# The cost of things at scale

Robert Graham

@ErrataRob

https://blog.erratasec.com

# The cost of things

- How fast can CPUs execute instructions
- How fast can CPUs access-memory
- How fast are kernel system calls
- How fast are synchronization primitives
- How fast are "context-switches"

# Code

- https://github.com/robertdavidgraham/c10mbench

# C10M defined

- 10 million concurrent connections
- 1 million connections/second
- 10 gigabits/second
- 10 million packets/second
- 10 microsecond latency
- 10 microsecond jitter
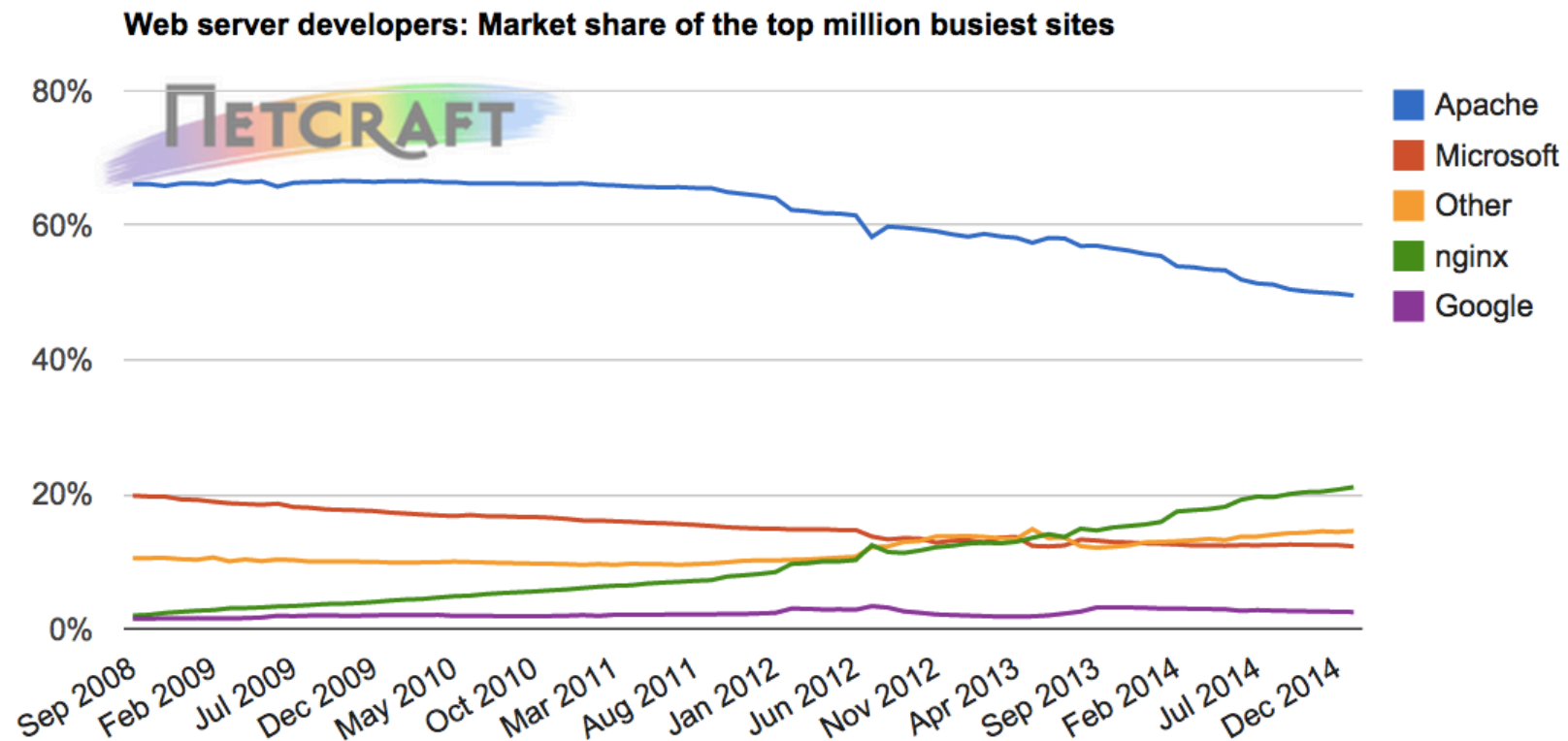- 10 coherent CPU cores

# Classic definition: Context-switch

- Process/thread context switches

# ..but process context switches becoming rare

- NodeJS
- Nginx
- Libevent
- Java user-mode threads
- Lua coroutines

# …but context switches becoming rare



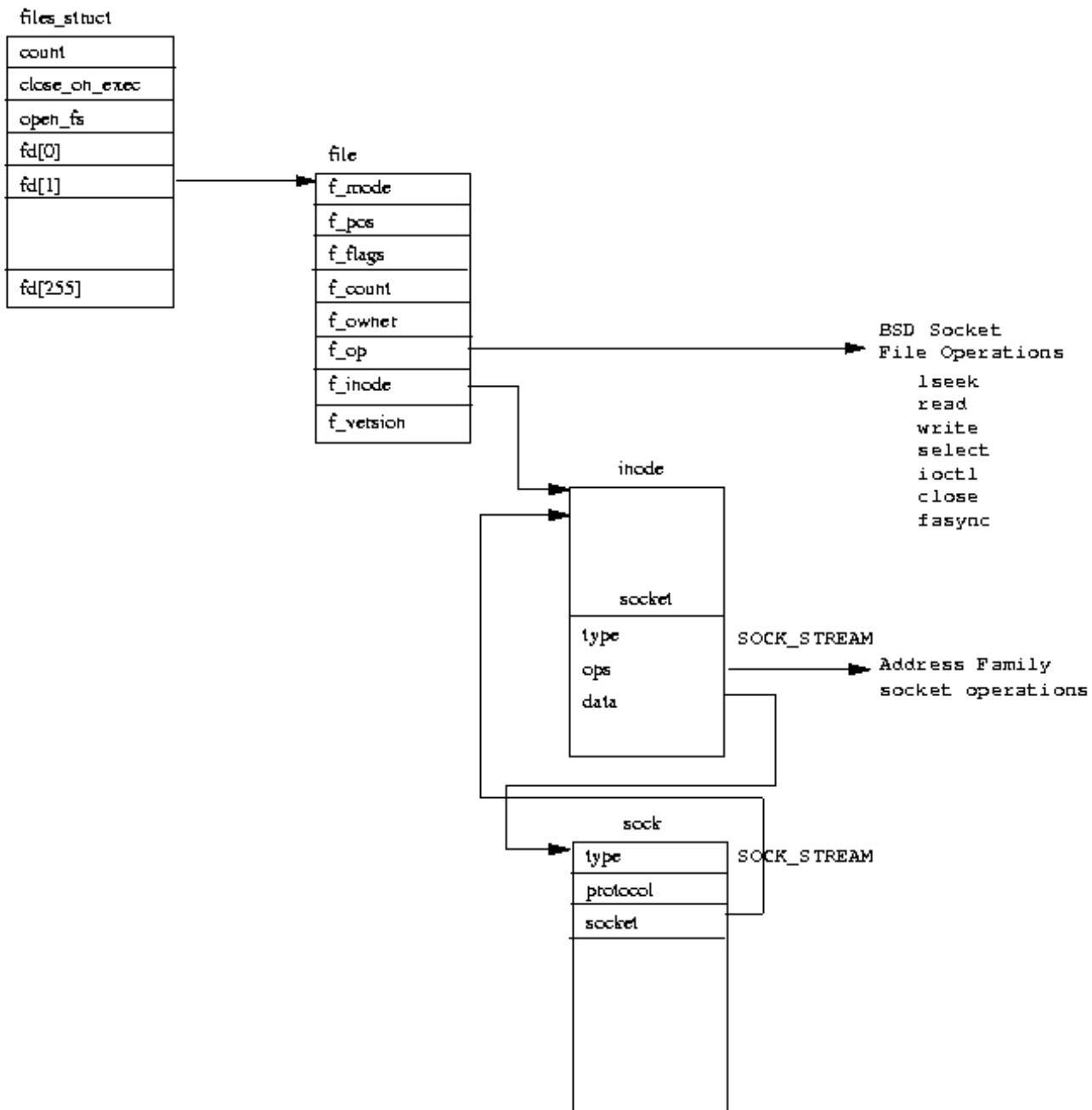Web server developers: Market share of the top million busiest sites

# Real definition: Context-switch

- Each TCP connection is a task, with context
  - Whether you assign a thread to it, a closure, or a data structure
- Each incoming packet causes a random context switch
- A lot of small pieces of memory must be touched – *sequentially*
  - "pointer-chasing"

**files_struct**

| |
|---|
| count |
| close_on_exec |
| open_fs |
| fd[0] |
| fd[1] |
| |
| fd[255] |

**file**

| |
|---|
| f_mode |
| f_pos |
| f_flags |
| f_count |
| f_owner |
| f_op |
| f_inode |
| f_version |

BSD Socket
File Operations

```
    lseek
    read
    write
    select
    ioctl
    close
    fasync
```

**inode**

| |
|---|
| |
| |

**socket**

| |
|---|
| type |
| ops |
| data |

SOCK_STREAM

Address Family
socket operations

**sock**

| |
|---|
| type |
| protocol |
| socket |
| |

SOCK_STREAM
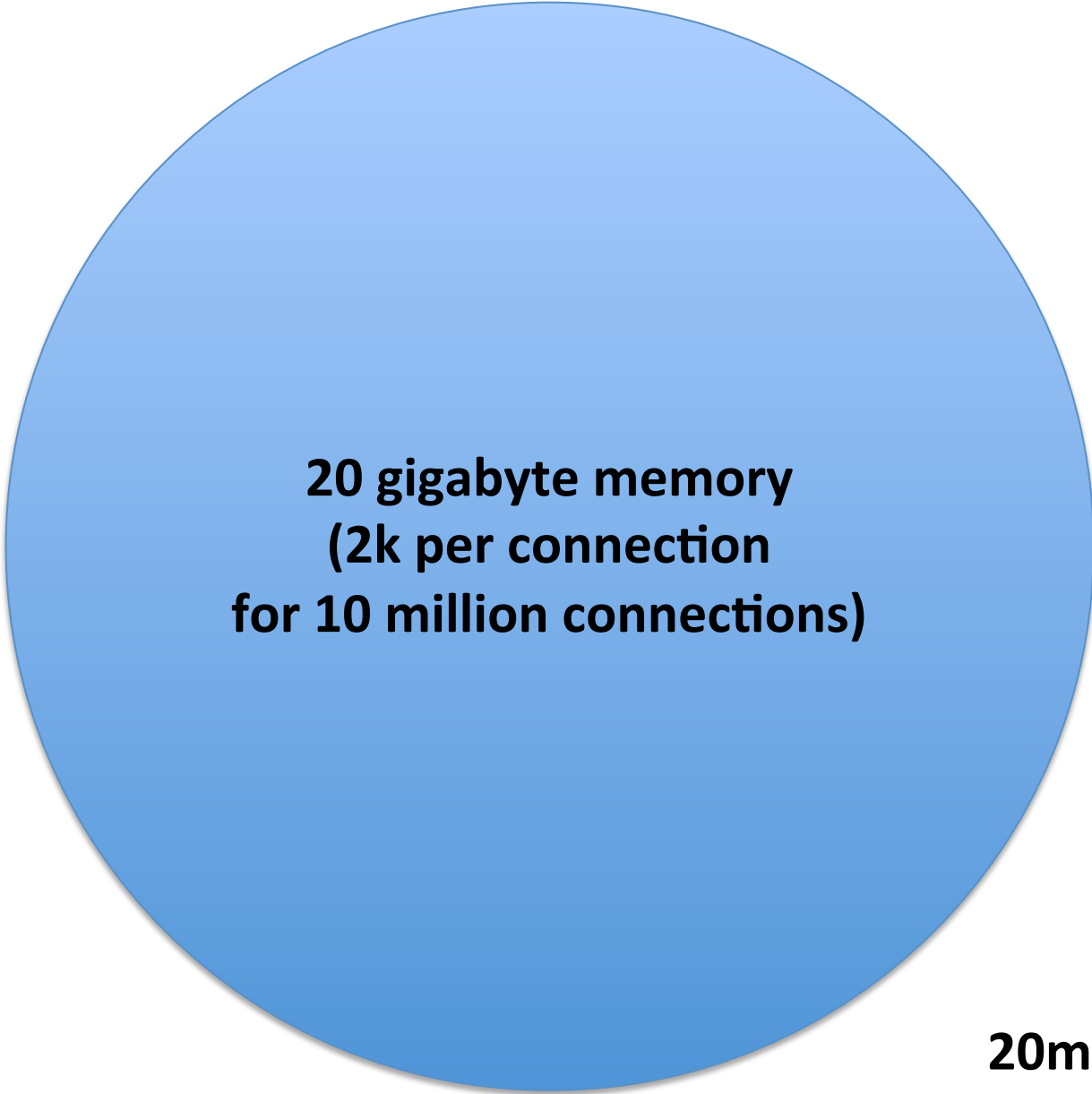
CPU

L1 cache—4 cycles

L2 cache—12 cycles

L3 cache—30 cycles

main memory—300 cycles

20 gigabyte memory
(2k per connection
for 10 million connections)

20meg L3 cache

# Measured latency: 85ns

**Concurrent memory latency**

# budget

10 million packets/second
divided by 10 cores
by 100 nanoseconds/miss
------------------
 ~10 cache misses per packet

# Now for user-mode

- Apps written in C have few data structures
- Apps written in high-level languages (Java, Ruby, Lua, JavaScript) have bits of memory strewn around

# User-mode memory is virtual

- Virtual addresses are translated to physical addresses on every memory access
  - Walk a chain of increasingly smaller page table entries
- But TLB cache makes it go fast
  - But not at scale
  - TLB cache is small
  - Page tables themselves may not fit in the cache

# Small Page Diagram for x64 Virtual Memory

64　　　　　48　　　39　　　30　　　21　　　12　　　0

```
page_walk(addr) {
 L4 = CSR3;
 L3 = L4[(addr>>39)&0x1FF];
 L2 = L3[(addr>>30)&0x1FF];
 if (L2[(addr>>21)&0x1FF] & 1) {
  page = L2[(addr>>21)&0x1FF] ^ 1;
  return page | addr & 0x1FFFFF;
 } else {
  L1 = L2[(addr>>21)&0x1FF];
  page = L1[(addr>>12)&0x1FF];
  return page | addr & 0xFFF;
 }
}
```

Page Map Level 4

Page Directory Pointer Table

Page Directory

Page Table

Memory Page

0

CSR3

L4　　　L3　　　L2　　　L1

4096

# Large Page Diagram for x64 Virtual Memory

```
page_walk(addr) {
 L4 = CSR3;
 L3 = L4[(addr>>39)&0x1FF];
 L2 = L3[(addr>>30)&0x1FF];
 if (L2[(addr>>21)&0x1FF] & 1) {
  page = L2[(addr>>21)&0x1FF] ^ 1;
  return page | addr & 0x1FFFFF;
 } else {
  L1 = L2[(addr>>21)&0x1FF];
  page = L1[(addr>>12)&0x1FF];
  return page | addr & 0xFFF;
 }
}
```

64   48   39   30   21   0

Page Map Level 4

Page Directory Pointer Table

Page Directory

Memory Page

CSR3

L4    L3    L2

2-megabytes

20 gigabyte memory
(2k per connection
for 10 million connections)

10k hugepage tables

20meg L3 cache

40meg small page tables

# User-mode latency



**Concurrent memory latency**

# QED:

- Memory latency becomes a big scalability problem for high-level languages

# How to solve

- Hugepages to avoid page translation
- Break the chain
  - Add "void *prefetch[8]" to the start of every TCP control block.
  - Issue prefetch instructions on them as soon as packet arrives
  - Get all the memory at once

# Memory access is parallel

- CPU
  - Each core can track 72 memory reads at the same time
  - Entire chip can track ?? reads at the same time
- DRAM
  - channels X slots X ranks X banks
  - My computer: 3 * 2 * 1 * 4 = 24 concurrent accesses
  - Measured: 190-million/sec = 15 concurrent accesses

# Some reading

- "What every programmer should know about memory" by Ulrich Draper
- http://www.akkadia.org/drepper/cpumemory.pdf

# Multi-core

# Multi-threading is not the same as multi-core

- Multi-threading
  - More than one thread per CPU core
  - Spinlock/mutex must therefore stop one thread to allow another to execute
  - Each thread a different task (multi-tasking)
- Multi-core
  - One thread per CPU core
  - When two threads/cores access the same data, they can't stop and wait for the other
  - All threads part of the same task

# Most code doesn't scale past 4 cores

# #1 rule of multi-core: don't share memory

- People talk about ideal mutexes/spinlocks, but they still suffer from shared memory

- There is exist data structures, "lock free", that don't require them

# Let's measure the problem

- A "locked add" simulates the basic instructions behind spinlocks, futexes, etc.

```c
static void
worker_thread(void *parms)
{
    size_t i;
    for (i=0; i<BENCH_ITERATIONS2; i++) {
        pixie_locked_add_u32(&result, 1);
    }
}
```

# Total additions per second

**Incrementing a shared memory**

# Latency per addition per thread



**Latency per addition operation per core**

# Two things to note

- ~5 nanoseconds
  - Cost of an L3 cache operation (~10ns)
  - Minus the out-of-order execution by the CPU (~5ns)
  - ...and I'm still not sure
- ~100 nanoseconds
  - When many thread contending, it becomes as expensive as a main memory operation

# Syscalls

- Mutexes often done with system calls
- So what's the price of a such a call?
  - On my machine
  - ~30 nanoseconds is minimum
  - ~60 ns is more typical idealized cases
  - ~400 ns in more practical cases

# Solution: lock-free ring-buffers

- No mutex/spinlock

- No syscalls

- Since head and tail are separate, no sharing of cache lines

- Measured on my machine:
  - 100-million msgs/second
  - ~10ns per msg

Fetch most recent Data

Add new Data

Ring Buffer

Release old Data

# Shared ring vs. pipes

- Pipes
  - ~400ns per msg
  - 2.5 m-msgs/sec
- Ring
  - ~10ns per msg
  - 100 m-msgs/sec

```c
static void
reader(void *parms)
{
    int fd = *(int*)parms;
    size_t i;

    for (i=0; i<BENCH_ITERATIONS; i++) {
        int x;
        char c;

        x = read(fd, &c, 1);
        if (x != 1)
            break;
    }
}
```

# Function call overhead

- ~1.8ns
- Note the jump for "hyperthreading"
  - My machine has 6 hyperthreaded cores
- 6 clock cycles

**Function pointer latency**

# DMA isn't

# Linux kernel map

| functions layers | system | processing | memory | storage | networking | human interface |
|---|---|---|---|---|---|---|
| | kernel/ | kernel/ | mm/ | fs/ | net/ | |



**user space interfaces**

**system** — kernel/
**system interfaces**
linux/syscalls.h · system files
asm-i386/uaccess.h · /proc /dev
copy_from_user · /sysfs
cdev = cdev_add
register_chrdev
cdev_map · sysfs_ops
sys_reboot · sys_init_module

**processing** — kernel/
**processes**
sys_execve · sys_fork
fs/exec.c · sys_vfork
· sys_clone
linux_binfmt
sys_nanosleep

**memory** — mm/
**memory access**
sys_brk
sys_mmap2
/proc/self/maps
do_path_lookup · sys_open
· sys_read
· sys_write
sys_mount
sys_sync

**storage** — fs/
**files & directories access**
sys_open
sys_read
sys_write
sys_mount
sys_sync

**networking** — net/
**sockets access**
sys_socketcall
sys_socket
socket_file_ops

**human interface**
**HI char devices**
cdev · kmsg
· sys_syslog
input_fops · snd_fops
console_fops · video_fops
fb_fops
input

**virtual**

**Device Model**
drivers/base/ · kobject
kset
subsystem_register
bus_type
class
device
class_device
class_device_create
device_create
device_driver
probe · driver_register

**threads**
work_struct · workqueue_struct
create_workqueue
kthread_create
kernel_thread
current
thread_info · do_fork

**Virtual memory**
virtually continues memory
find_vma_prepare
vmalloc
vmlist
vm_struct
virt_to_page

**Virtual File System**
file · vfs_read
· vfs_write
inode · vfs_create
inode_operations
file_operations
file_system_type
get_sb
super_block

**protocol families**
sock_create · socket
inet_family_ops
inet_create · unix_family_ops
udp_ops
inet_dgram_ops inet_stream_ops

**security**
security/ · linux/security.h
may_open
security_socket_create
security_inode_create
security_ops
selinux_ops
© 2007 Constantine Shulyupin
www.LinuxDriver.co.il/kernel_map
printk

**bridges**

**Synchronization**
wake_up · completion
mutex_lock
add_timer · mutex · down
semaphore
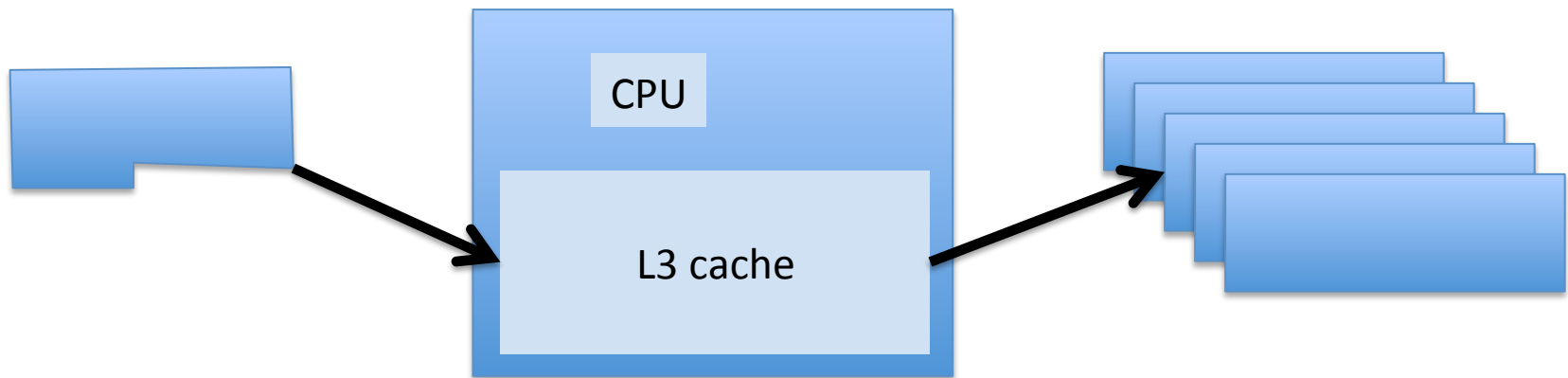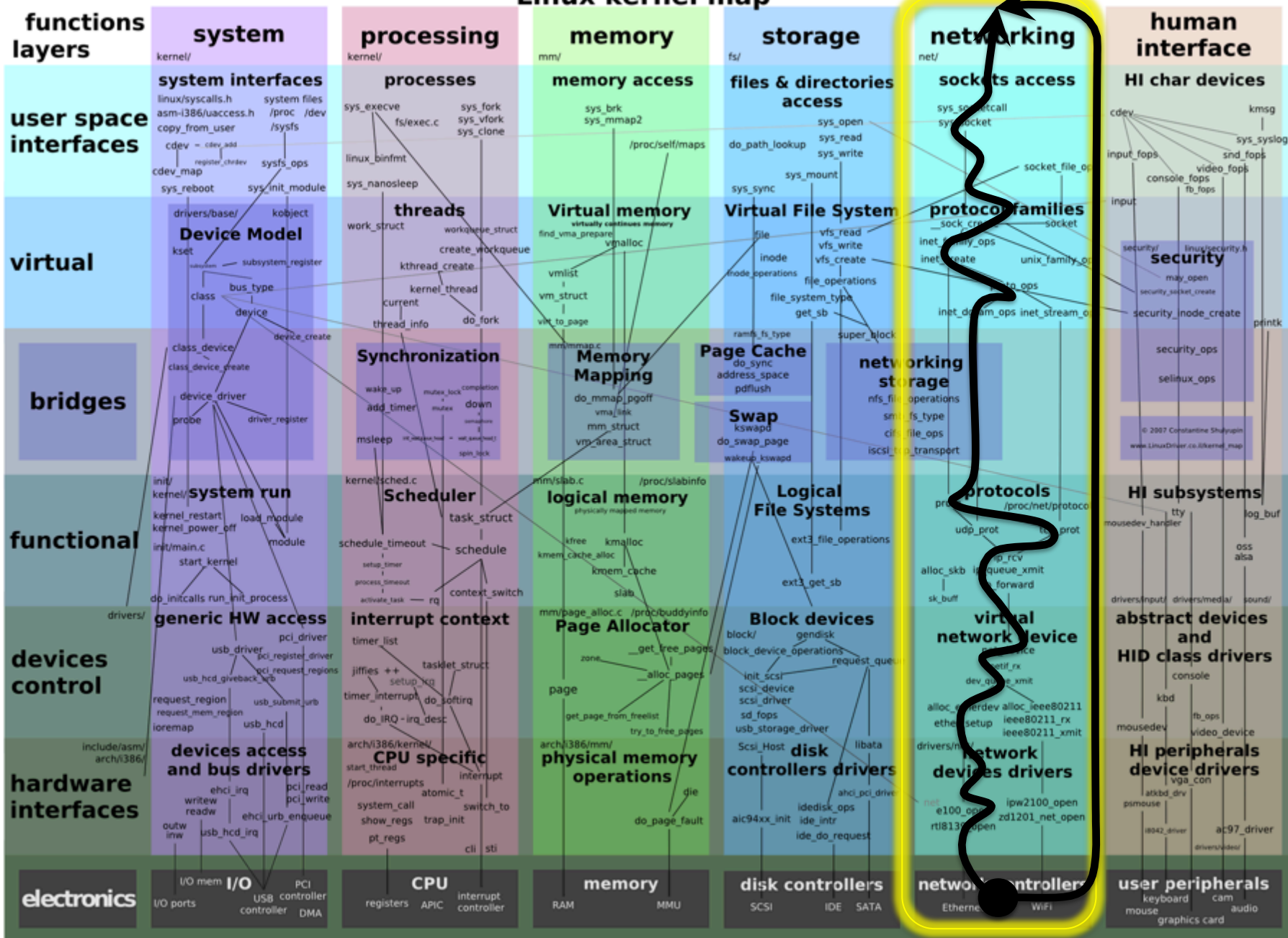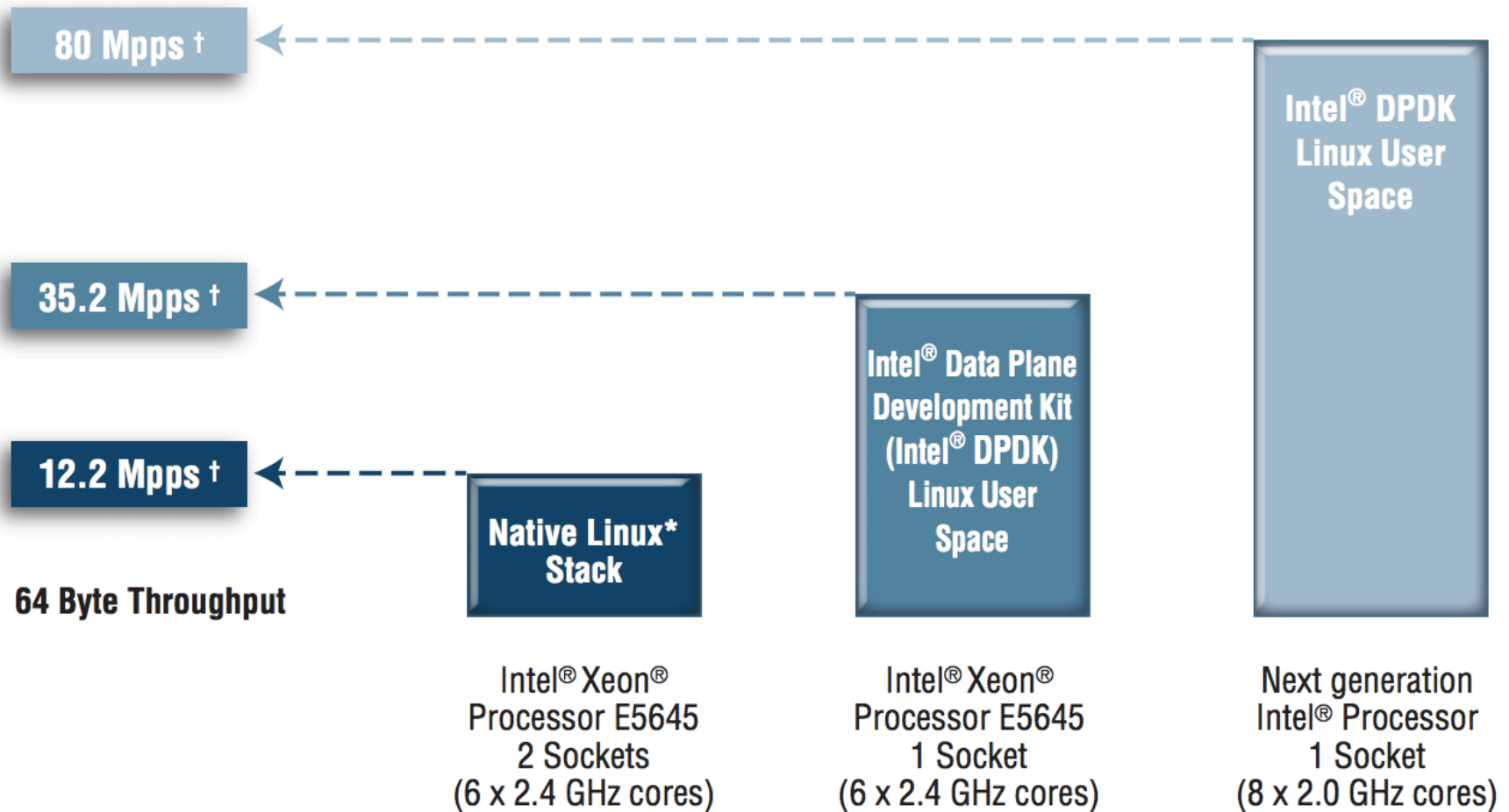msleep · init_mutex_head · init_sem_head
spin_lock

**Memory Mapping**
mm/mmap.c
do_mmap_pgoff
vma_link
mm_struct
vm_area_struct

**Page Cache**
do_sync
address_space
pdflush

**Swap**
kswapd
do_swap_page
wakeup_kswapd

**networking storage**
nfs_file_operations
smb_fs_type
cifs_file_ops
iscsi_tcp_transport

**functional**

**system run**
init/kernel/
kernel_restart · load_module
kernel_power_off
init/main.c · module
start_kernel
do_initcalls run_init_process
drivers/

**Scheduler**
kernel/sched.c
task_struct
schedule_timeout · schedule
setup_timer
process_timeout
activate_task · rq
context_switch

**logical memory**
physically mapped memory
kfree · kmalloc
kmem_cache_alloc
kmem_cache
slab

**Logical File Systems**
ext3_file_operations
ext3_get_sb

**protocols**
/proc/net/protocols
udp_prot · tcp_prot
tcp_rcv
alloc_skb · ip_queue_xmit
· forward
sk_buff

**HI subsystems**
tty · log_buf
mousedev_handler
oss
alsa

**devices control**

**generic HW access**
pci_driver
usb_driver · pci_register_driver
usb_register_driver · pci_request_regions
usb_hcd_giveback_urb
request_region · usb_submit_urb
request_mem_region
ioremap · usb_hcd

**interrupt context**
timer_list
jiffies ++ · tasklet_struct
timer_interrupt · setup_irq
do_softirq
do_IRQ - irq_desc

**Page Allocator**
block/ · gendisk
__get_free_pages
zone
__alloc_pages
page
get_page_from_freelist
try_to_free_pages

**Block devices**
block/ · gendisk
block_device_operations
request_queue
init_scsi
scsi_device
scsi_driver
sd_fops
usb_storage_driver

**virtual network device**
netif_rx
dev_queue_xmit
alloc_etherdev alloc_ieee80211
ether_setup · ieee80211_rx
ieee80211_xmit

**abstract devices and HID class drivers**
console
kbd
fb_ops
mousedev · video_device

**hardware interfaces**

**devices access and bus drivers**
include/asm/
arch/i386/
ehci_irq · pci_read
writew · pci_write
readw
request_region · ehci_urb_enqueue
outw · usb_hcd_irq
inw

**CPU specific**
arch/i386/kernel/
start_thread · interrupt
/proc/interrupts
system_call · atomic_t
show_regs · switch_to
· trap_init
pt_regs
cli sti

**physical memory operations**
arch/i386/mm/
die
do_page_fault

**disk controllers drivers**
Scsi_Host · libata
ahci_pci_driver
aic94xx_init · ide_intr
idedisk_ops
ide_do_request

**network devices drivers**
drivers/net/
net_device
e100_open
ipw2100_open
rtl8139_open zd1201_net_open

**HI peripherals device drivers**
vga_con
atkbd_drv
psmouse
i8042_driver · ac97_driver
drivers/video/

**electronics**

| **I/O** | **CPU** | **memory** | **disk controllers** | **network controllers** | **user peripherals** |
|---|---|---|---|---|---|
| I/O mem · PCI controller | registers · APIC interrupt controller | RAM · MMU | SCSI · IDE · SATA | Ethernet · WiFi | keyboard · cam audio · mouse · graphics card |
| I/O ports · USB controller · DMA | | | | | |

Ver 0.6, 1/1/2008

# Where can I get some?

- PF_RING
  - Linux
  - open-source

- Netmap
  - FreeBSD
  - open-source

- Intel DPDK
  - Linux
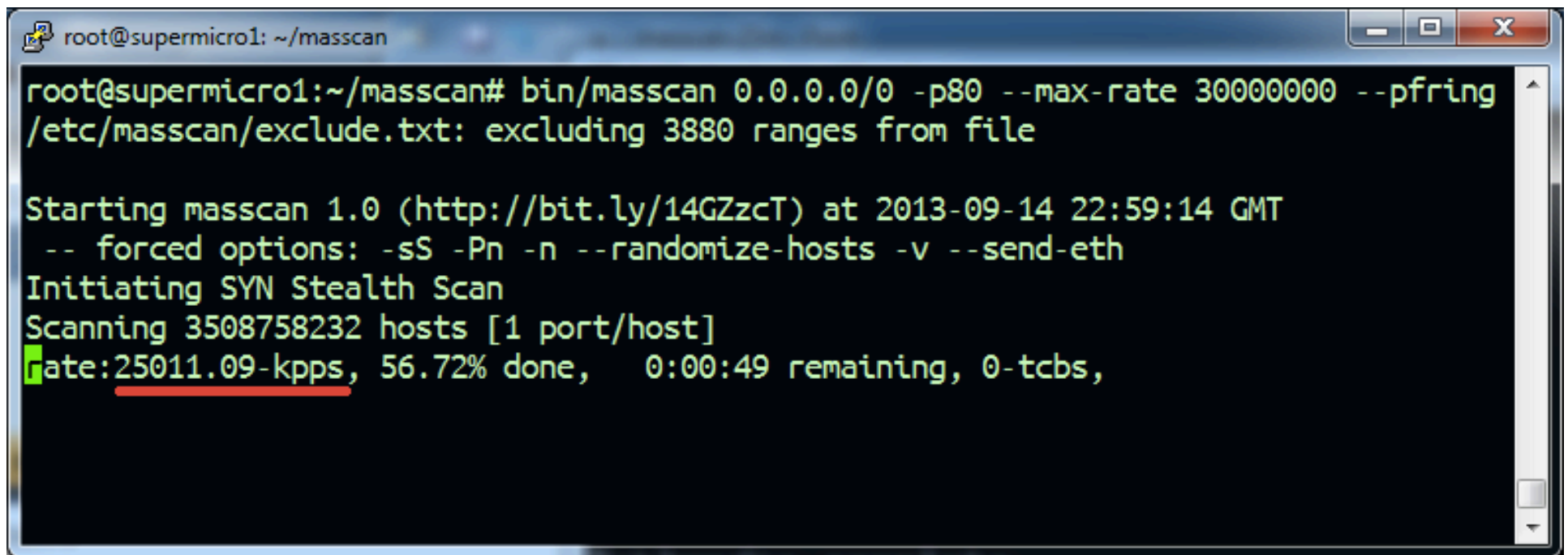  - License fees
  - Third party support
    - 6WindGate

# 200 CPU clocks per packet



80 Mpps †

35.2 Mpps †

12.2 Mpps †

**64 Byte Throughput**

**Intel® DPDK Linux User Space**

**Intel® Data Plane Development Kit (Intel® DPDK) Linux User Space**

**Native Linux* Stack**

Intel® Xeon® Processor E5645 2 Sockets (6 x 2.4 GHz cores)

Intel® Xeon® Processor E5645 1 Socket (6 x 2.4 GHz cores)

Next generation Intel® Processor 1 Socket (8 x 2.0 GHz cores)

# masscan

- Quad-core Sandy Bridge 3.0 GHz
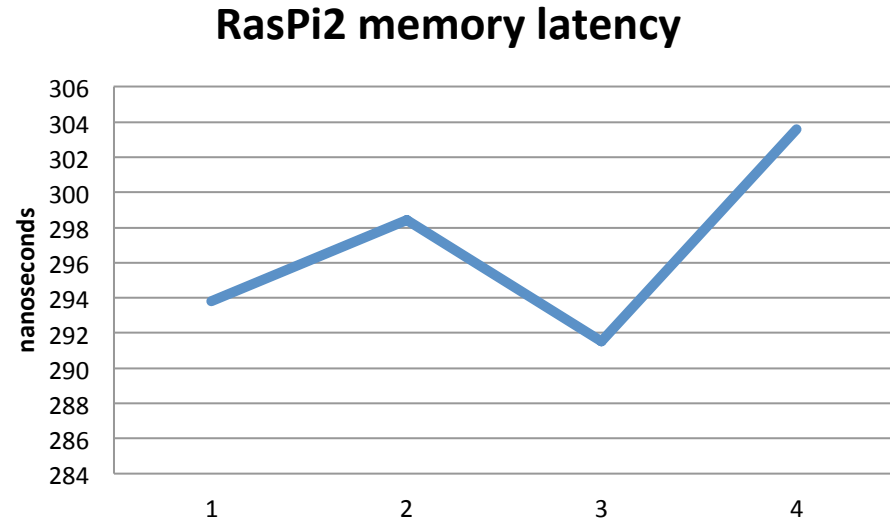
# Premature optimization is good

- Start with prototype that reaches theoretical max
  - Then work backwards
- Restate the problem so that it can be solved by the best solutions
  - Ring-buffers and RCU (read-copy-update) are the answers, find problems solved by them
- Measure and identify bottlenecks as they occur

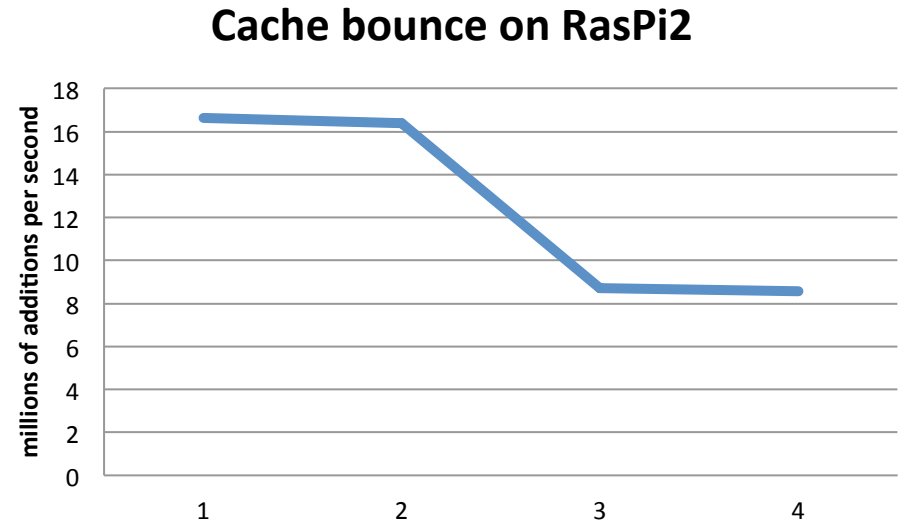# Raspberry PI 2

900 MHz quad core ARM w/ GPU

# Memory latency

- High latency Probably due to limited TLB resources

- Didn't test max outstanding transactions, but should be high for GPU

**RasPi2 memory latency**

# Cache Bounce

- Seems strange

- No performance
  loss for two threads

**Cache bounce on RasPi2**



- Answer: ARM Cortex-A8 comes in 2-cpu
  modules that share cache

# Compared to x86

- .

|         | ARM    | x86    | Speedup |
|---------|--------|--------|---------|
| Hz      | 0.900  | 3.2    | 3.6     |
| syscall | 0.99   | 2.5    | 2.6     |
| funcall | 59.90  | 556.4  | 9.3     |
| pipe    | 0.17   | 2.5    | 14.8    |
| ring    | 3.90   | 74.0   | 19.0    |

# Todo:

- C10mbench work
  - More narrow benchmarks to test things
  - Improve benchmarks
  - Discover exactly why benchmarks have the results they do
  - Benchmark more systems
    - Beyond ARM and x86