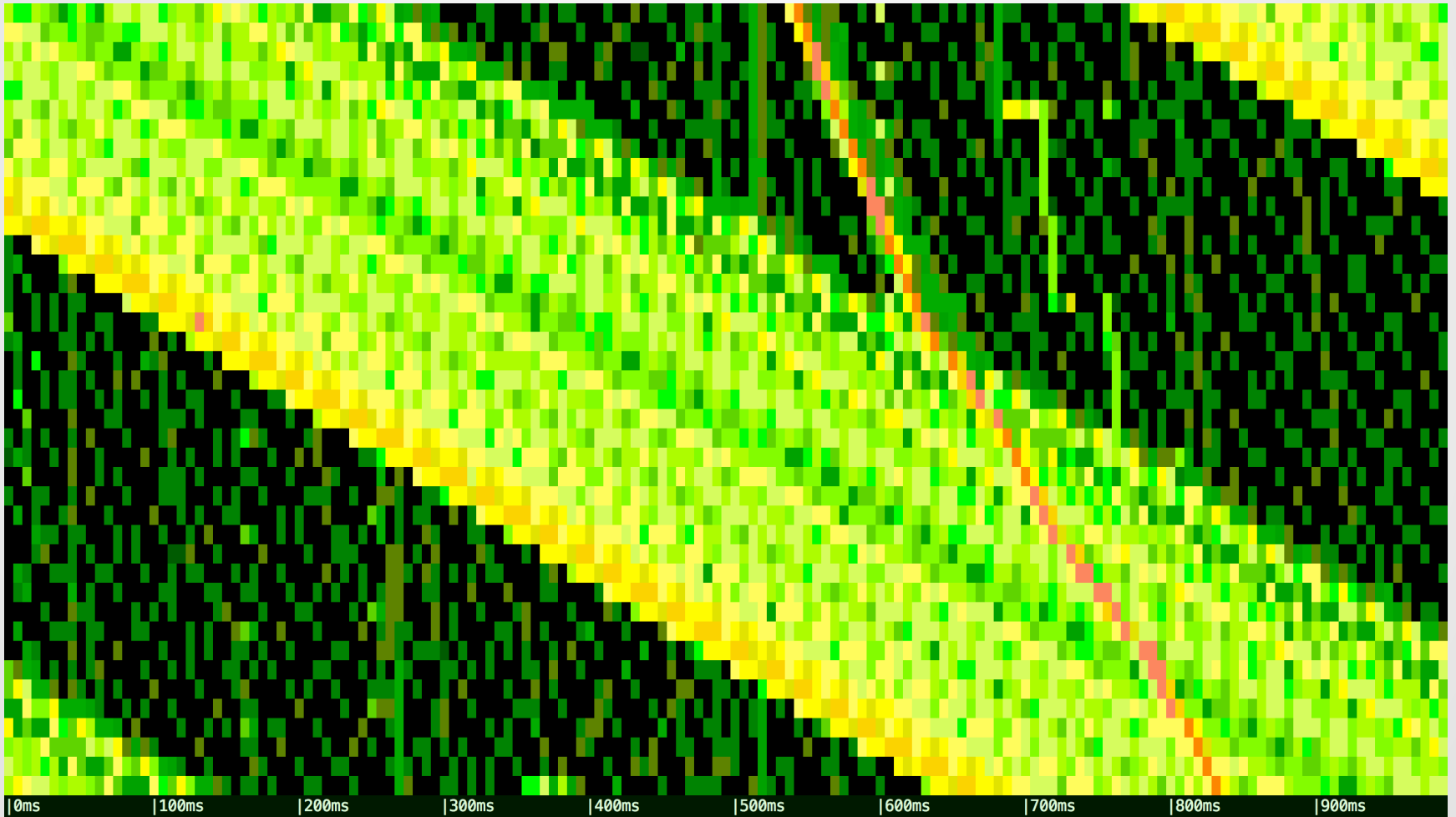
Template 5: Histogram Summary

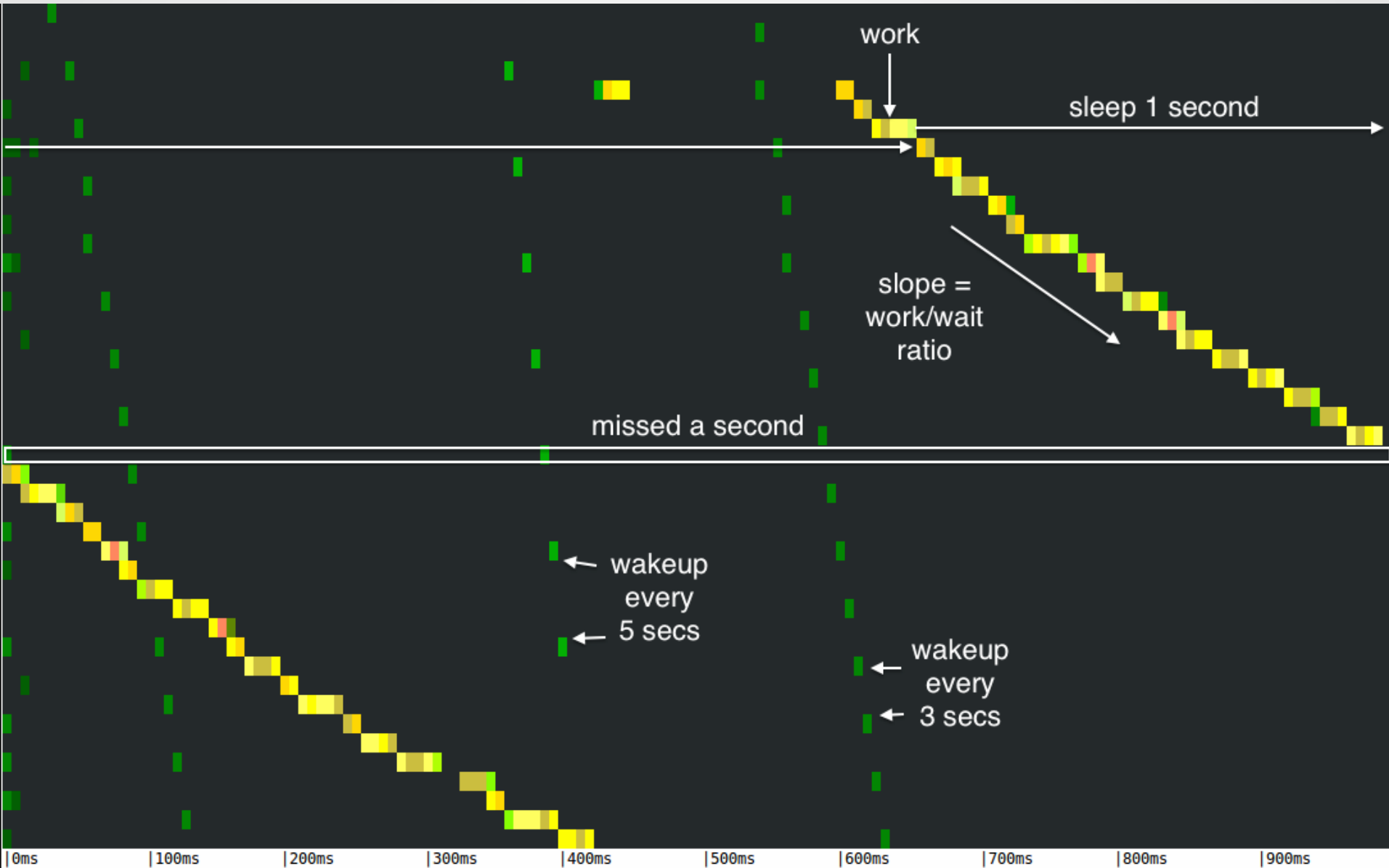
Examples: `*dist`, `*latency`

```
# biolatency
Tracing block device I/O... Hit Ctrl-C to end.
^C
      usecs           : count      distribution
      4 -> 7          : 0
      8 -> 15         : 0
     16 -> 31         : 0
     32 -> 63         : 0
     64 -> 127        : 1
    128 -> 255        : 12      *****
    256 -> 511        : 15      *****
    512 -> 1023       : 43      *****
   1024 -> 2047       : 52      *****
   2048 -> 4095       : 47      *****
   4096 -> 8191       : 52      *****
   8192 -> 16383      : 36      *****
  16384 -> 32767      : 15      *****
  32768 -> 65535      : 2        *
 65536 -> 131071     : 2        *
```

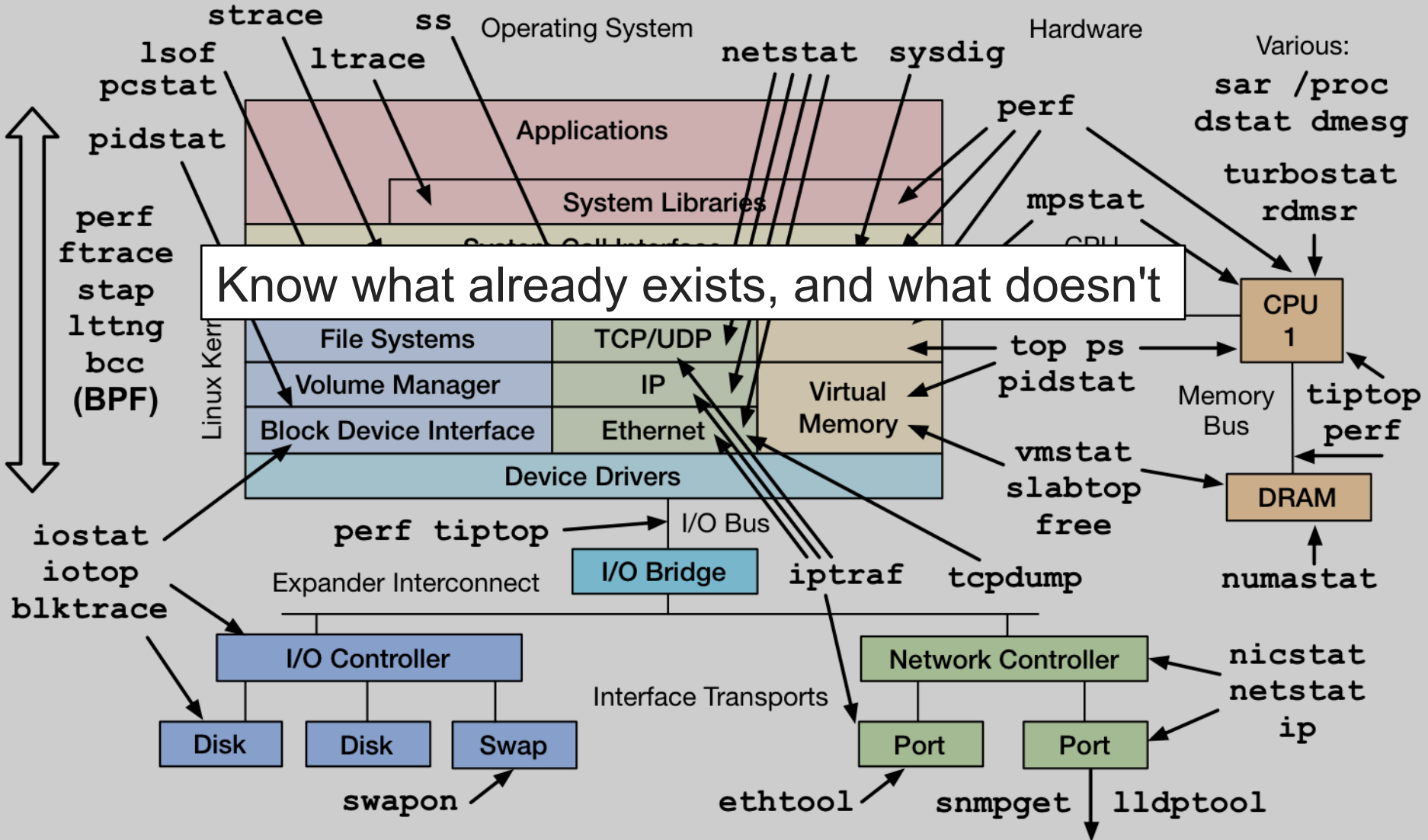
Template 6: Heatmap Summary

Example: `subsecoffset.lua` (aka "spectrogram")

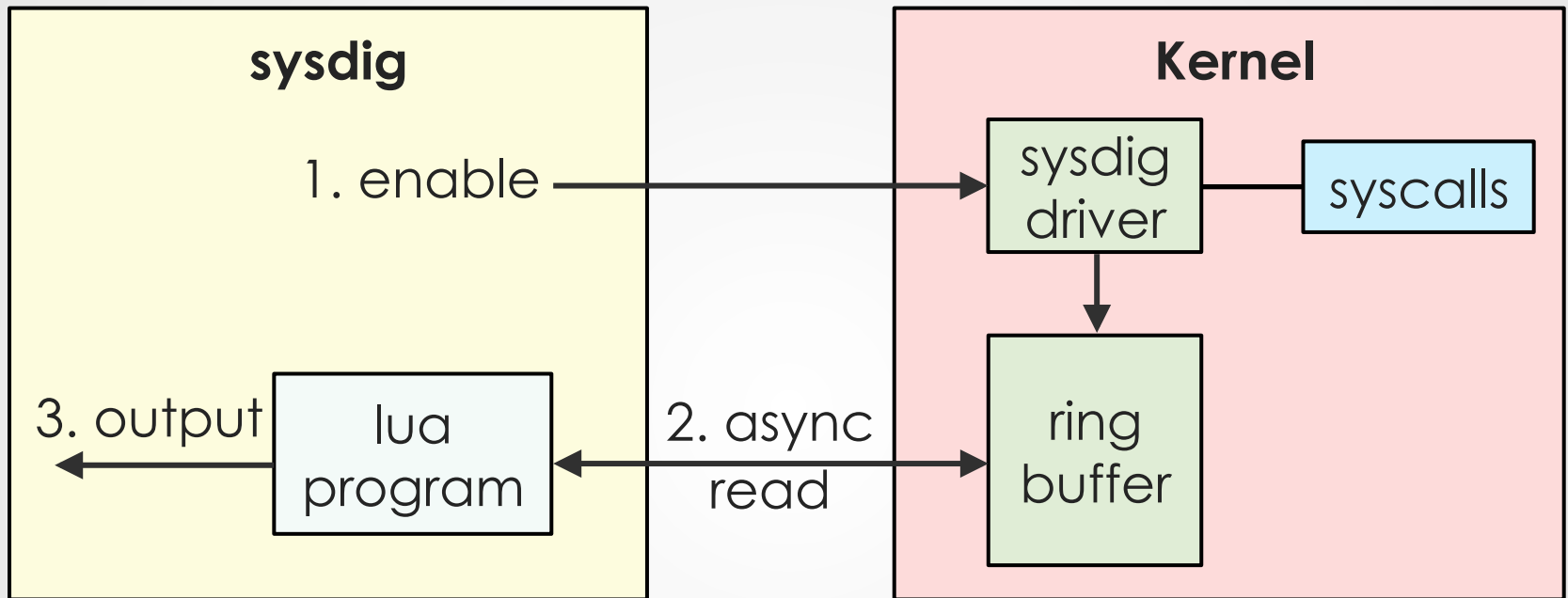




Valuable



Low Overhead (or documented)



- Understand tracing internals
 - For example, sysdig's design has ~20x lower overhead than strace (it still has overhead: test and measure to see if it's acceptable)
 - Tracing overhead is usually relative to event rate
- Design for low overhead, and document expectations

Documentation

- Good tools have 3 docs:

1. Code comments
2. Man page
3. Examples file

- Man page

– troff, docbook, ...

- Examples file:

common man macros (see `groff_man(7)`)

.TH	Title heading
.SH	Section heading
.IP	Indented paragraph
.TP	Indented paragraph with label
.B	Bold
\-	-

Demonstrations of biosnoop, the Linux eBPF/bcc version.

biosnoop traces block device I/O (disk I/O), and prints a line of output per I/O. Example:

```
# ./biosnoop
TIME(s)          COMM          PID    DISK    T  SECTOR    BYTES    LAT(ms)
0.000004001    supervise    1950    xvda1   W  13092560  4096      0.74
[...]
```

Concise, intuitive, self-explanatory

```
# ./iolatency
```

```
Tracing block I/O. Output every 1 seconds. Ctrl-C to end.
```

<code>>=(ms)</code>	<code>..</code>	<code><(ms)</code>	<code>:</code>	<code>I/O</code>	<code>Distribution</code>
0	->	1	:	4381	#####
1	->	2	:	9	#
2	->	4	:	5	#
4	->	8	:	0	
8	->	16	:	1	#

```
[...]
```

- Useful startup message
 - What I'm tracing, when there's output, when I'll end
- Vigorous tooling is concise
 - No wasted text; leave less useful output for non-default options
- Unix philosophy: do one thing and do it well

POSIX-style Arguments

```
# ./biolateness -h
usage: biolateness [-h] [-T] [-Q] [-m] [-D] [interval] [count]
```

Summarize block device I/O latency as a histogram

positional arguments:

interval	output interval, in seconds
count	number of outputs

optional arguments:

-h, --help	show this help message and exit
-T, --timestamp	include timestamp on output
-Q, --queued	include OS queued time in I/O time
-m, --milliseconds	millisecond histogram
-D, --disks	print a histogram per disk device

examples:

```
./biolateness # summarize block I/O latency as a histogram
./biolateness 1 10 # print 1 second summaries, 10 times
./biolateness -mT 1 # 1s summaries, milliseconds, and timestamps
./biolateness -Q # include OS queued time in I/O time
./biolateness -D # show each disk device separately
```

POSIX-style Arguments

Option	Alternate	Expectation
-a	--all	all events
-c CMD	--cmd ...	run this command
-d SECONDS	--duration ...	duration of tool execution
-h	--help	help
-i FILE	--input ...	input file
-i SECONDS	--interval ...	summary interval
-n name	--name ...	this process name only
-o FILE	--output ...	output file
-p PID	--pid ...	this process ID only
-P	--by-process	per-process ID breakdown
-P PORT	--port ...	this TCP port only
-t or -T	--[no]timestamp	include or exclude timestamps
-v	--verbose	verbose output
-x	--extended, --errors	extended output, or only failures
[interval [count]]	-	summary interval, and # of outputs

Testing, Testing, Testing

- If you can't write the workload, you can't write the tool
 - Be it 10 lines of C, some shell, or dd
 - `dd if=/dev/urandom of=/dev/null bs=1k count=23`
- Known workload testing: create 23 events
- Testing can be time consuming
 - I spend more time testing a tool than writing it
 - Automatic tool testing is a difficult problem

Example: gethostlatency

```
# gethostlatency
```

TIME	PID	COMM	LATms	HOST
06:10:24	28011	wget	90.00	www.iovisor.org
06:10:28	28127	wget	0.00	www.iovisor.org
06:10:41	28404	wget	9.00	www.netflix.com
06:10:48	28544	curl	35.00	www.netflix.com.au
06:11:10	29054	curl	31.00	www.plumgrid.com
06:11:16	29195	curl	3.00	www.facebook.com
06:11:25	29404	curl	72.00	foo
06:11:28	29475	curl	1.00	foo

Example: ext4slower

```
# ext4slower 1
```

```
Tracing ext4 operations slower than 1 ms
```

TIME	COMM	PID	T	BYTES	OFF_KB	LAT(ms)	FILENAME
06:49:17	bash	3616	R	128	0	7.75	cksum
06:49:17	cksum	3616	R	39552	0	1.34	[
06:49:17	cksum	3616	R	96	0	5.36	2to3-2.7
06:49:17	cksum	3616	R	96	0	14.94	2to3-3.4
06:49:17	cksum	3616	R	10320	0	6.82	411toppm
06:49:17	cksum	3616	R	65536	0	4.01	a2p
06:49:17	cksum	3616	R	55400	0	8.77	ab
06:49:17	cksum	3616	R	36792	0	16.34	aclocal-1.14
06:49:17	cksum	3616	R	15008	0	19.31	acpi_listen
06:49:17	cksum	3616	R	6123	0	17.23	add-apt-repository
06:49:17	cksum	3616	R	6280	0	18.40	addpart
06:49:17	cksum	3616	R	27696	0	2.16	addr2line
06:49:17	cksum	3616	R	58080	0	10.11	ag
06:49:17	cksum	3616	R	906	0	6.30	ec2-meta-data

[...]

Example: biolatency

```
# biolatency -m 1 5
Tracing block device I/O... Hit Ctrl-C to end.
```

msecs	:	count	distribution
0 -> 1	:	36	*****
2 -> 3	:	1	*
4 -> 7	:	3	***
8 -> 15	:	17	*****
16 -> 31	:	33	*****
32 -> 63	:	7	*****
64 -> 127	:	6	*****

[...]

GUI Tracing Tools



GUI Visualizations

1. Event logs
2. Tables
3. Line graphs
4. Histograms
5. Heatmaps (spectrographs)
6. Waterfall charts
7. Directed graphs
8. Flame graphs
9. Flame charts

Visualization 1: Event Logs

The screenshot shows the Wireshark interface with the following details:

- Packet List:**

No.	Time	Source	Destination	Protocol	Info
40	139.931167	Wistron_07:07:ee	Broadcast	ARP	who has 192.168.1.254? tell 192.168.1.68
47	139.931463	ThomsonT_08:35:4f	Wistron_07:07:ee	ARP	192.168.1.254 is at 00:90:d0:08:35:4f
48	139.931466	192.168.1.68	192.168.1.254	DNS	Standard query A www.google.com
49	139.975406	192.168.1.254	192.168.1.68	DNS	Standard query response CNAME www.l.google.com A 66.102.9.99
50	139.976811	192.168.1.68	66.102.9.99	TCP	62216 > http [SYN] Seq=0 Win=8192 Len=0 MSS=1460 WS=2
51	140.079578	66.102.9.99	192.168.1.68	TCP	http > 62216 [SYN, ACK] Seq=0 Ack=1 Win=5720 Len=0 MSS=1430
52	140.079583	192.168.1.68	66.102.9.99	TCP	62216 > http [ACK] Seq=1 Ack=1 Win=65780 Len=0
53	140.080278	192.168.1.68	66.102.9.99	HTTP	GET /complete/search?hl=en&client=suggest&js=true&q=m&cp=1 H
54	140.086765	192.168.1.68	66.102.9.99	TCP	62216 > http [FIN, ACK] Seq=805 Ack=1 Win=65780 Len=0
55	140.086921	192.168.1.68	66.102.9.99	TCP	62218 > http [SYN] Seq=0 Win=8192 Len=0 MSS=1460 WS=2
56	140.197484	66.102.9.99	192.168.1.68	TCP	http > 62216 [ACK] Seq=1 Ack=805 Win=7360 Len=0
57	140.197777	66.102.9.99	192.168.1.68	TCP	http > 62216 [FIN, ACK] Seq=1 Ack=806 Win=7360 Len=0
58	140.197811	192.168.1.68	66.102.9.99	TCP	62216 > http [ACK] Seq=806 Ack=2 Win=65780 Len=0
59	140.219219	66.102.9.99	192.168.1.68	TCP	http > 62218 [SYN, ACK] Seq=0 Ack=1 Win=5720 Len=0 MSS=1430

- Packet Details:**

 - Frame 1 (42 bytes on wire, 42 bytes captured)
 - Ethernet II, Src: Vmware_38:eb:0e (00:0c:29:38:eb:0e), Dst: Broadcast (ff:ff:ff:ff:ff:ff)
 - Address Resolution Protocol (request)

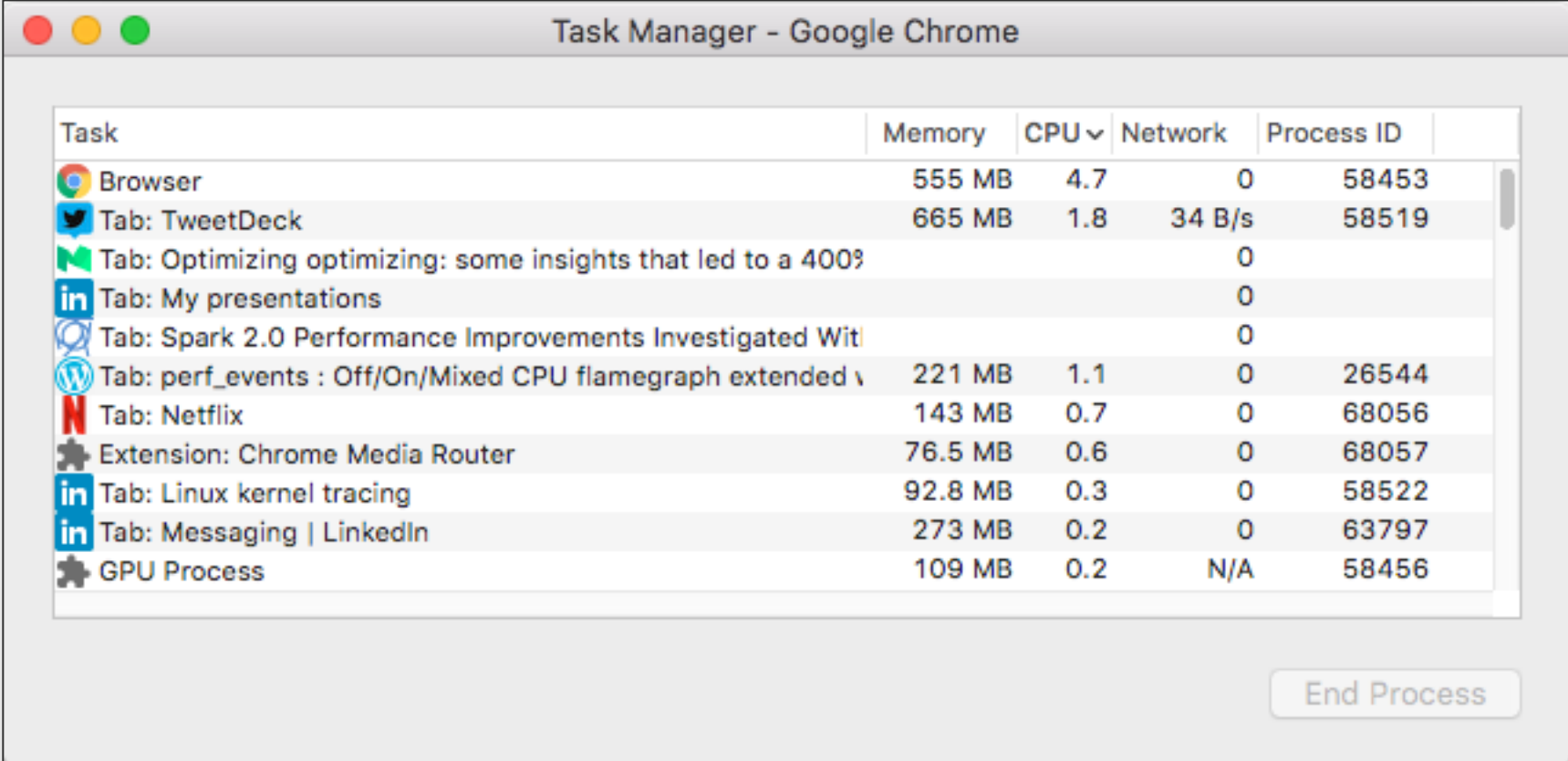
- Raw Data:**

```
0000  ff ff ff ff ff ff 00 0c 29 38 eb 0e 08 06 00 01  ..... )8.....
0010  08 00 06 04 00 01 00 0c 29 38 eb 0e c0 a8 39 80  ..... )8....9.
0020  00 00 00 00 00 00 00 c0 a8 39 02  ..... 9.
```

eth0: <live capture in progress> Fil... Packets: 445 Displayed: 445 Marked: 0 Profile: Default

https://commons.wikimedia.org/wiki/File:Wireshark_screenshot.png

Visualization 2: Tables

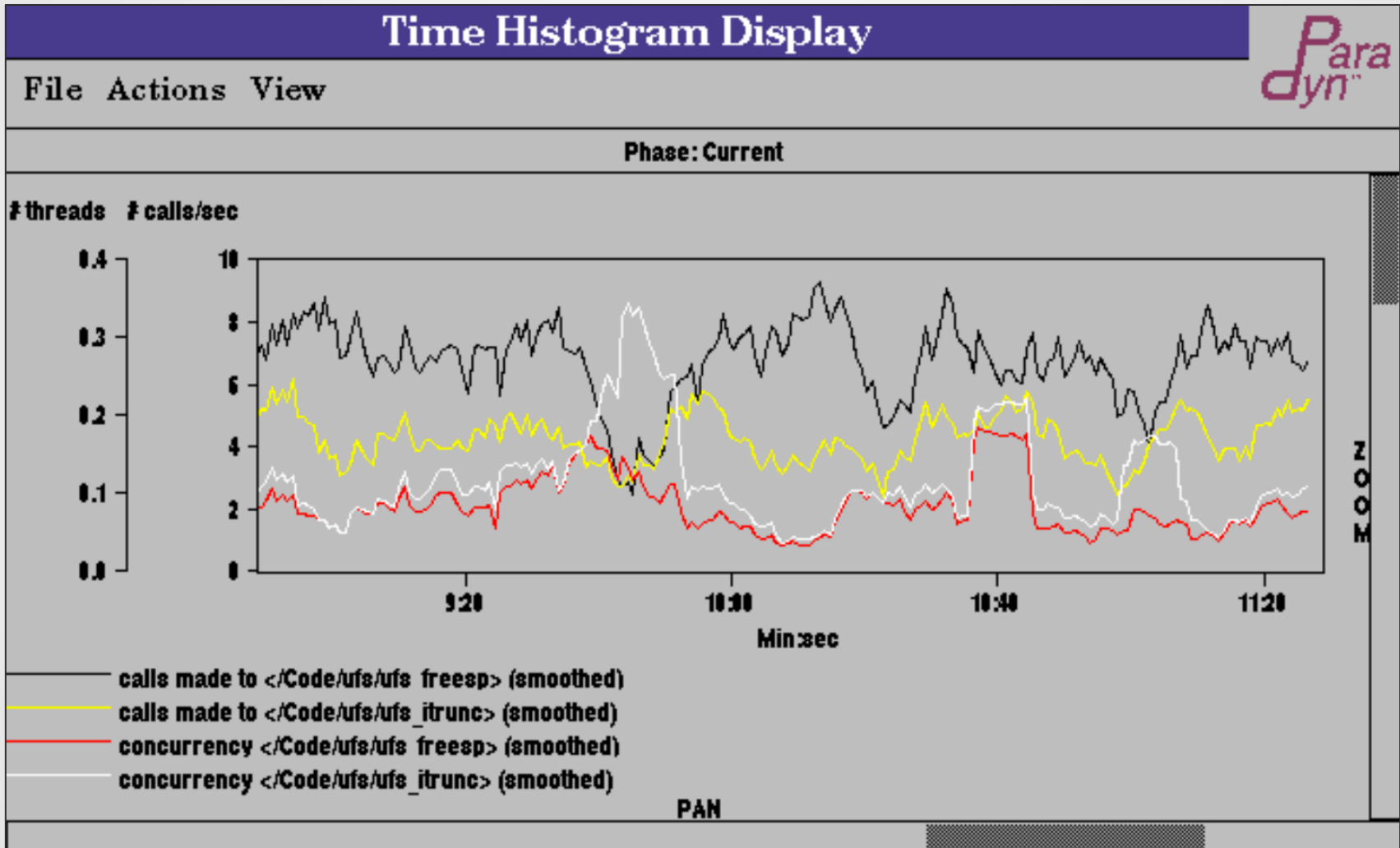


Task Manager - Google Chrome

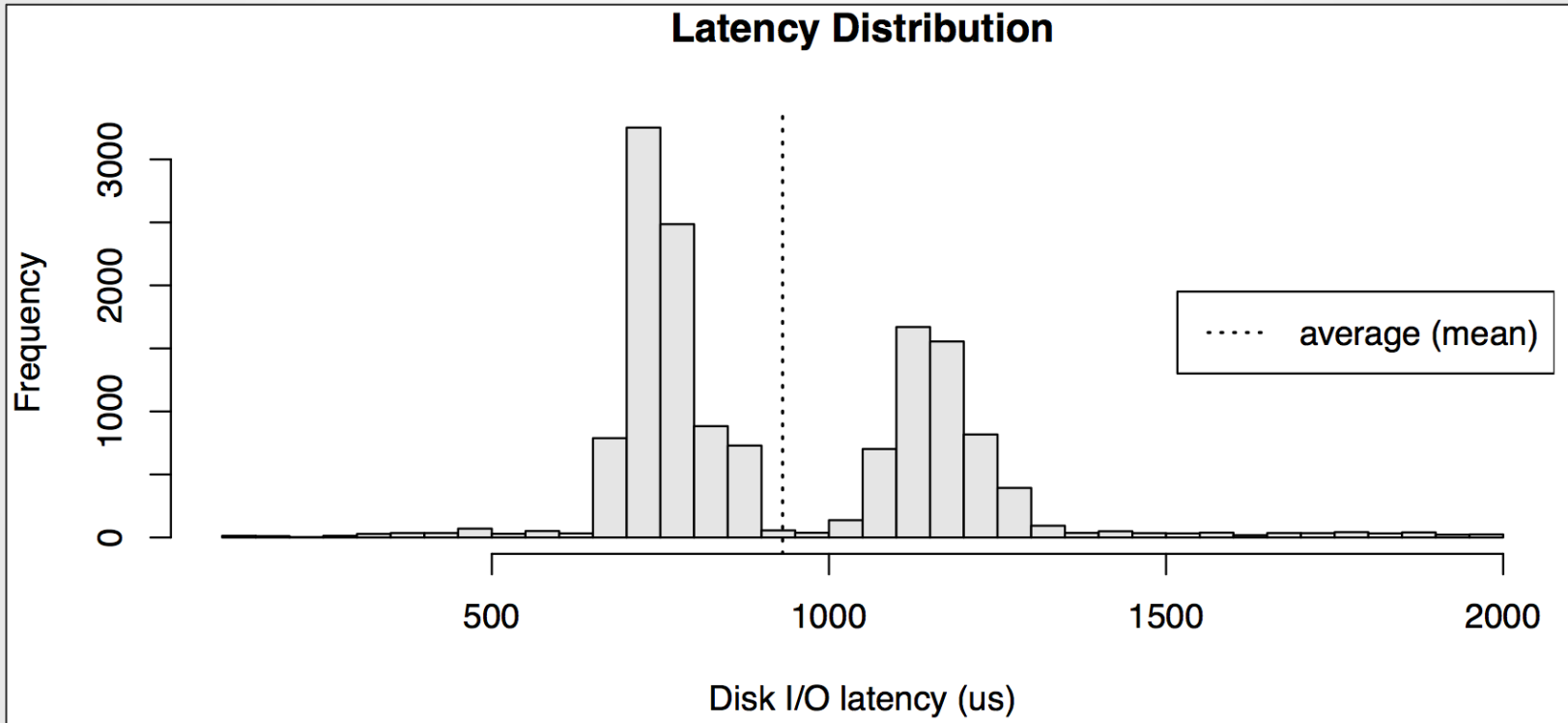
Task	Memory	CPU	Network	Process ID
Browser	555 MB	4.7	0	58453
Tab: TweetDeck	665 MB	1.8	34 B/s	58519
Tab: Optimizing optimizing: some insights that led to a 400%			0	
Tab: My presentations			0	
Tab: Spark 2.0 Performance Improvements Investigated Wit			0	
Tab: perf_events : Off/On/Mixed CPU flamegraph extended v	221 MB	1.1	0	26544
Tab: Netflix	143 MB	0.7	0	68056
Extension: Chrome Media Router	76.5 MB	0.6	0	68057
Tab: Linux kernel tracing	92.8 MB	0.3	0	58522
Tab: Messaging LinkedIn	273 MB	0.2	0	63797
GPU Process	109 MB	0.2	N/A	58456

End Process

Visualization 3: Line Graphs



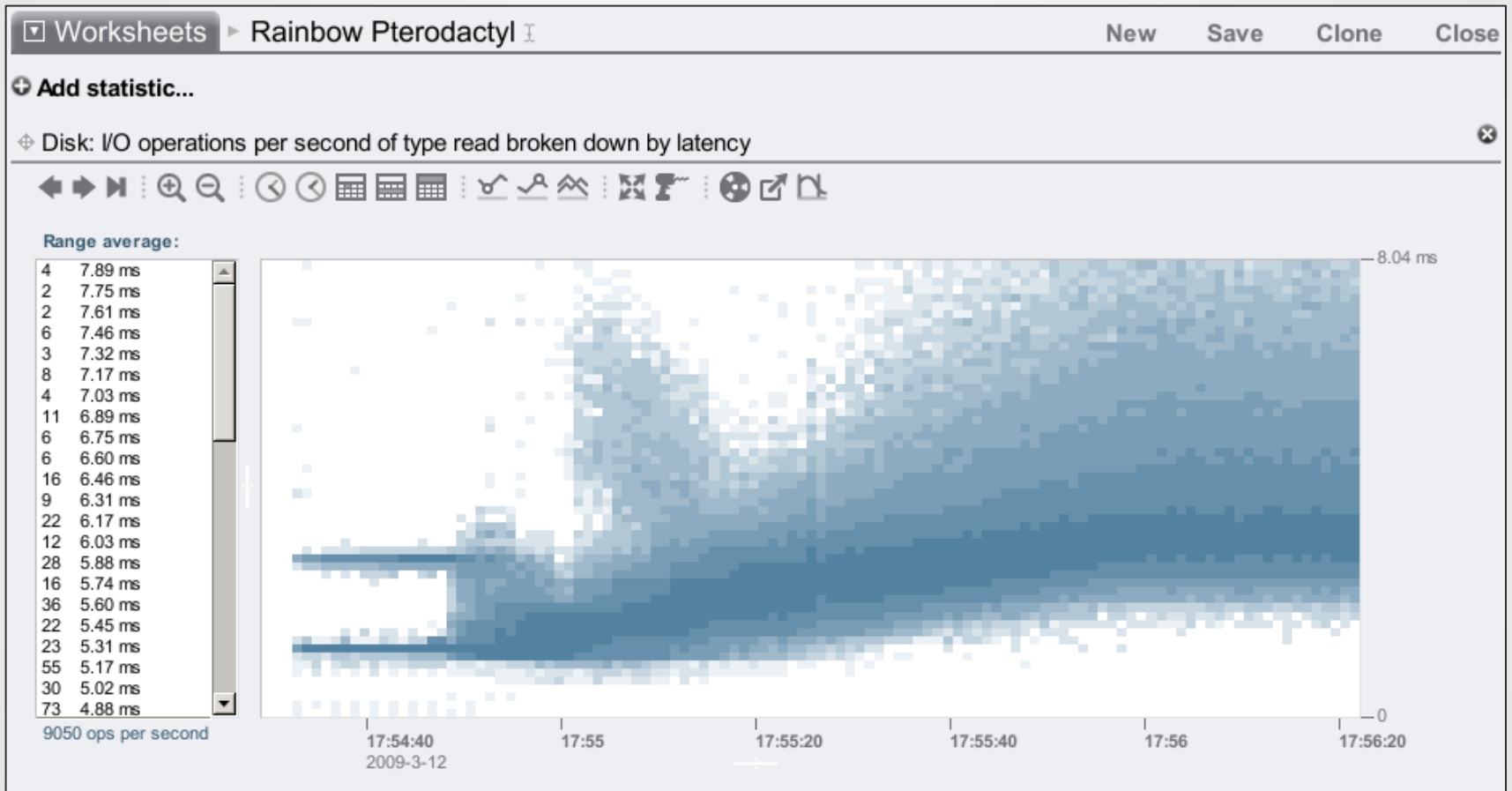
Visualization 4: Histograms



Or a density plot

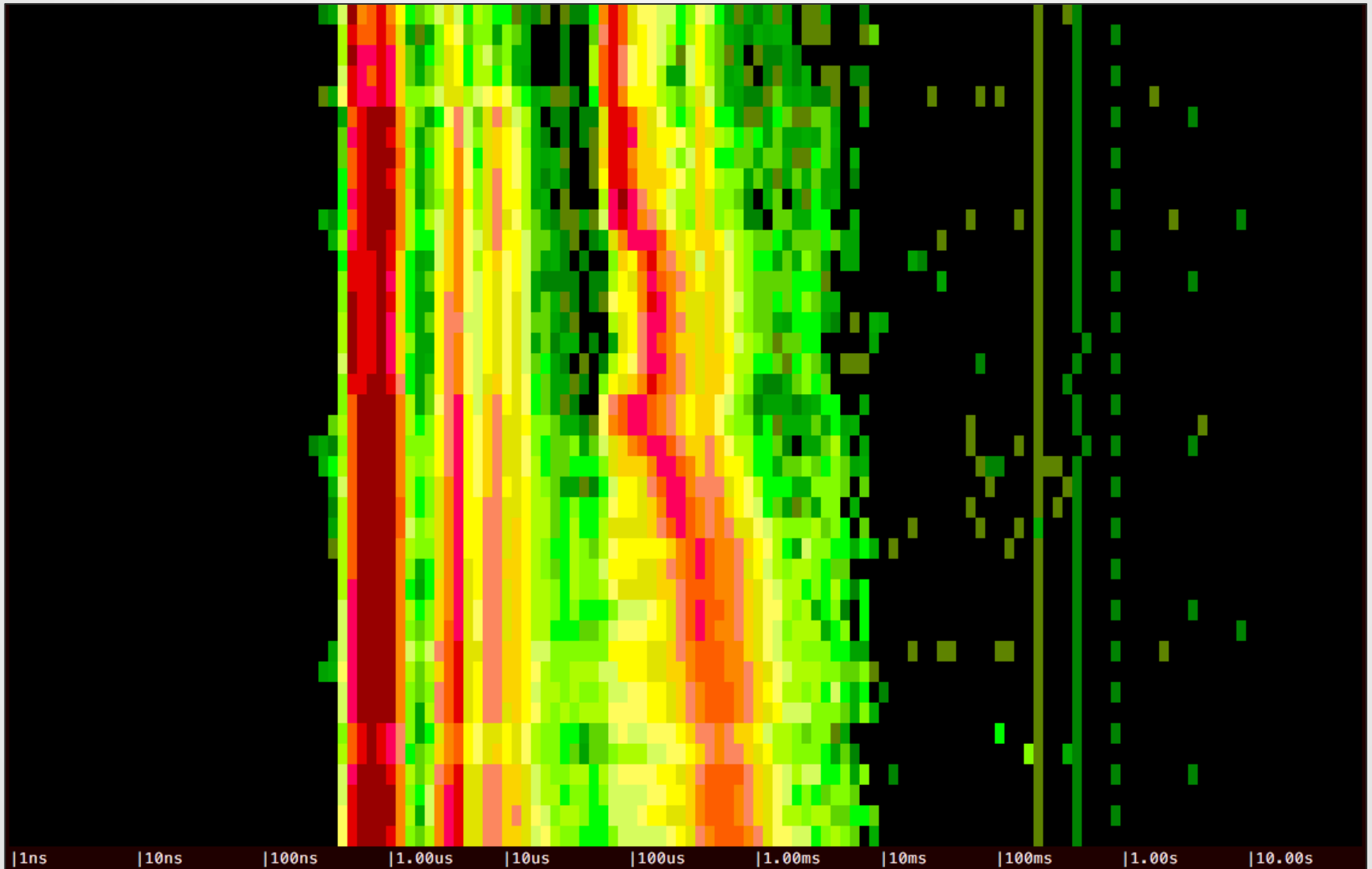
Or as a frequency trail (can cascade)

Visualization 5: Heat Maps

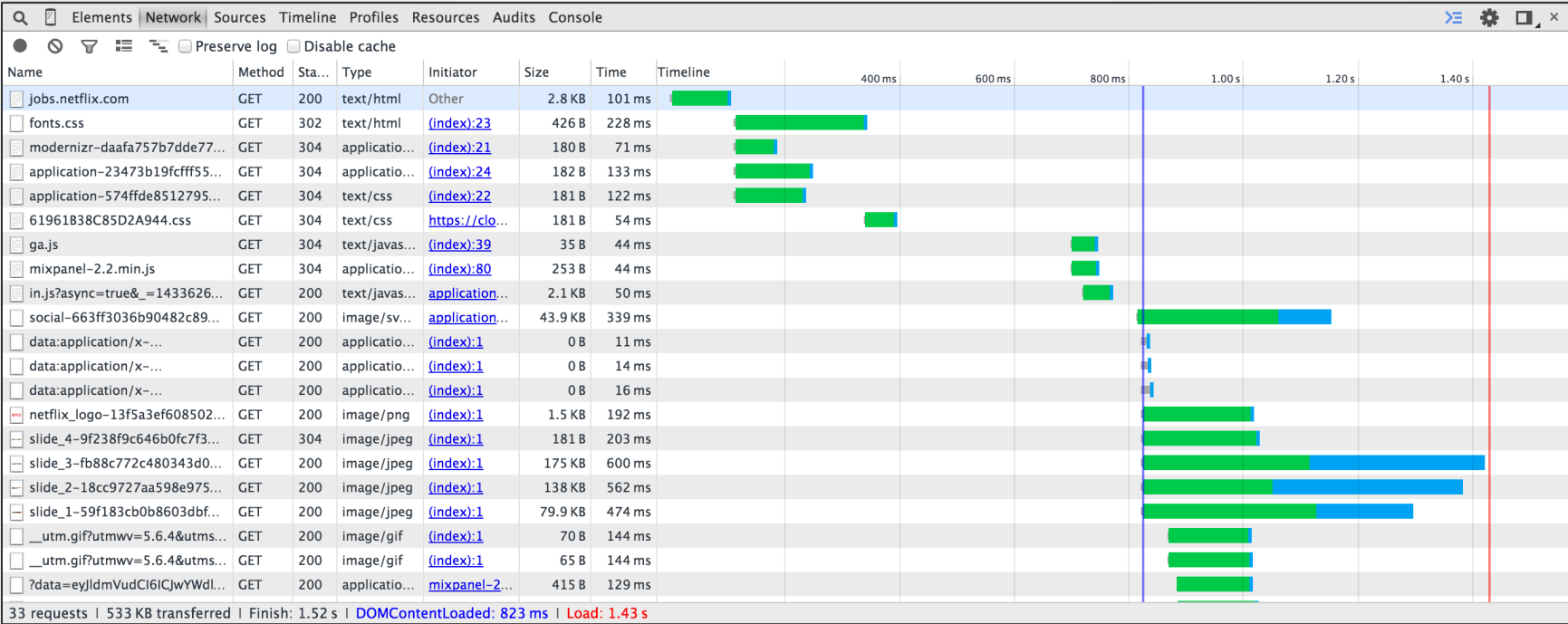


eg, Oracle ZFS Storage Appliance Analytics (DTrace-based)

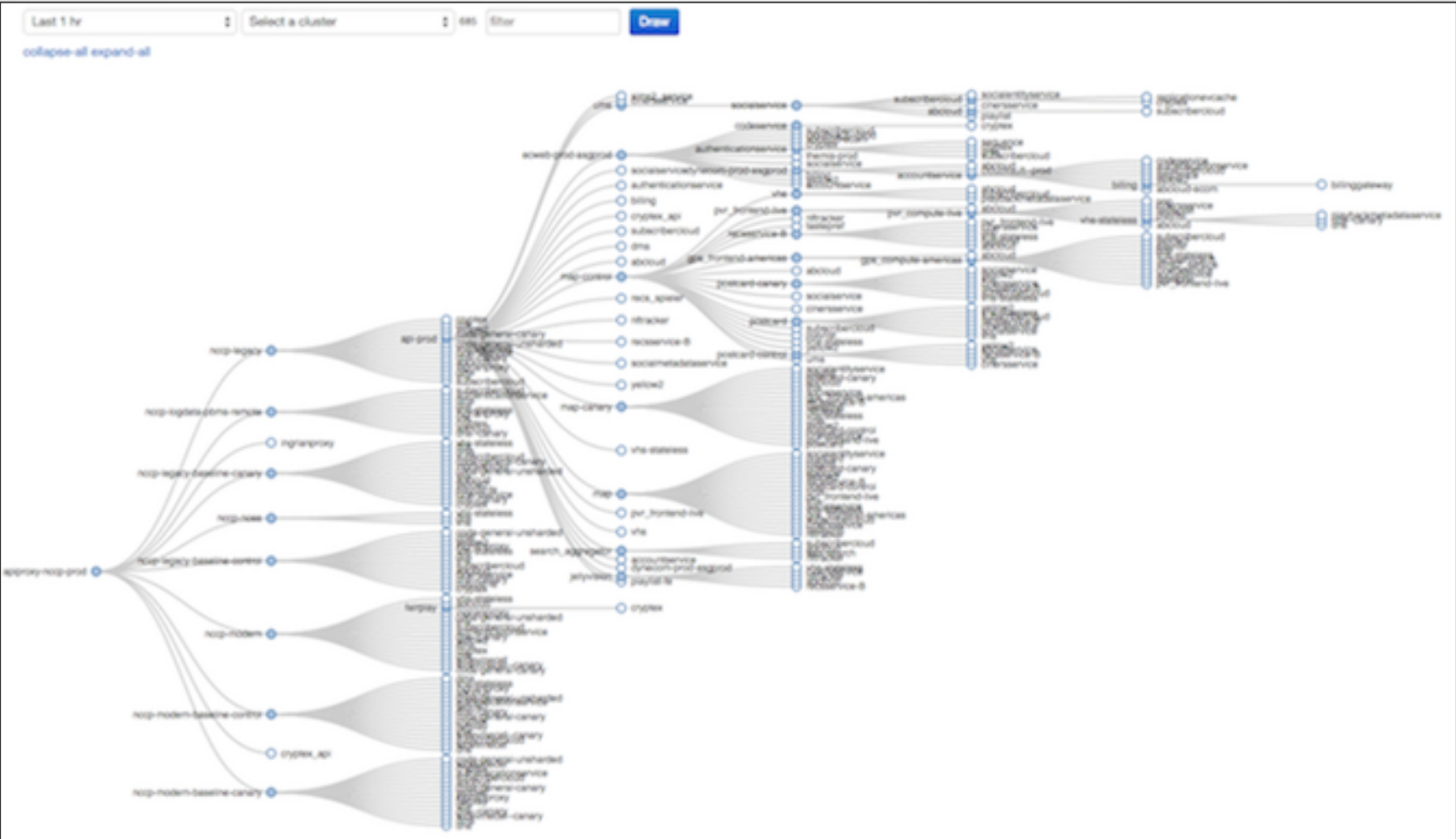
Visualization 5: Spectrograms



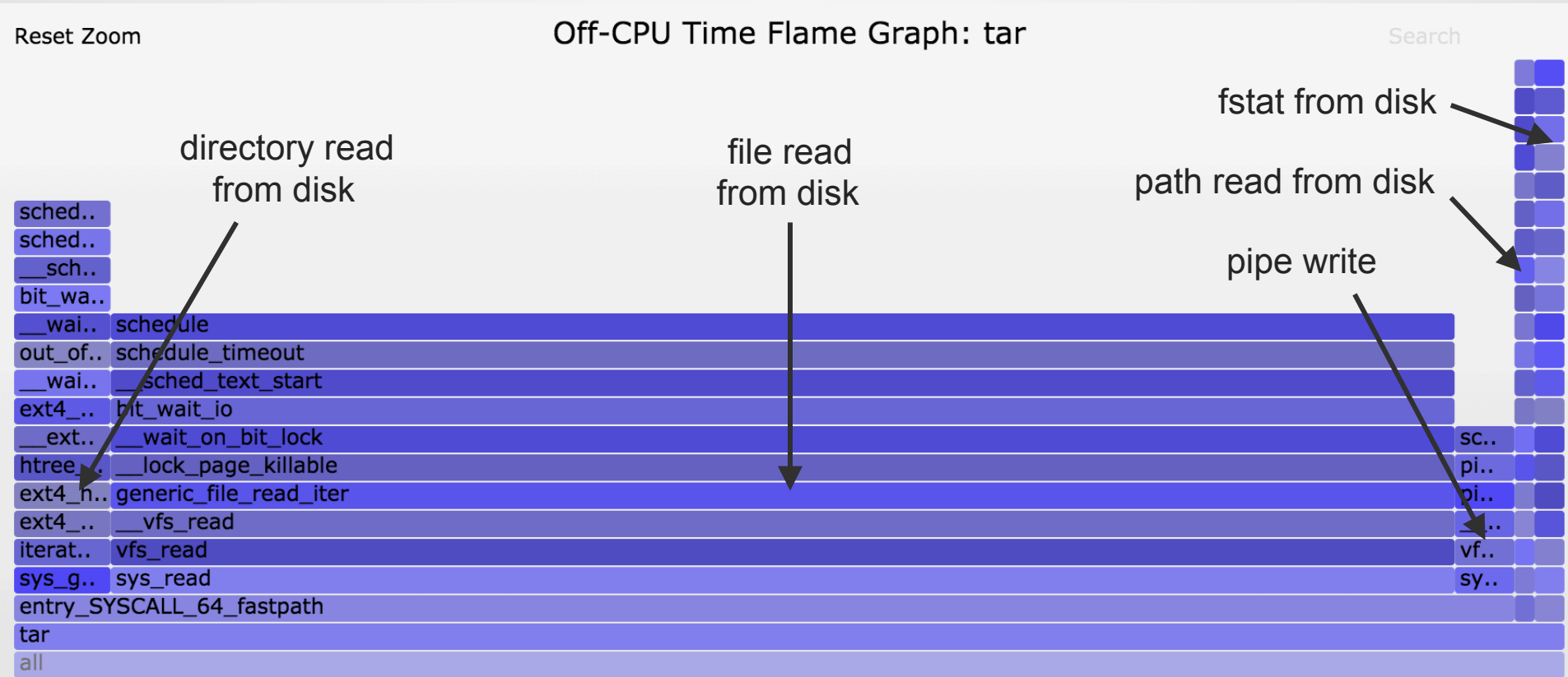
Visualization 6: Waterfall Charts



Visualization 7: Directed Graphs

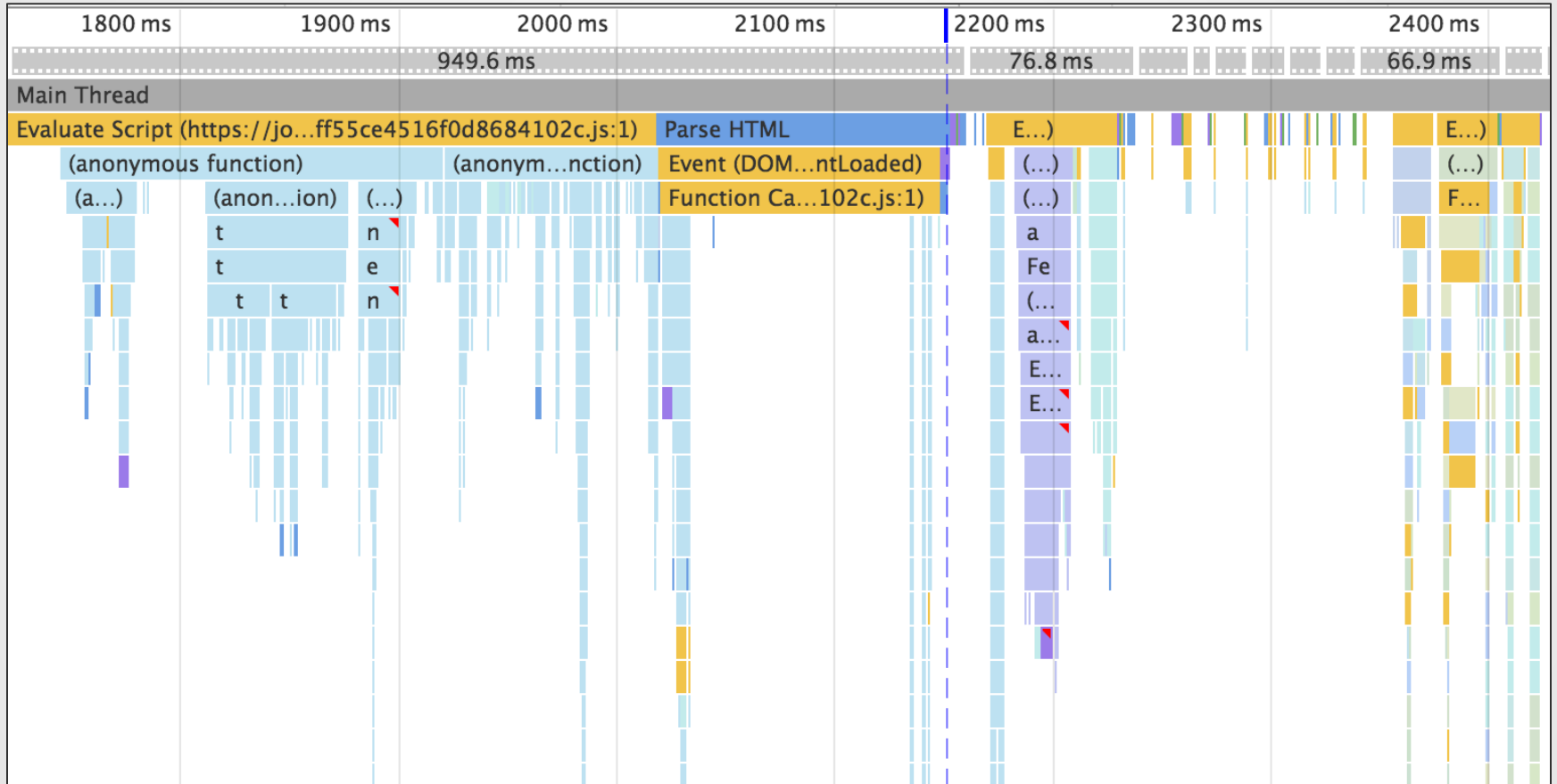


Visualization 8: Flame Graphs



Commonly used with CPU profilers. Also useful for tracers: off-CPU time, ...

Visualization 9: Flame Charts



Desirable Attributes

- Valuable
 - Methodologies provide ideas for purposeful metrics
- Documented
 - Tool tips, wikis
- Tested
- Real Time
- Dashboards
 - To support methodologies

Thank You!

<http://www.brendangregg.com>
<http://slideshare.net/brendangregg>
bgregg@netflix.com
[@brendangregg](#)

References & Links:

- <http://www.brendangregg.com/heatmaps.html>
- <http://www.brendangregg.com/flamegraphs.html>
- <http://www.slideshare.net/brendangregg/monitorama-2015-netflix-instance-analysis>

