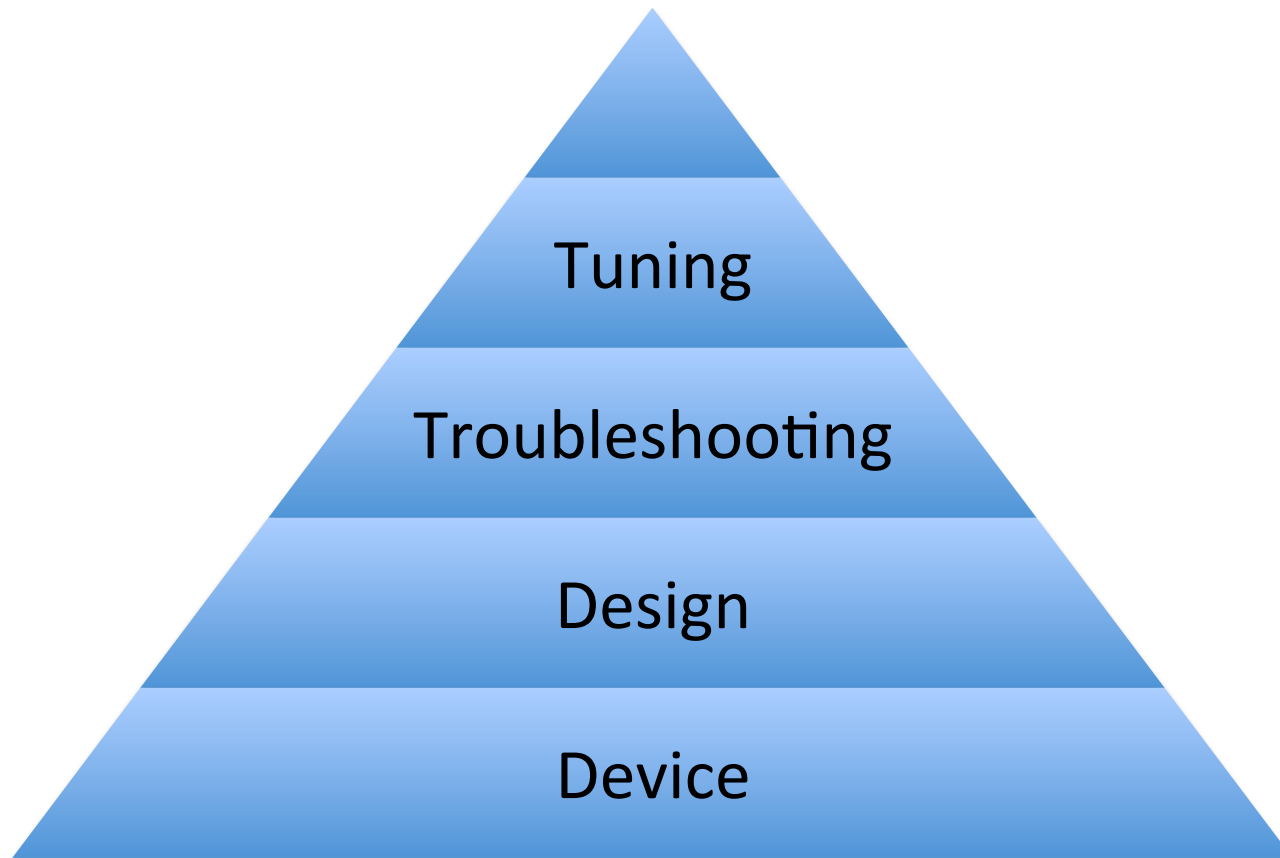


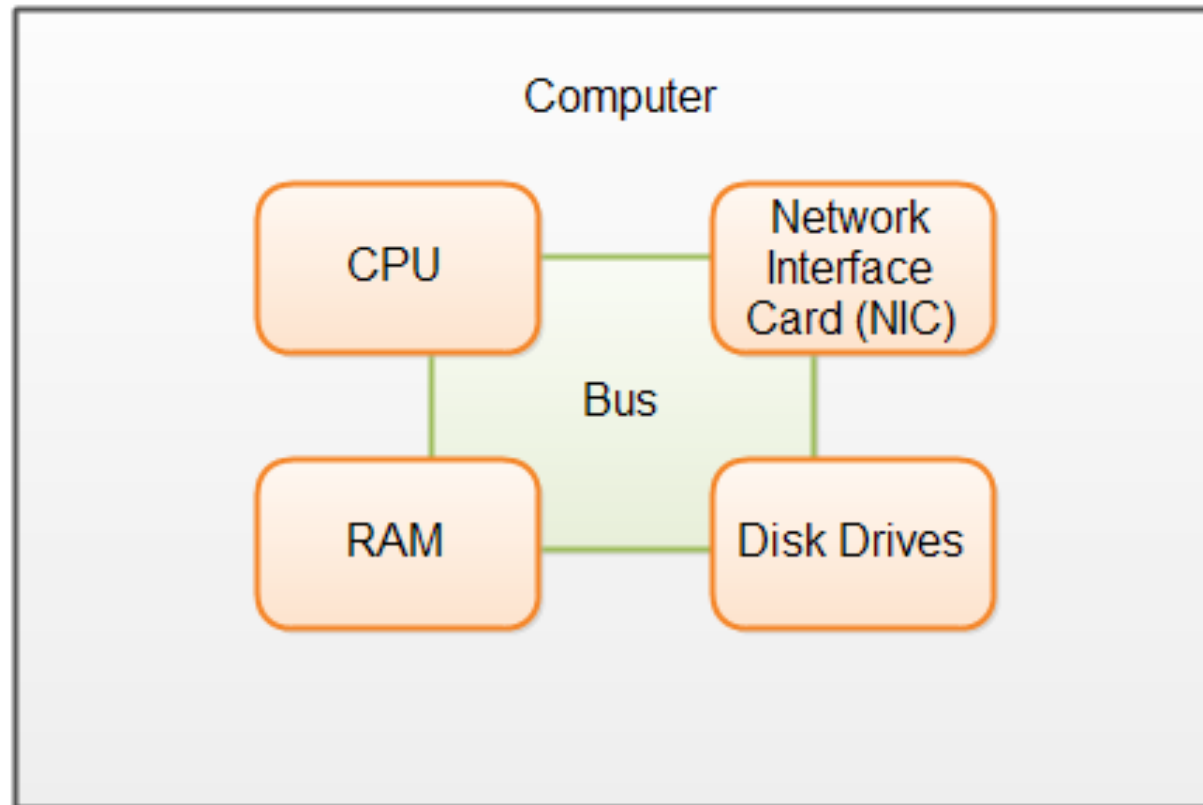
高性能服务器架构设计和调优

千石

Agenda

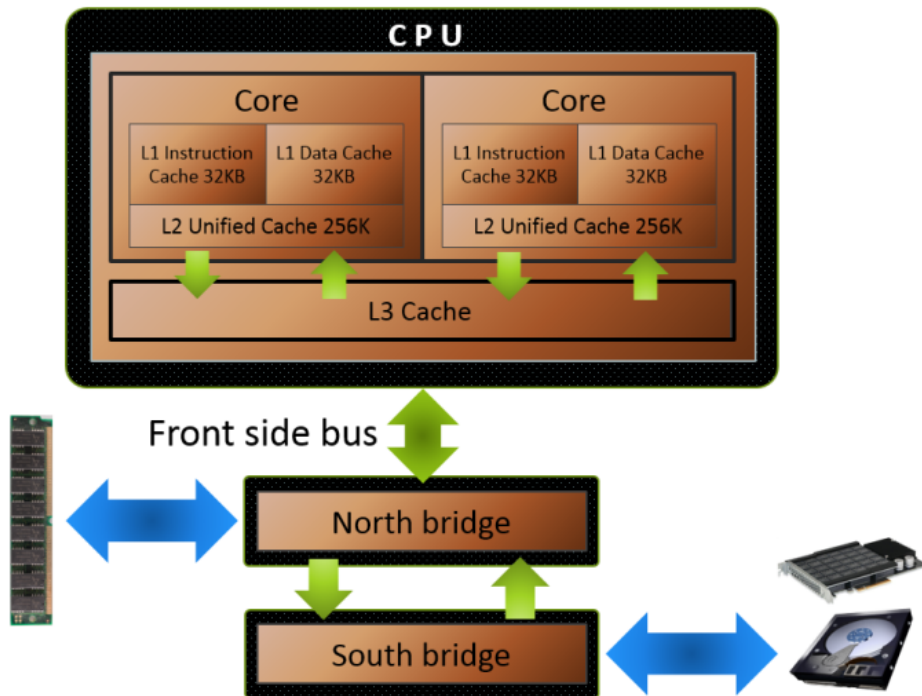


Device

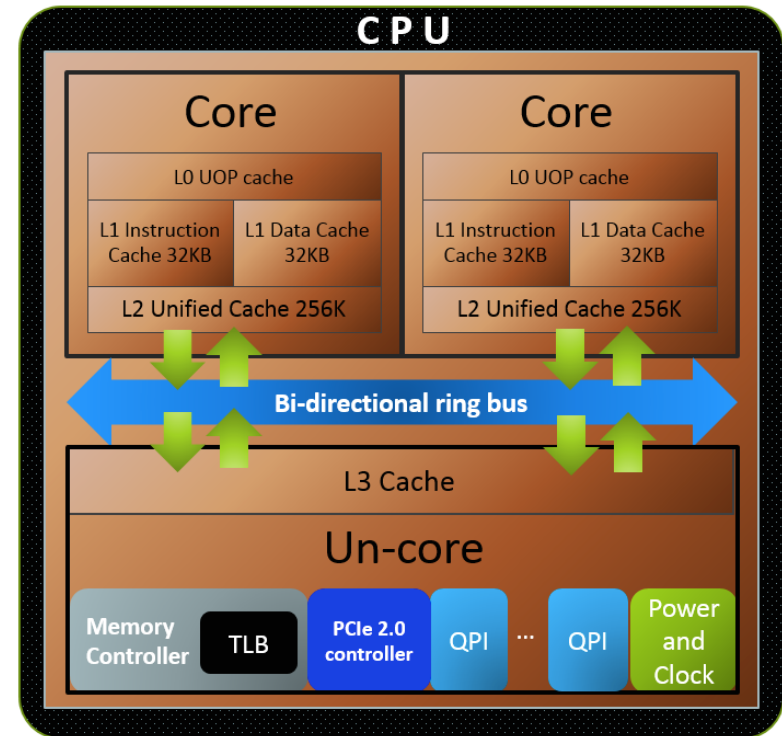


CPU

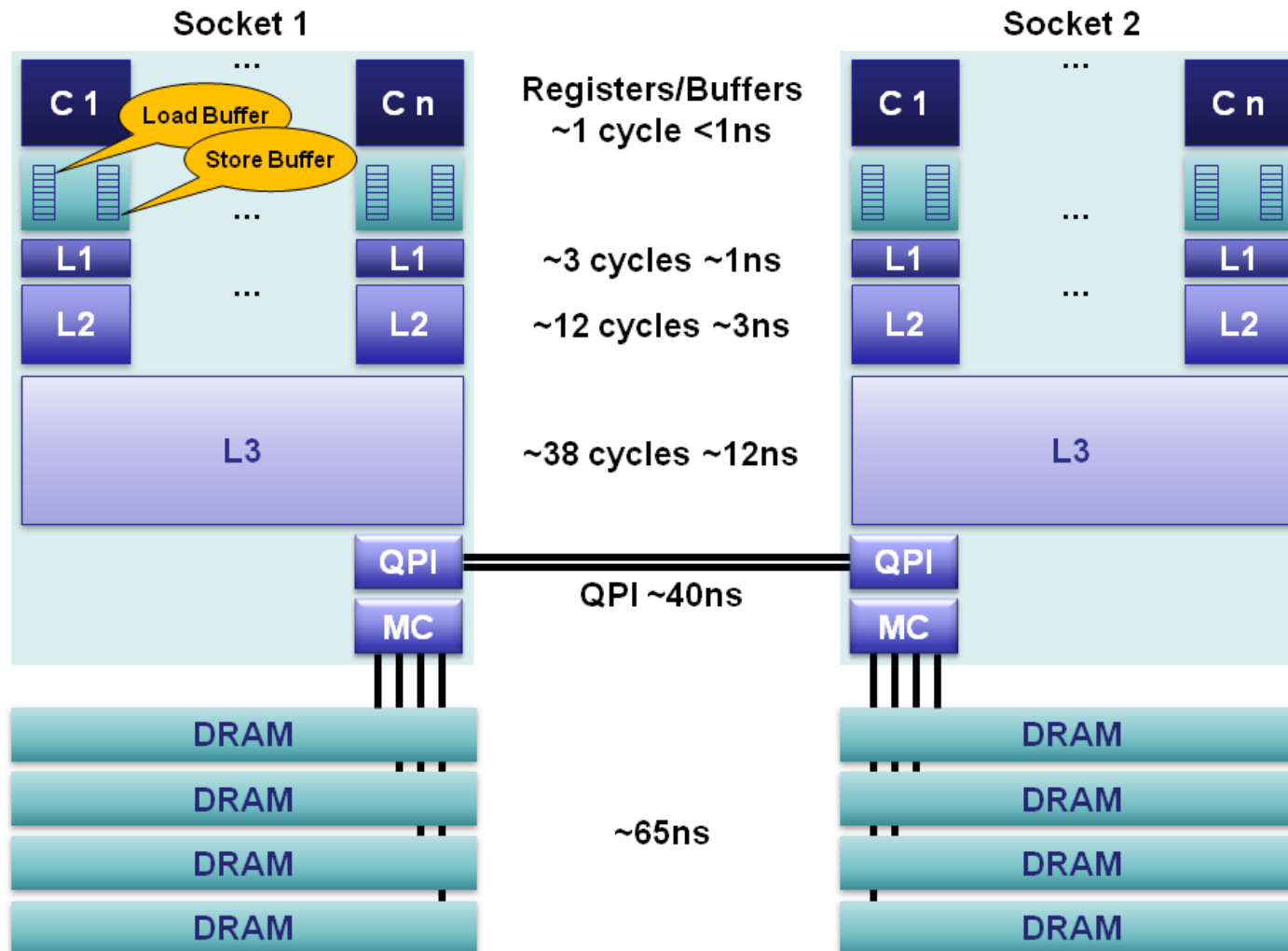
Intel Core -2



Intel Sandybridge



CPU Cache

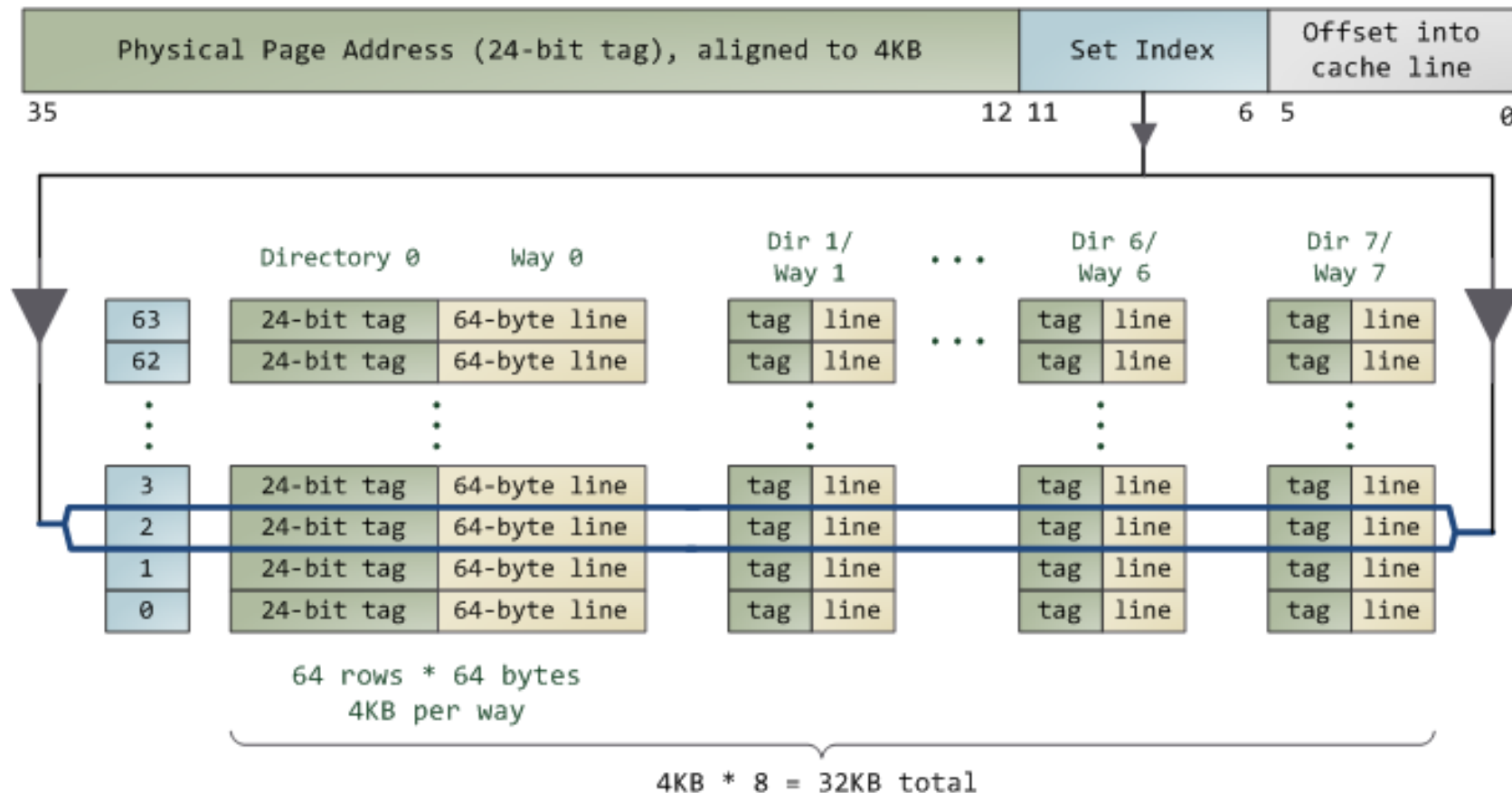


CPU Cache line

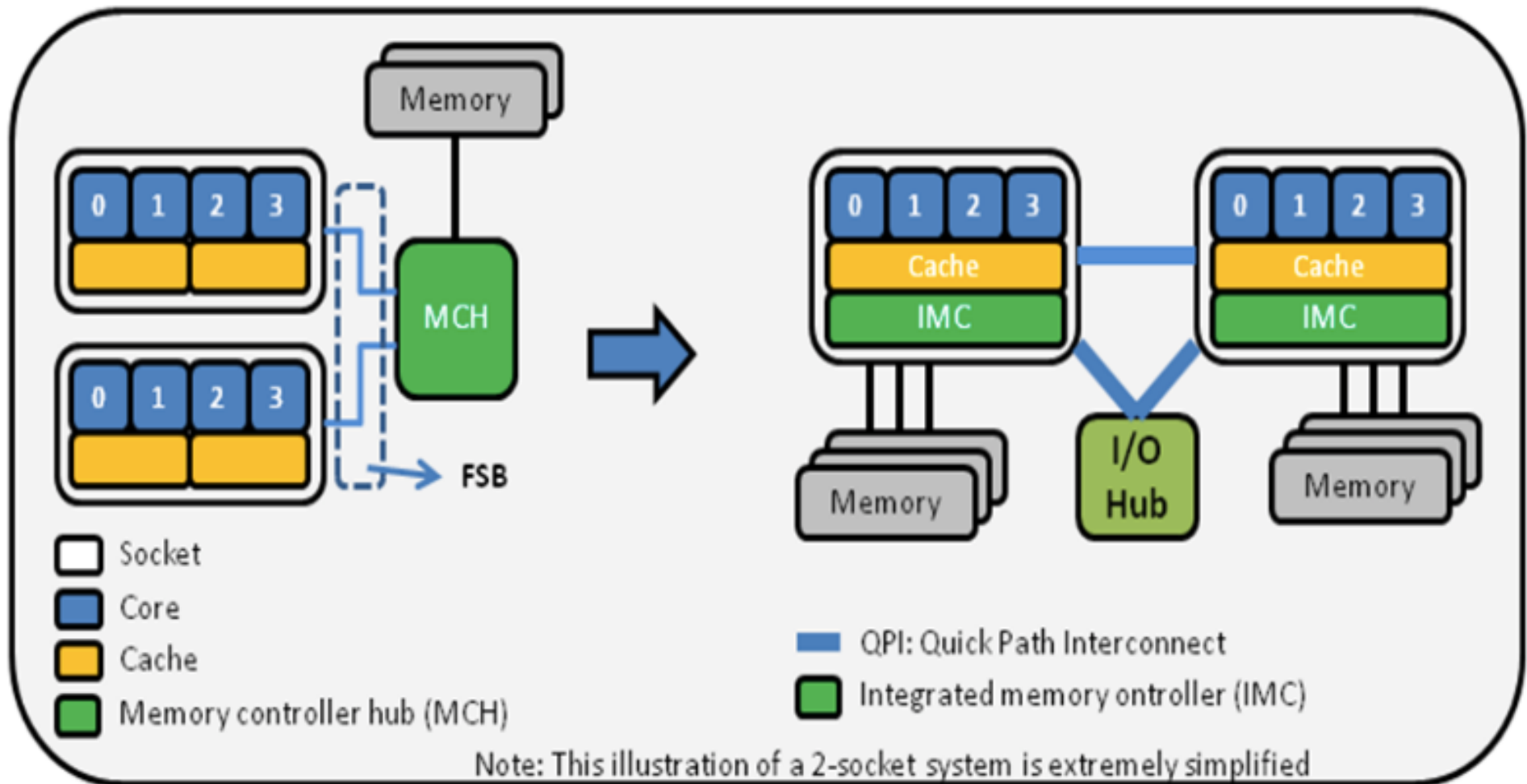
L1 Cache – 32KB, 8-way set associative, 64-byte cache lines

1. Pick cache set (row) by index

36-bit memory location as interpreted by the L1 cache:

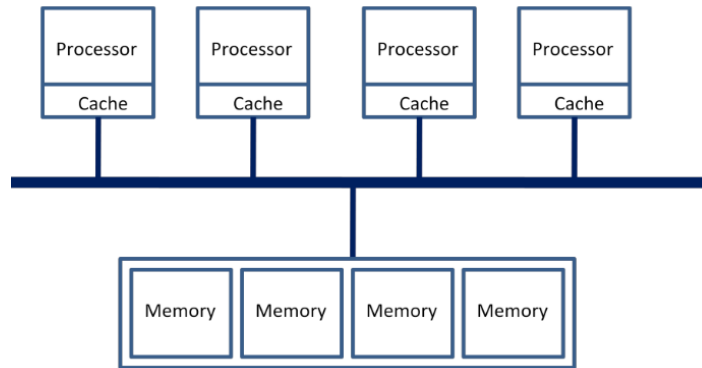


Memory Controller

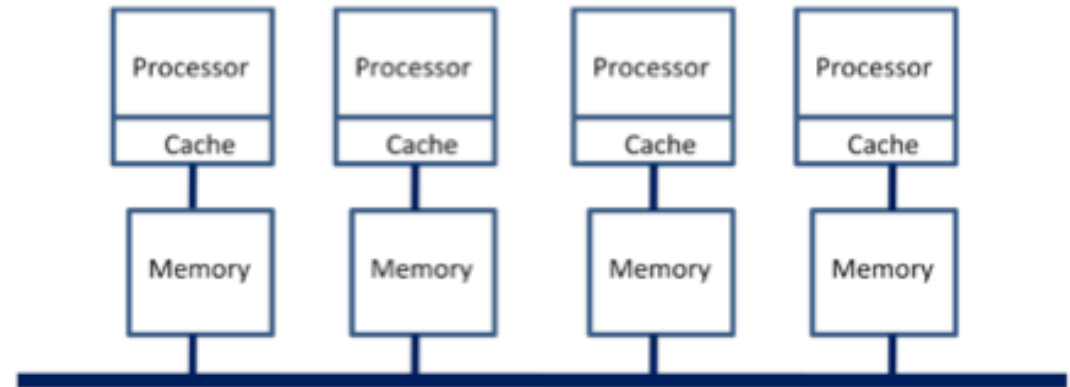


UMA/NUMA

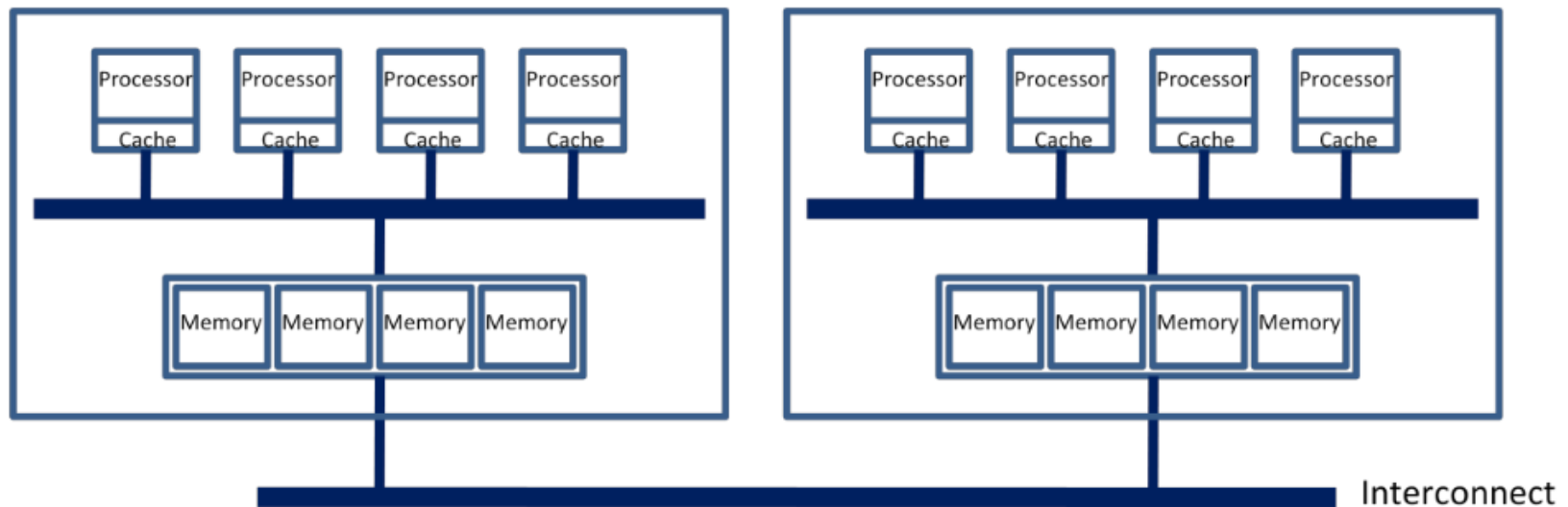
UMA



NUMA

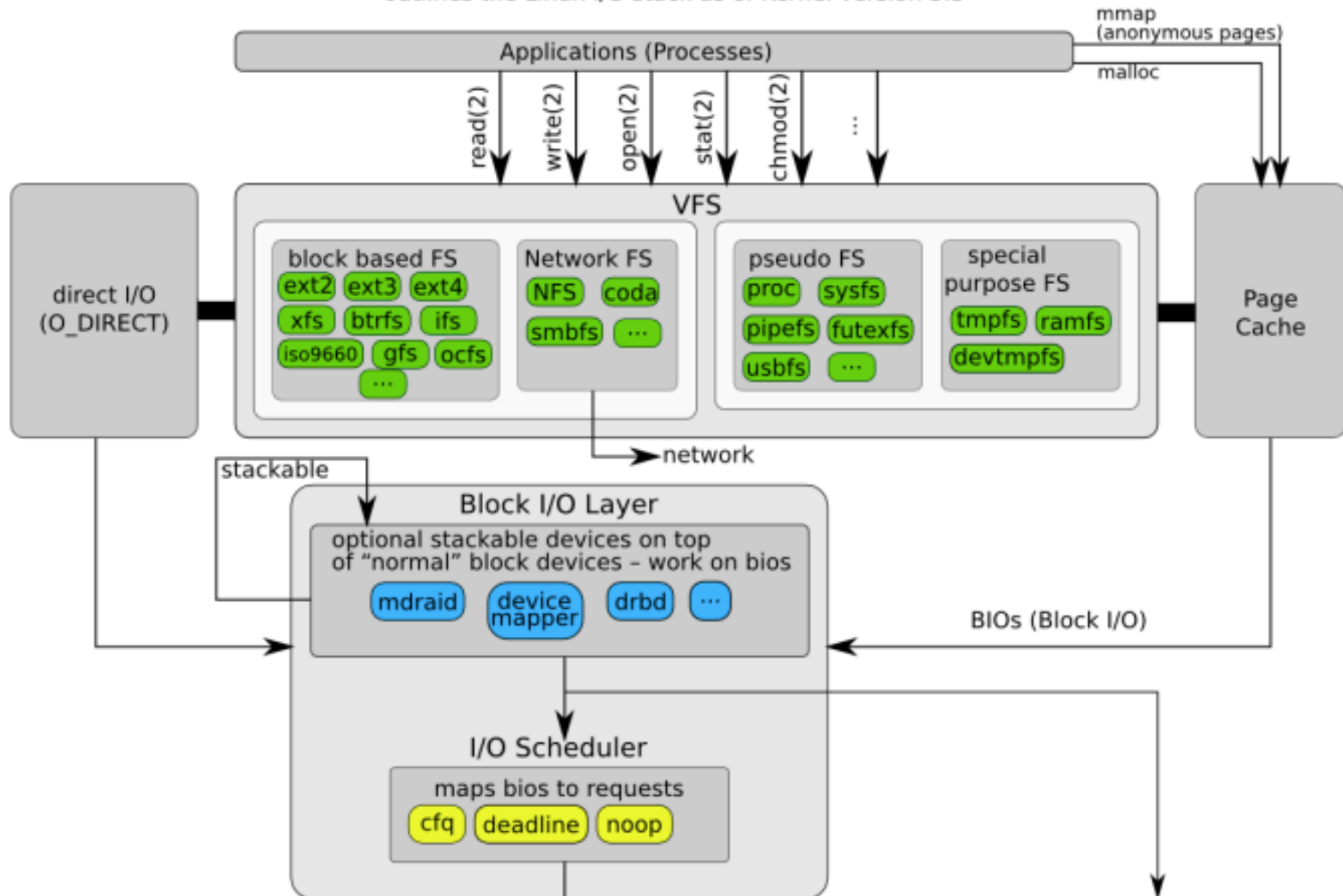


Mix architectures

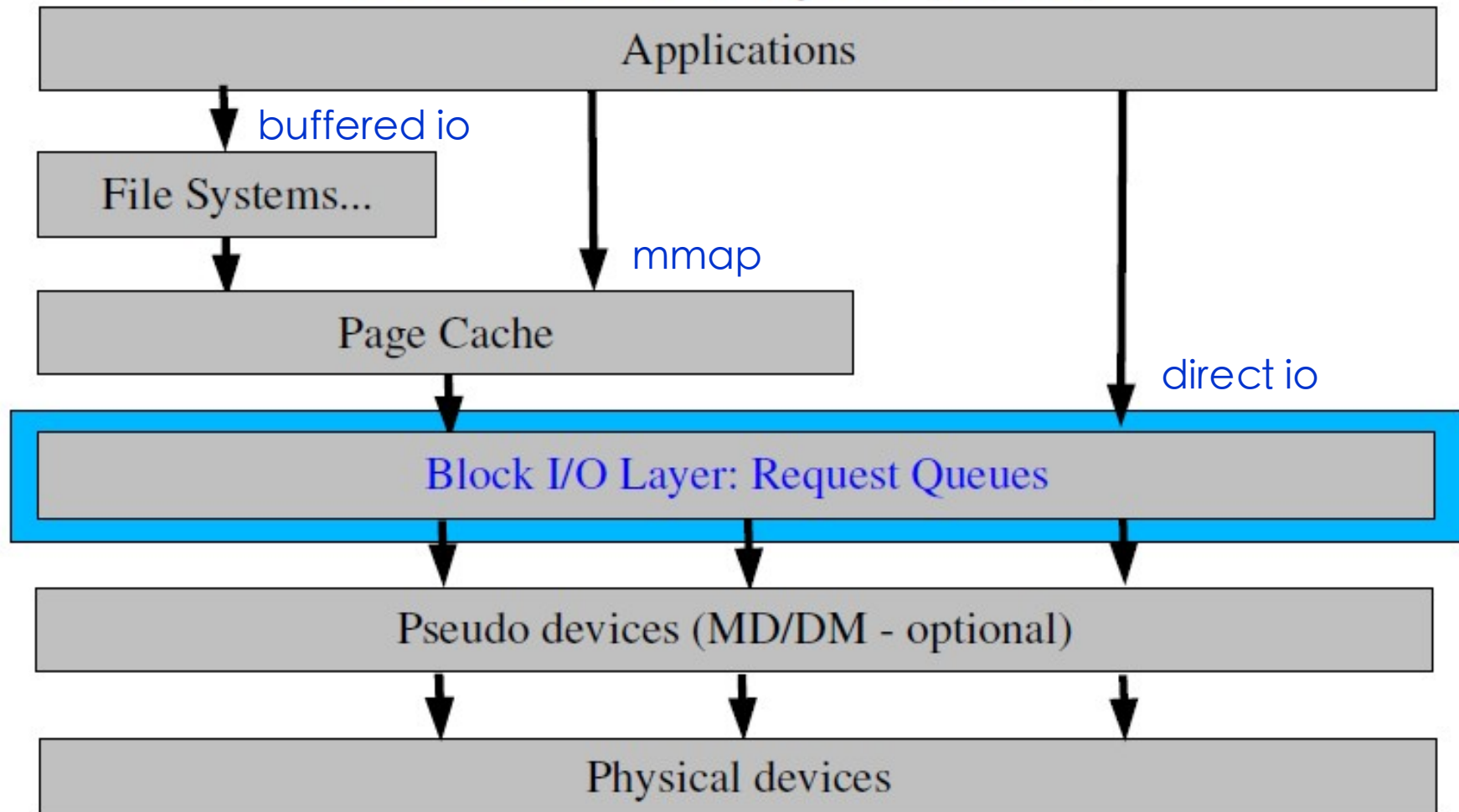


The Linux I/O Stack Diagram

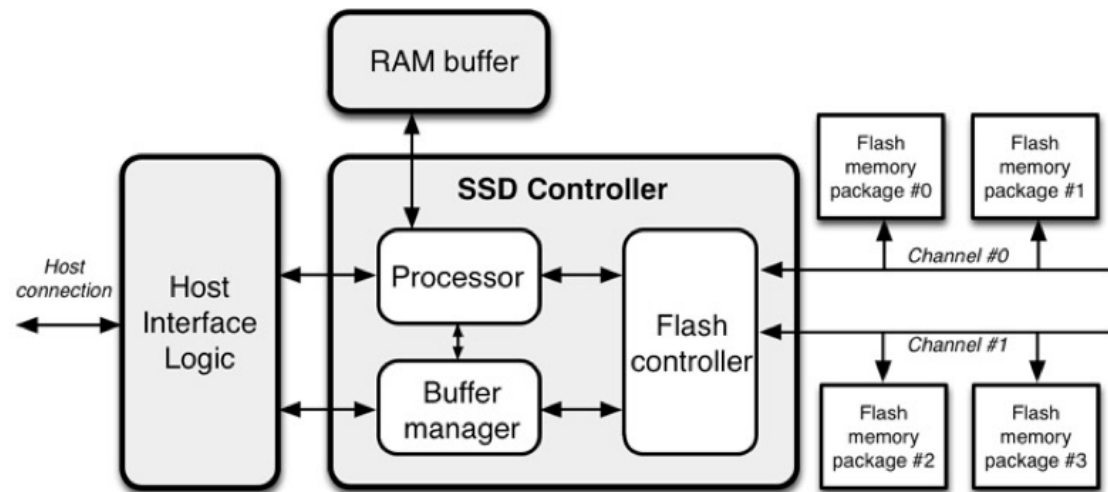
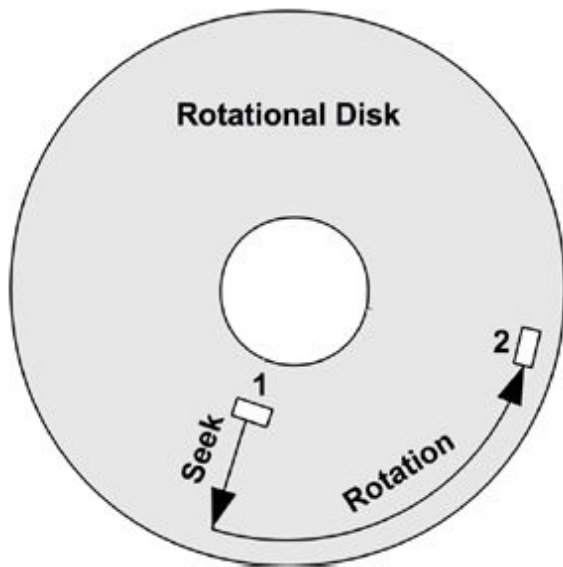
version 1.0, 2012-06-20
outlines the Linux I/O stack as of Kernel version 3.3



Block I/O Layer (simplified)



硬盘到SSD



Writing data to a solid-state drive

1. Initial configuration

Block 1000 (data)

PPN	data
0	x
1	y
2	z
3	

Block 2000 (free)

PPN	data
0	
1	
2	
3	

- Initially, block 2000 is free and block 1000 has three used pages at PPN = 0, 1, and 2 (Physical Page Number), and one free page at PPN = 3.

2. Writing a page

Block 1000 (data)

PPN	data
0	x
1	y
2	z
3	x'

Block 2000 (free)

PPN	data
0	
1	
2	
3	

- The data in block 1000 at PPN = 0 gets updated and becomes x'.
- Since pages cannot be overwritten, the page that contains x becomes stale (PPN = 0), and the new version of the data is stored in a free page, at PPN = 3.

3. Erasing a block (garbage collection)

Block 1000 (free)

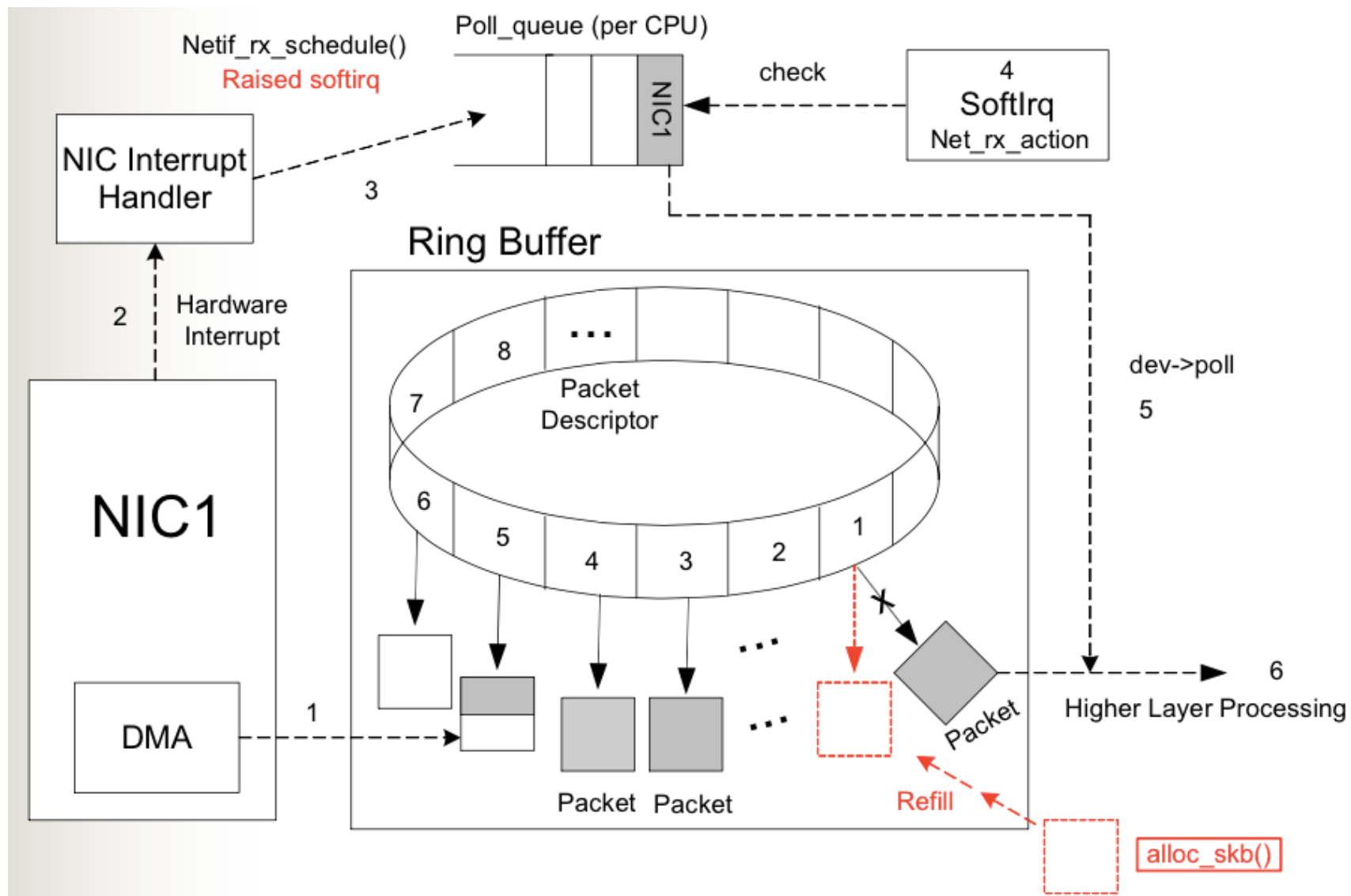
PPN	data
0	
1	
2	
3	

Block 2000 (data)

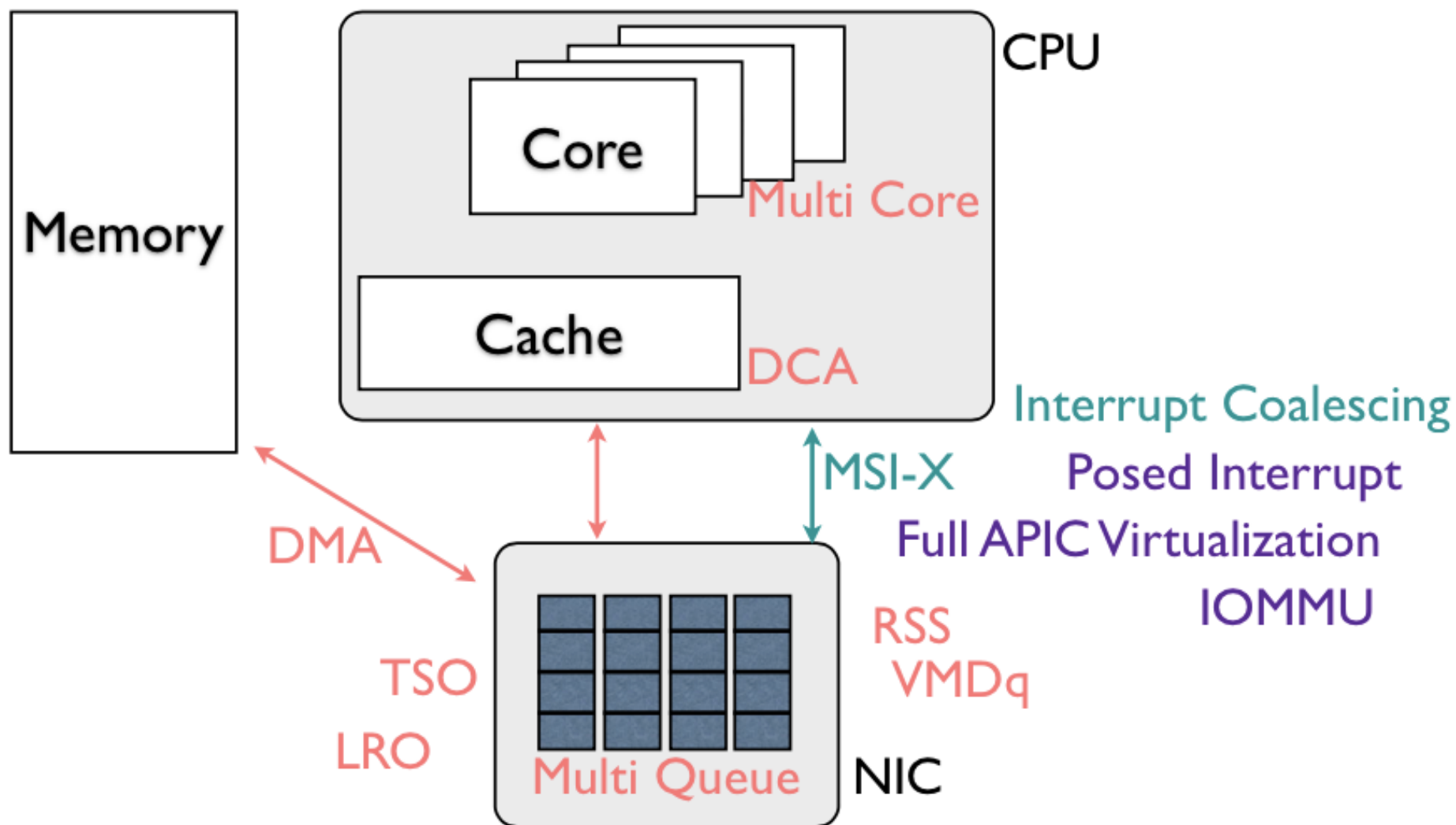
PPN	data
1	y
2	z
3	x'

- The garbage collection process copies all the valid pages from the data block 1000 into the free block 2000, leaving behind the stale pages.
- Block 1000 is erased, which makes it ready to receive new write operations. Blocks can only be erased a limited number of times (P/E cycles) until they wear off and become unusable.

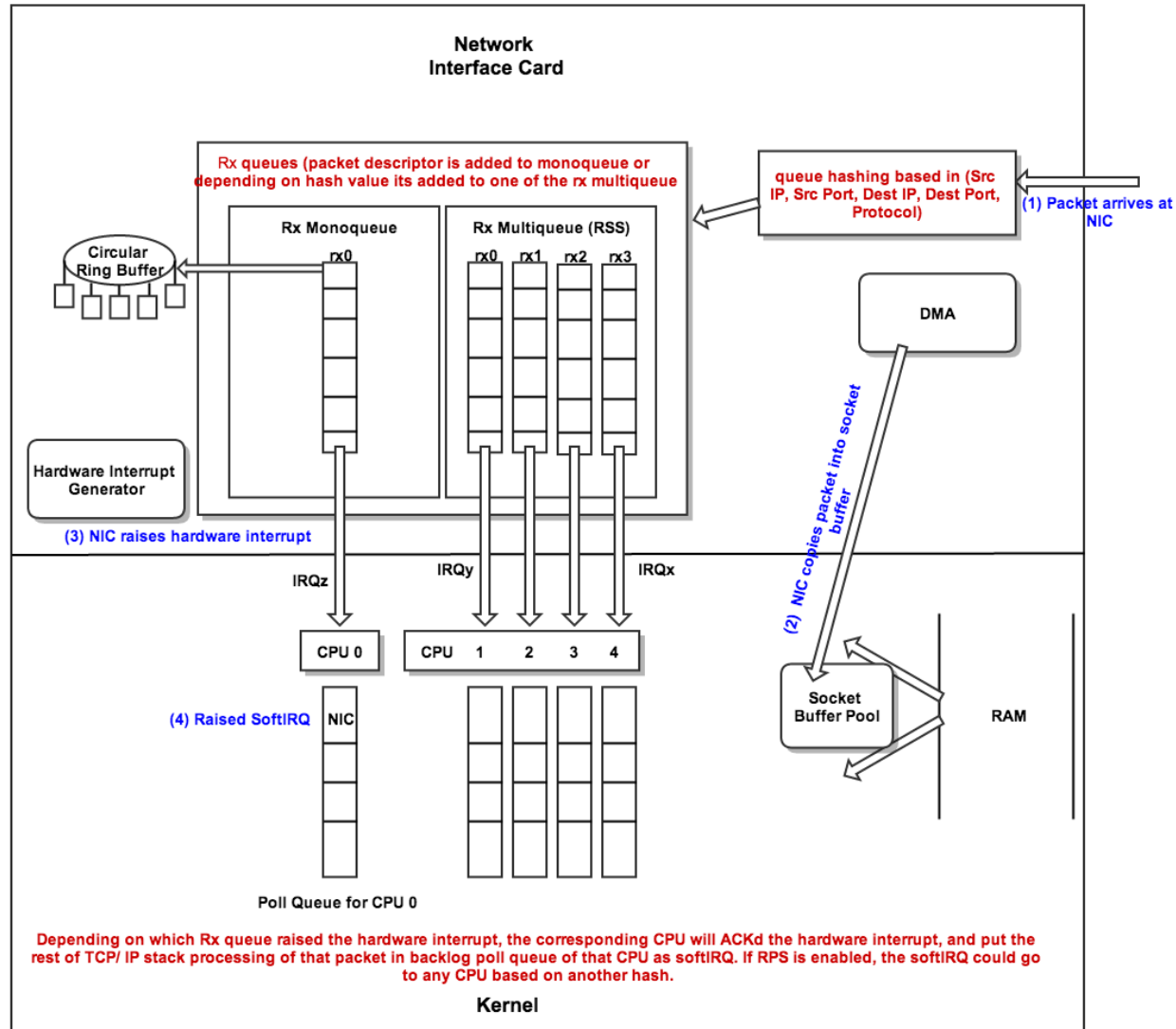
网卡收包



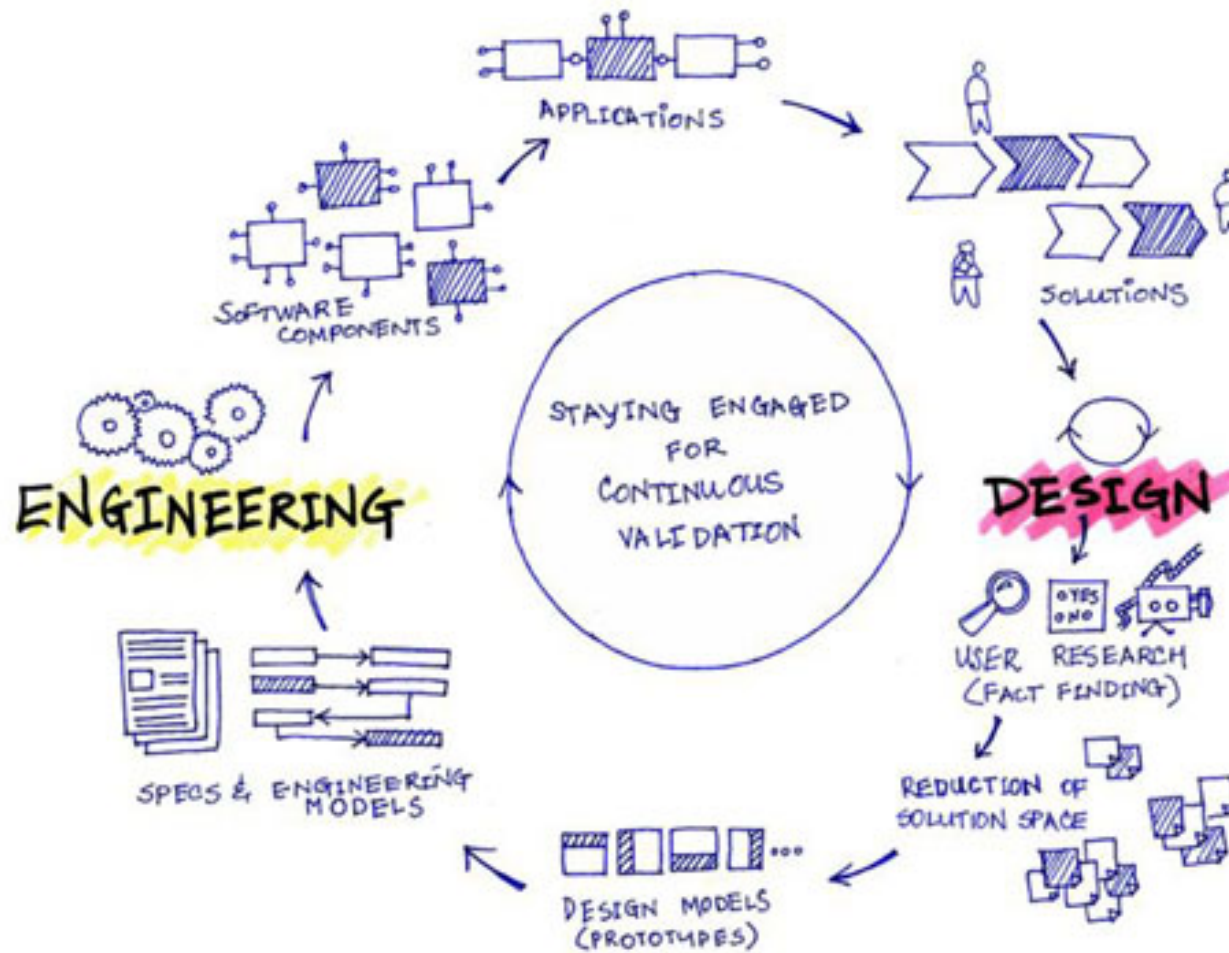
网卡多队列



网卡多队列绑定CPU



Design

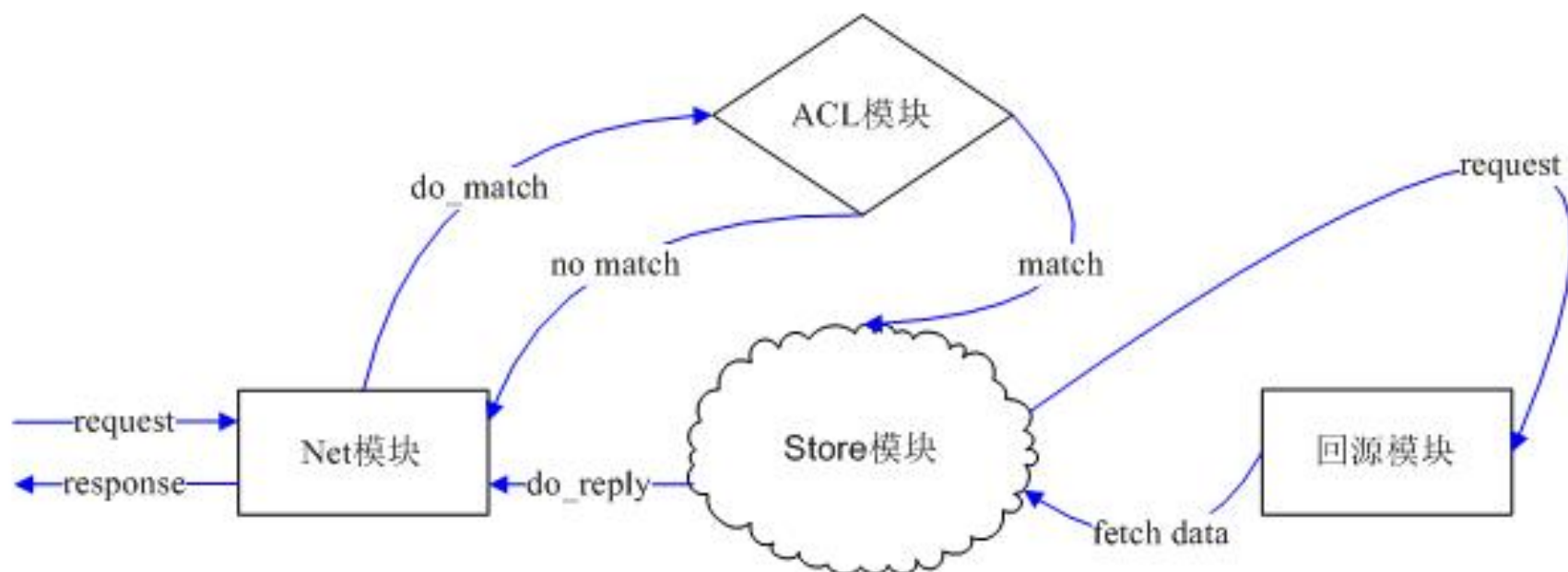


软件性能设计原则：不能成为硬件瓶颈

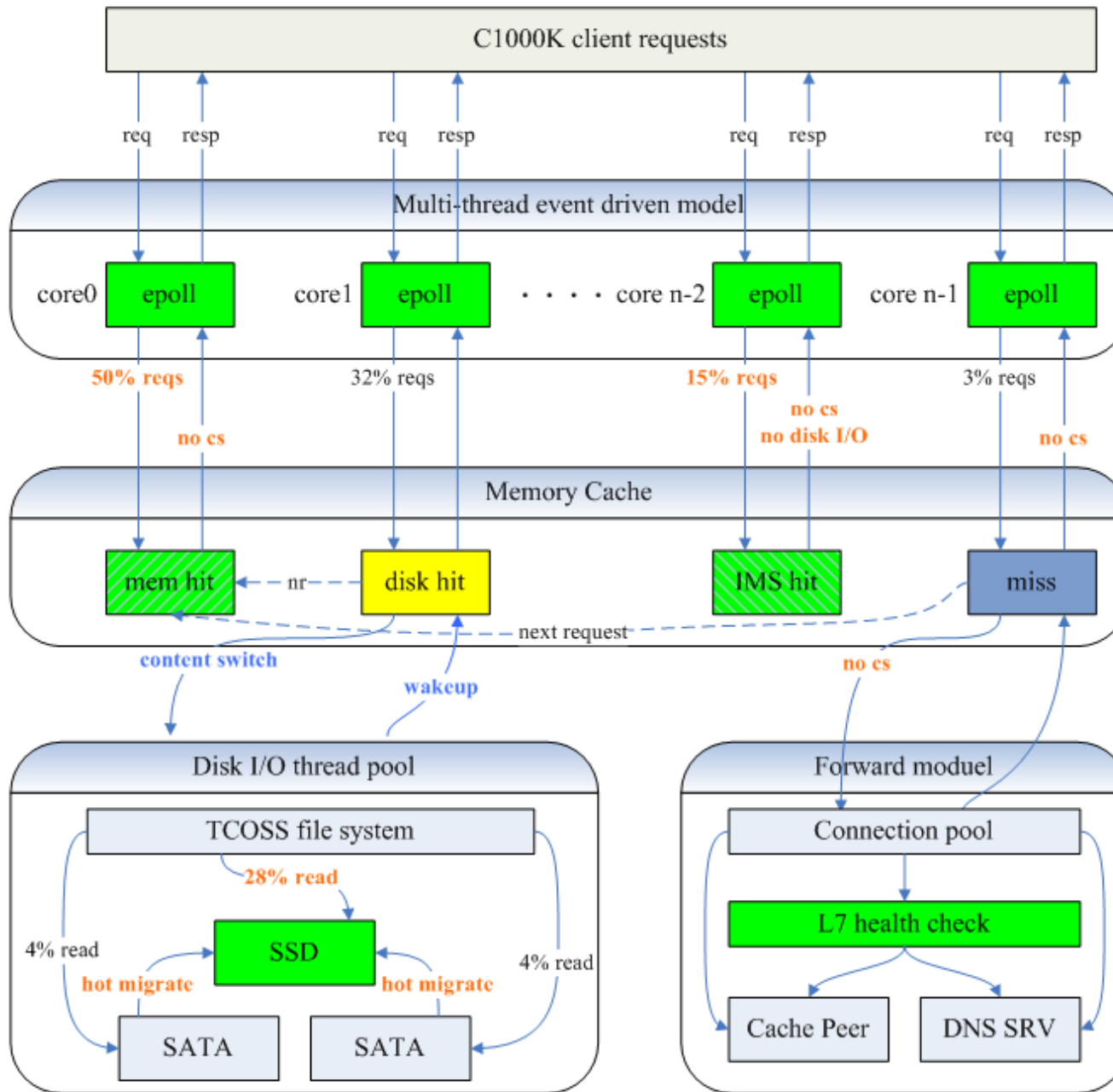
- 均衡使用CPU多核处理能力
- 高效合理地使用和控制内存
- 最大化磁盘IOPS和吞吐，异步化处理
- 小包跑满万兆网卡，中断平衡

CDN Cache系统模型

- Net模块： 支持大并发，跑满万兆网卡
- ACL模块： 高效匹配，减少CPU消耗
- Store模块： 提高命中率，高效利用磁盘IOPS
- 回源模块： L7-check，长连接保持



Swift -- High Performance Web Cache Architecture

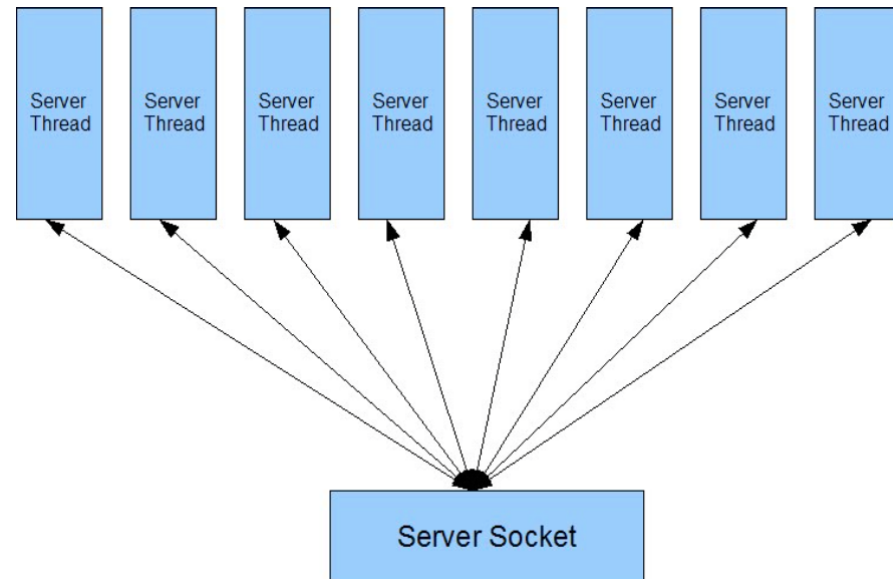


Net 模块

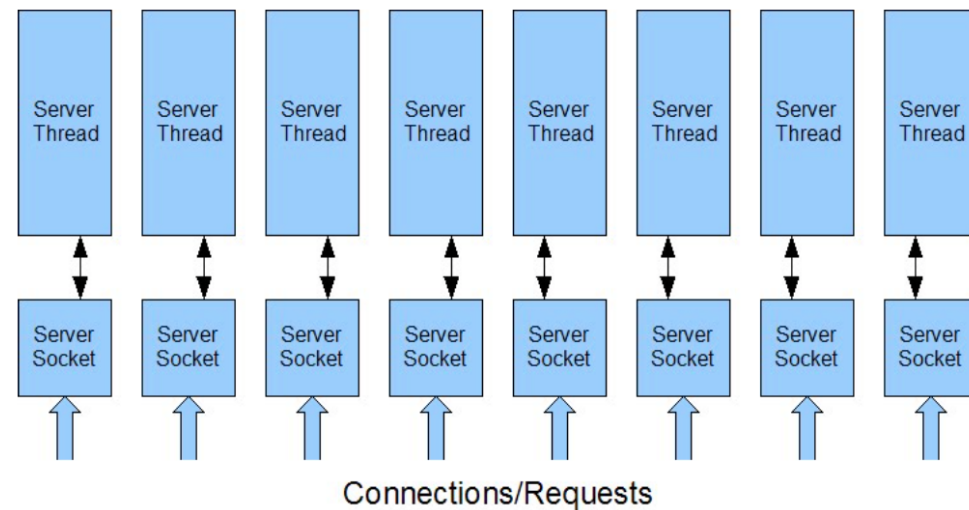
- I/O模型
 - epoll + O_NONBLOCK
- TCP选项
 - TCP_DEFER_ACCEPT / TCP_SYNCNT
 - TCP_CORK / TCP_NODELAY / TCP_QUICKACK
- 收发包的方式
 - RSS, SMP_AFFINITY
 - SO_REUSEPORT

SO_REUSEPORT

- listen 单个fd存在 accept 竞争问题

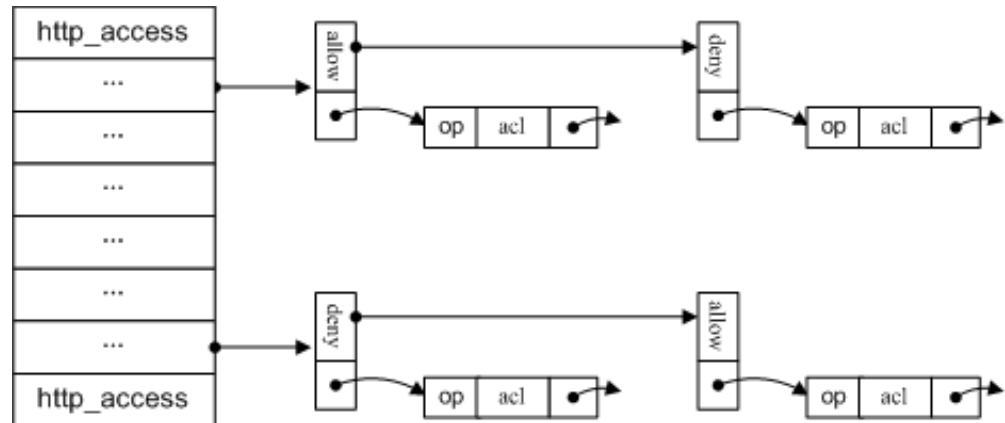


- SO_REUSEPORT listen 多个fd

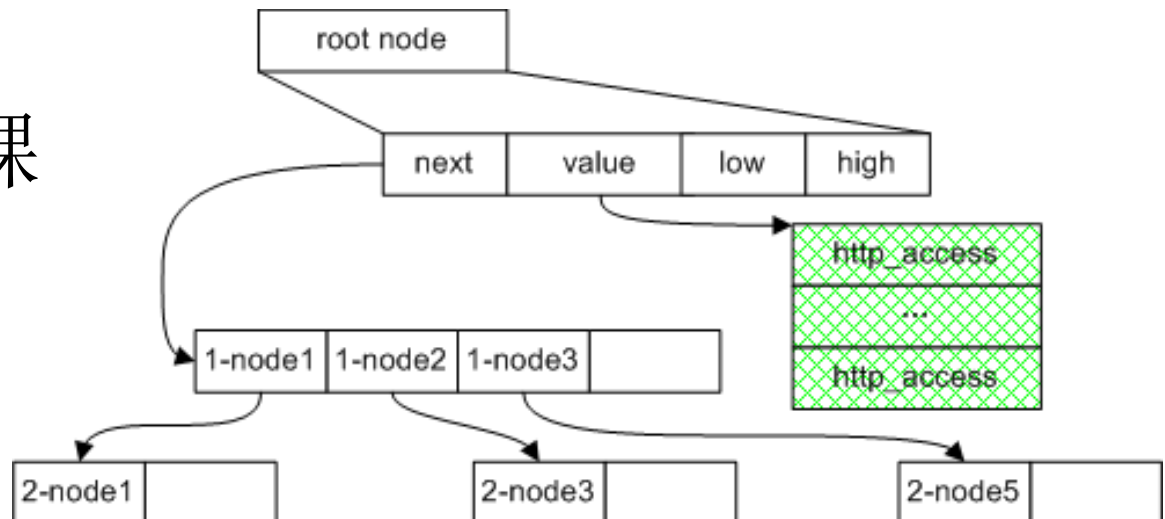


ACL优化

- squid根据http_access配置的顺序依次比较



- Swift将ACL按domain建立一棵Trie树

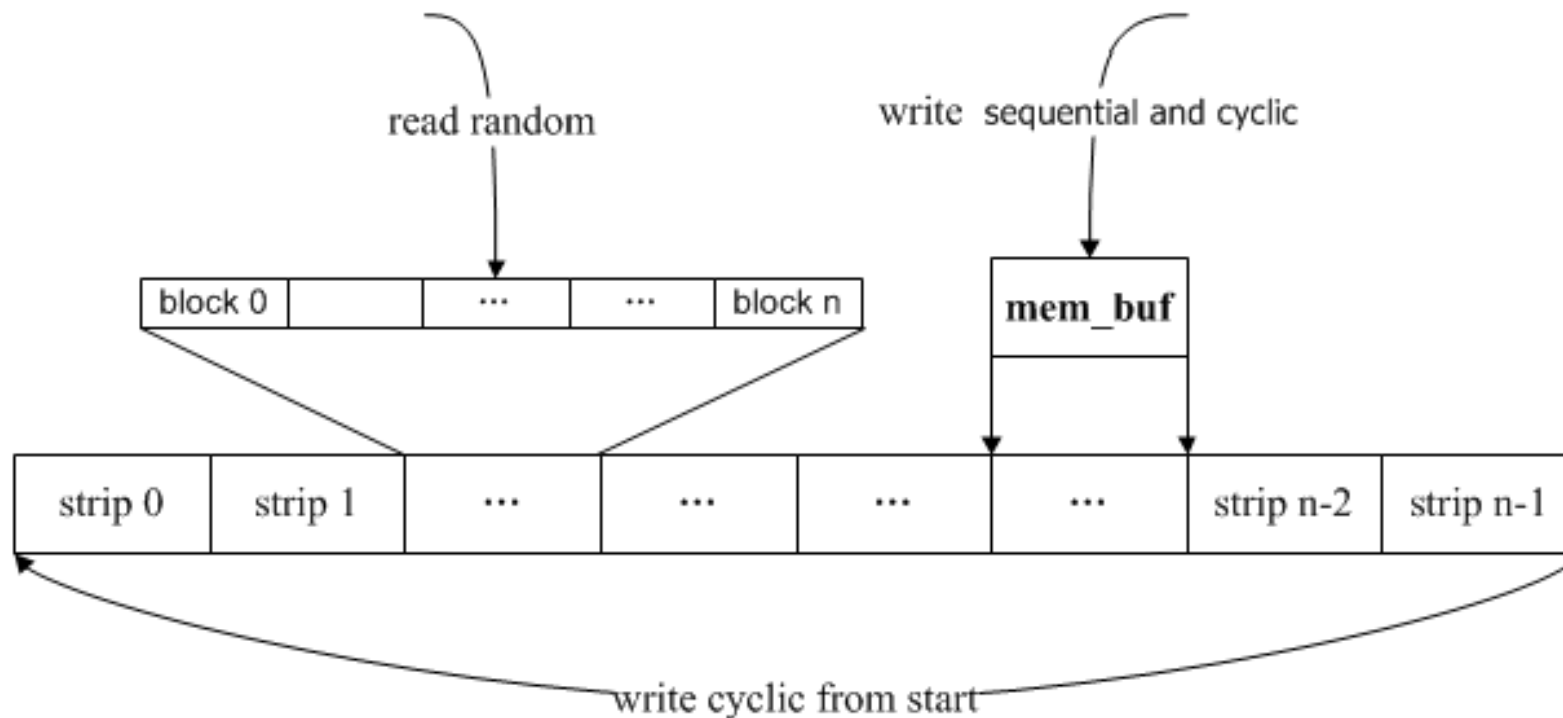


完美hash处理http header

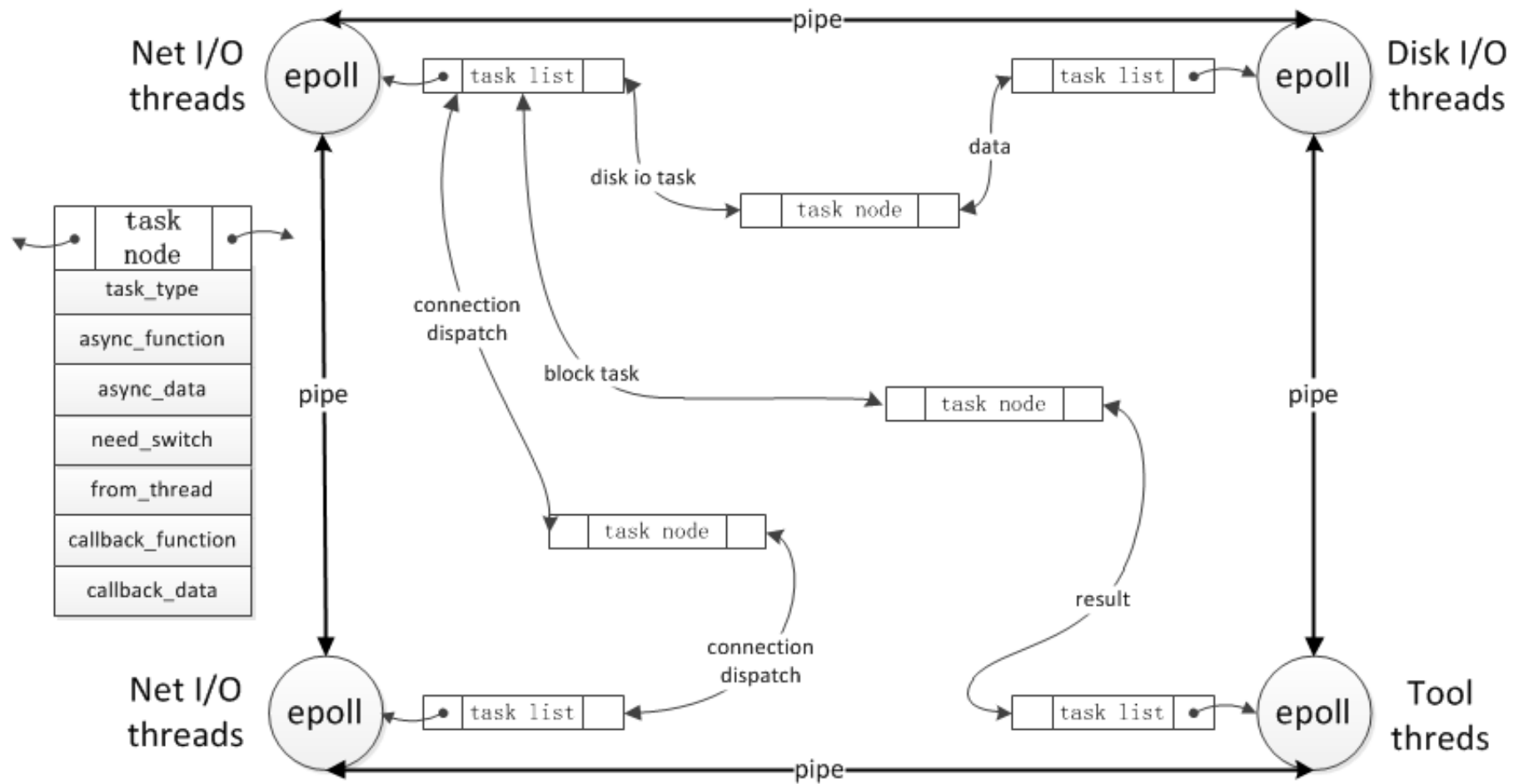
- 完美hash (Perfect Hash Function)
 - Hash table: key = value
 - PHF将key集合没有冲突地映射到一组整数
 - 查找key操作转换为索引整数表
- PHF场景&作用
 - 适合在key集合确定或不经常更新的情况
 - 主要作用是提高hash查找的速度

Store模块

- DIRECT IO写裸盘、绕过FS、不使用page cache
- 顺序写/随机读，stripe 8MB/block 512B
- mem_buf 在内存做写合并
- 省去open和close系统调用

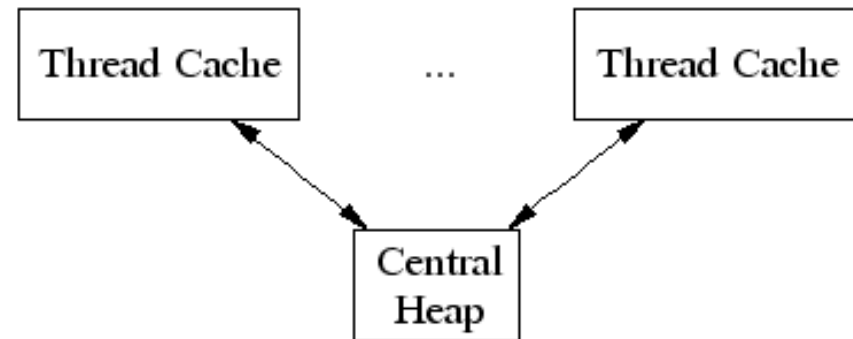


多线程任务交互模型

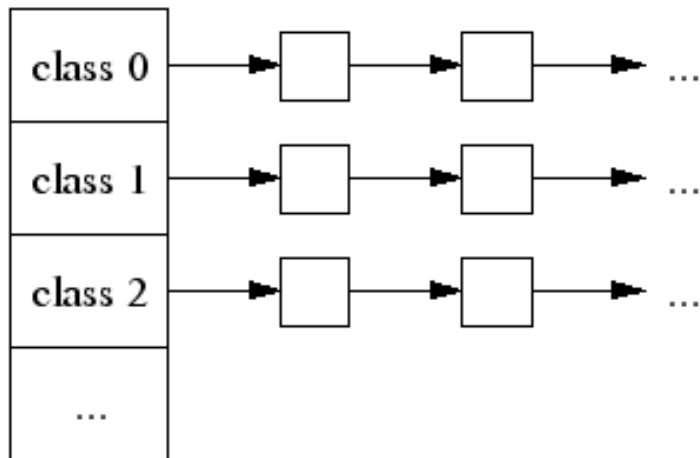


TCmalloc管理内存

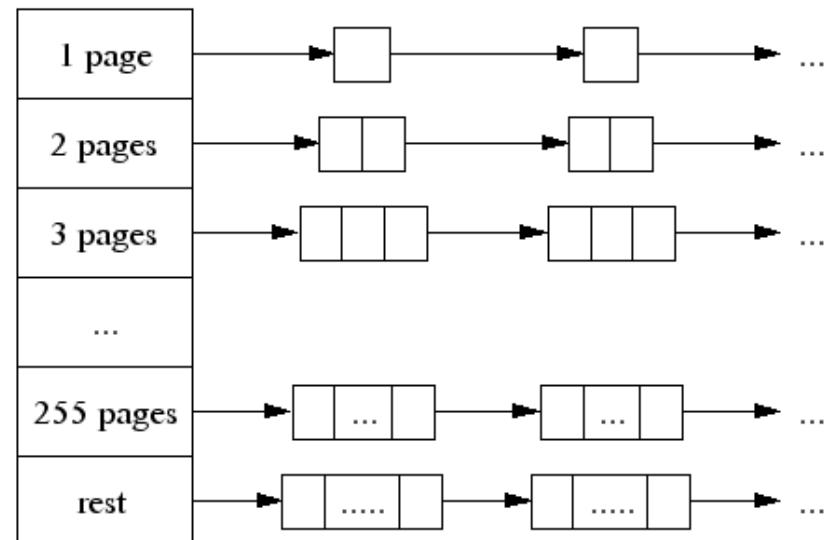
- thread free list
- central free list
- central heap



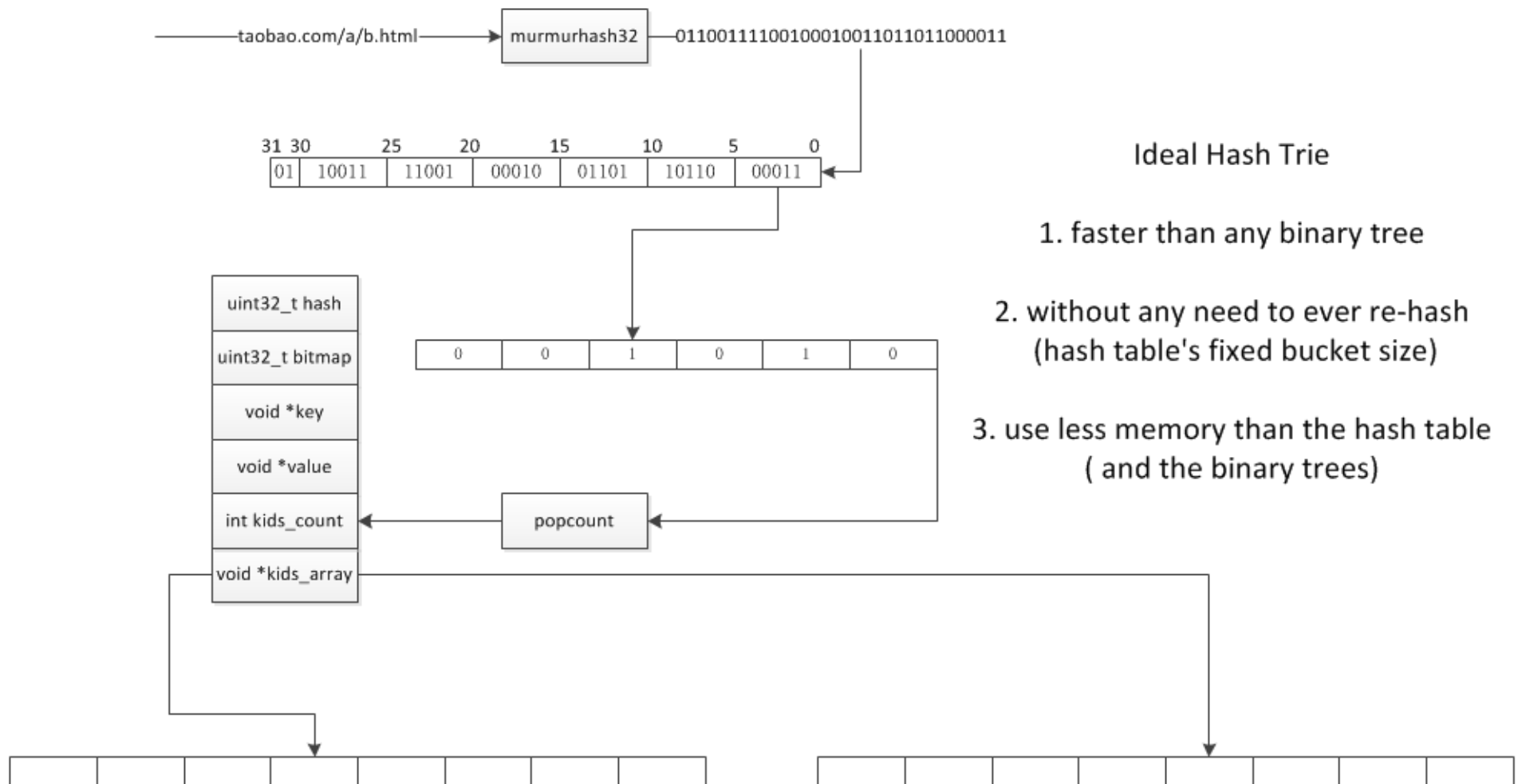
< 32KB



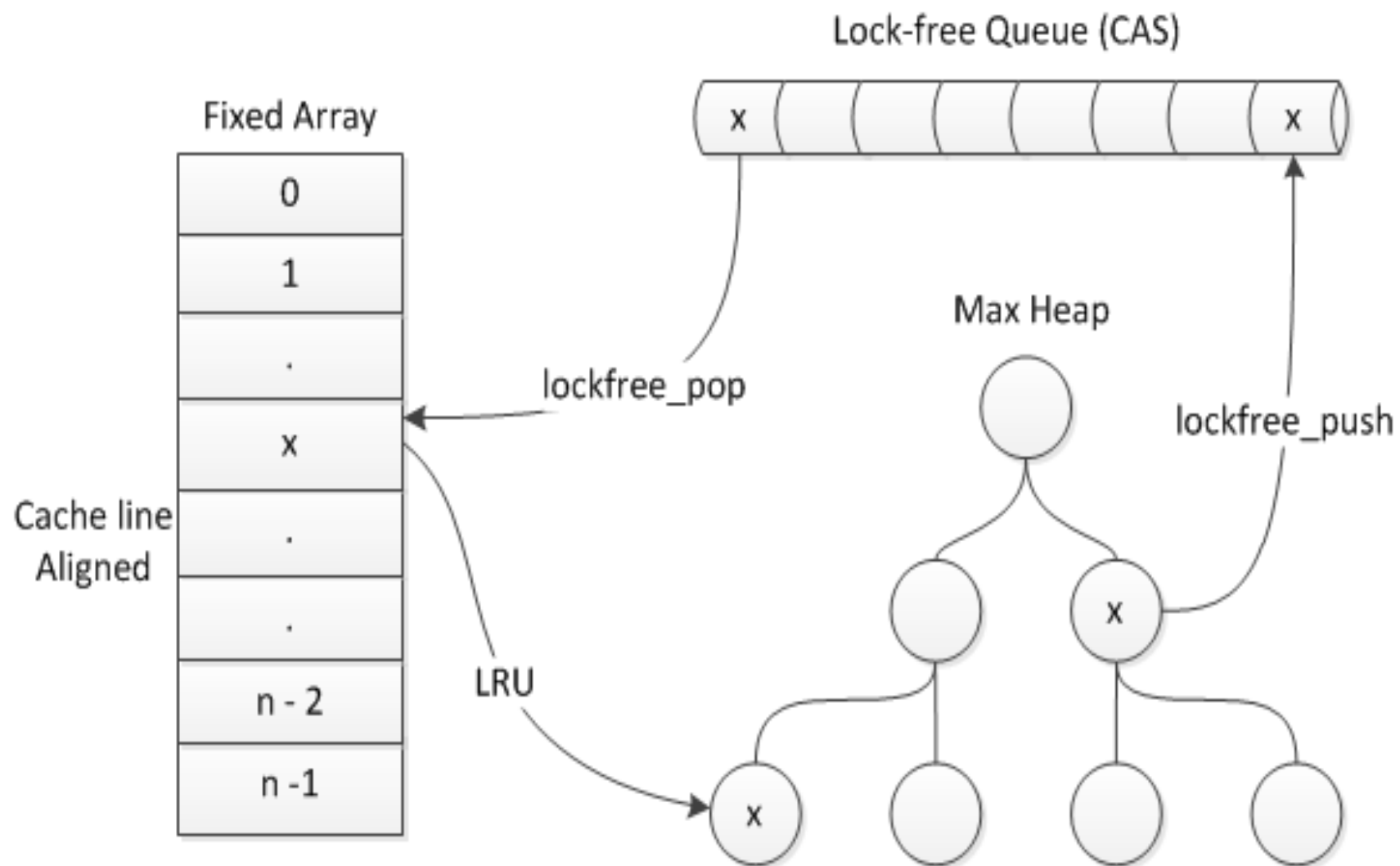
> 32KB



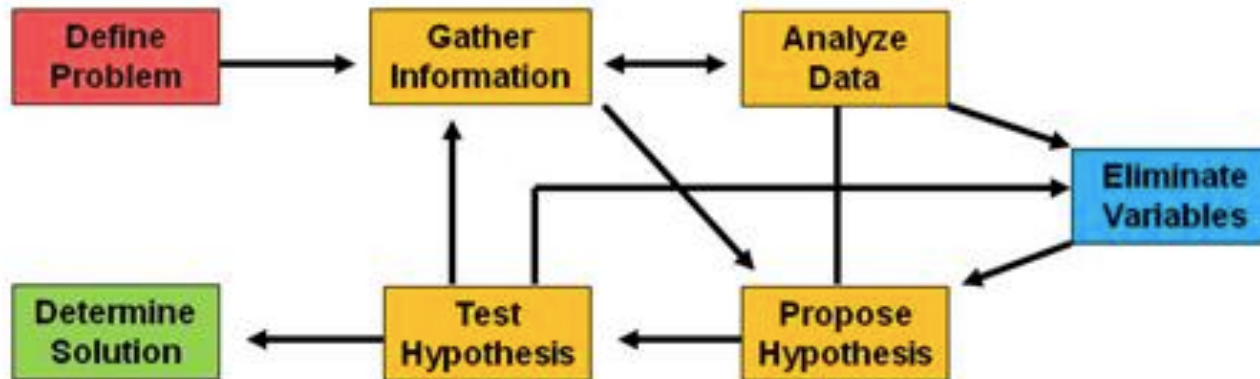
Hash Trie 替换 Hash table



无锁LRU链表



Troubleshooting



CPU分析工具

Linux	Solaris	Description
uptime	uptime	load averages
vmstat	vmstat	includes system-wide CPU averages
mpstat	mpstat	per-CPU statistics
sar	sar	historical statistics
ps	ps	process status
top	prstat	monitor per-process/thread CPU usage
pidstat	prstat	per-process/thread CPU breakdowns
time	ptime	time a command, with CPU breakdowns
DTrace, perf	DTrace	CPU profiling and tracing
perf	cpustat	CPU performance counter analysis

用 /proc/\$pid 定位问题

- # top -cbp \$pid
 - # strace -cp \$pid
 - # ps -flp \$pid
 - # pstack \$pid
-
- # cat /proc/\$pid/**wchan**
 - # cat /proc/\$pid/**status**
 - # cat /proc/\$pid/**sched**
 - # cat /proc/\$pid/**schedstat**
 - # cat /proc/\$pid/**syscall**
 - # cat /proc/\$pid/**stack**

Mem分析工具

Linux	Solaris	Description
vmstat	vmstat	virtual and physical memory statistics
sar	sar	historical statistics
slabtop	::kmastat	kernel slab allocator statistics
ps	ps	process status
top	prstat	monitor per-process memory usage
pmap	pmap	process address space statistics
DTrace	DTrace	allocation tracing

/proc/meminfo 项的关系

- $\text{MemTotal} = \text{LowTotal} + \text{HighTotal}$
- $\text{MemFree} = \text{LowFree} + \text{HighFree}$
- $\text{Slab} = \text{SReclaimable} + \text{SUnreclaimable}$
- $\text{Active} = \text{Active}(\text{anon}) + \text{Active}(\text{file})$
- $\text{Inactive} = \text{Inactive}(\text{anon}) + \text{Inactive}(\text{file})$
- $\text{AnonPages} + \text{?X?} = \text{Active}(\text{anon}) + \text{Inactive}(\text{anon})$
- $\text{Buffers} + \text{Cached} = \text{Active}(\text{file}) + \text{Inactive}(\text{file}) + \text{?X?}$
- $\text{AnonPages} + \text{Buffers} + \text{Cached} = \text{Active} + \text{Inactive}$
- $\text{SwapTotal} = \text{SwapFree} + \text{SwapCached}$

内存问题定位

- 内存泄露

- # env **HEAPPROFILE**=/home/qianshi/dev/swift.hprof ./swift -f swift.conf
- # pprof --pdf --base=swift.hprof.0001.heap . /swift swift.hprof.0002.heap > 1-2.pdf

- 内存写乱

- # **clang** -O1 -g -fsanitize=address -fno-omit-frame-pointer example_UseAfterFree.cc

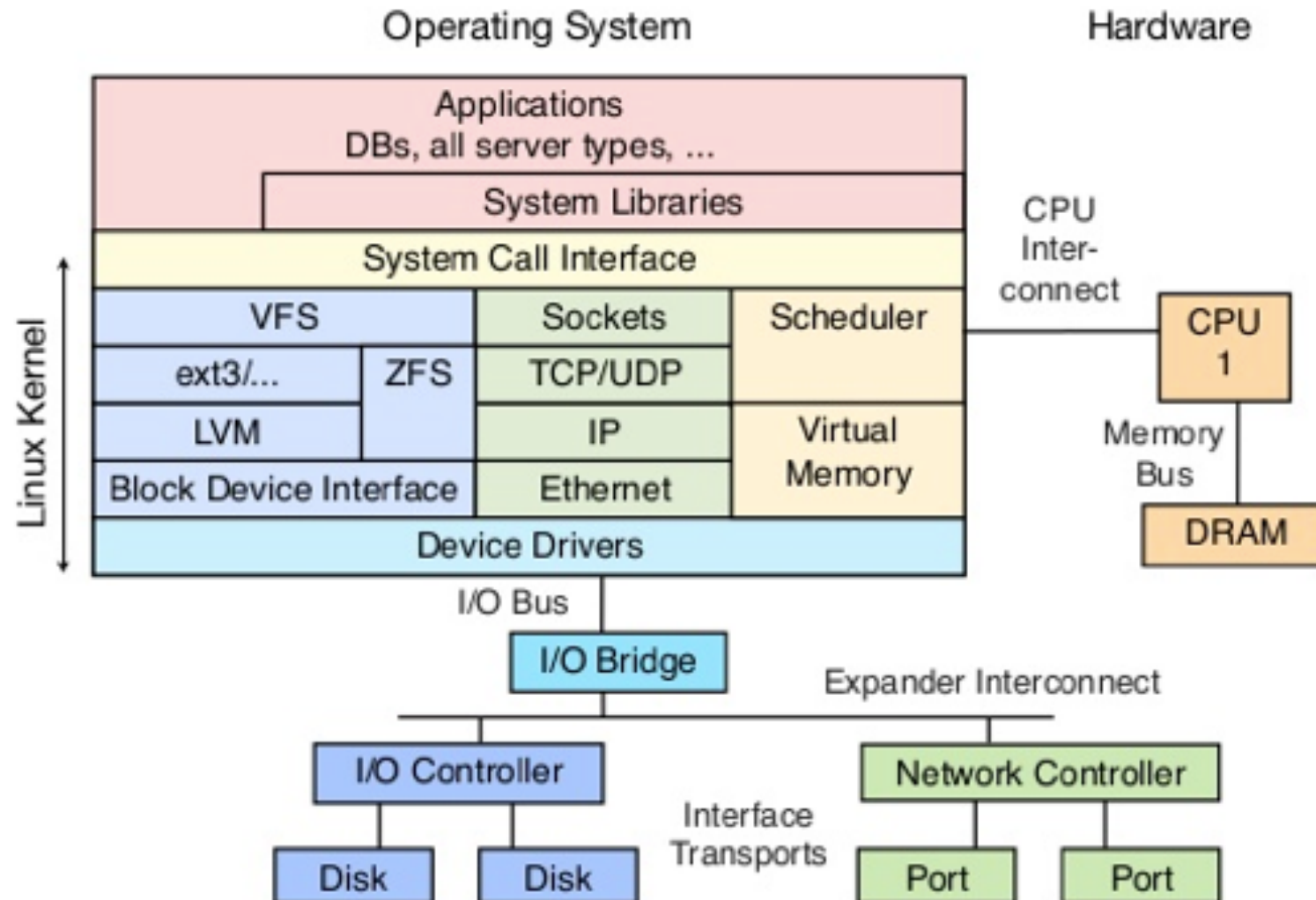
Disk分析工具

Linux	Solaris	Description
iostat	iostat	various per-disk statistics
sar	sar	historical disk statistics
pidstat, iotop	iotop	disk I/O usage by process
blktrace	iosnoop	disk I/O event tracing
DTrace	DTrace	custom static and dynamic tracing
MegaCli	MegaCli	LSI controller statistics
smartctl	smartctl	disk controller statistics

Network分析工具

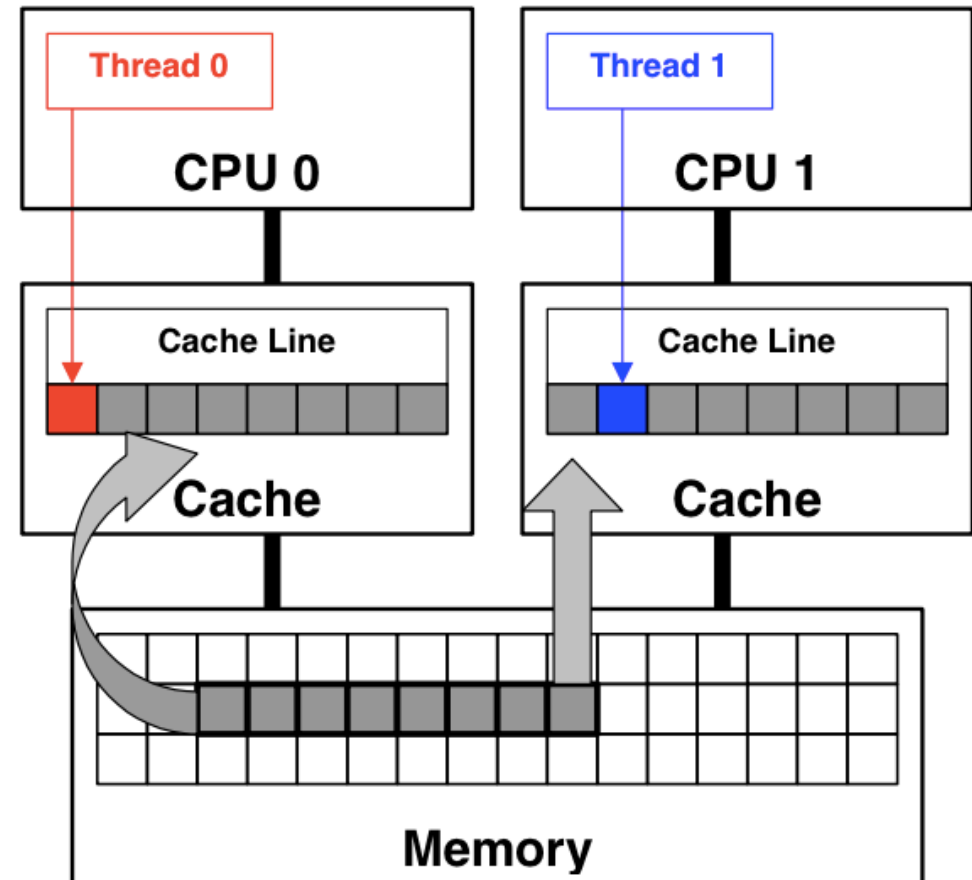
Linux	Solaris	Description
netstat	netstat	various network stack and interface statistics
sar	—	historical statistics
ifconfig	ifconfig	interface configuration
ip	dladm	network interface statistics
nicstat	nicstat	network interface throughput and utilization
ping	ping	test network connectivity
tracert	tracert	test network routes
pathchar	pathchar	determine network path characteristics
tcpdump	snoop/tcpdump	network packet sniffer
Wireshark	Wireshark	graphical network packet inspection
DTrace, perf	DTrace	TCP/IP stack tracing: connections, packets, drops, latency

Tuning



CPU tuning

- CPU亲和性
 - 提高cache命中率
 - 降低访问内存延迟
 - `taskset -c -p $pid`
- 避免false sharing
 - 编译器强制对齐
 - 填充结构体保证cache line对齐
 - 使用线程局部数据



Memory tuning

- 关掉SWAP
 - /proc/sys/vm/swappiness
 - swapoff -a
- OOM处理
 - /proc/\$pid/oom_adj
 - /proc/sys/vm/overcommit_memory
 - /proc/sys/vm/overcommit_ratio

Disk tuning

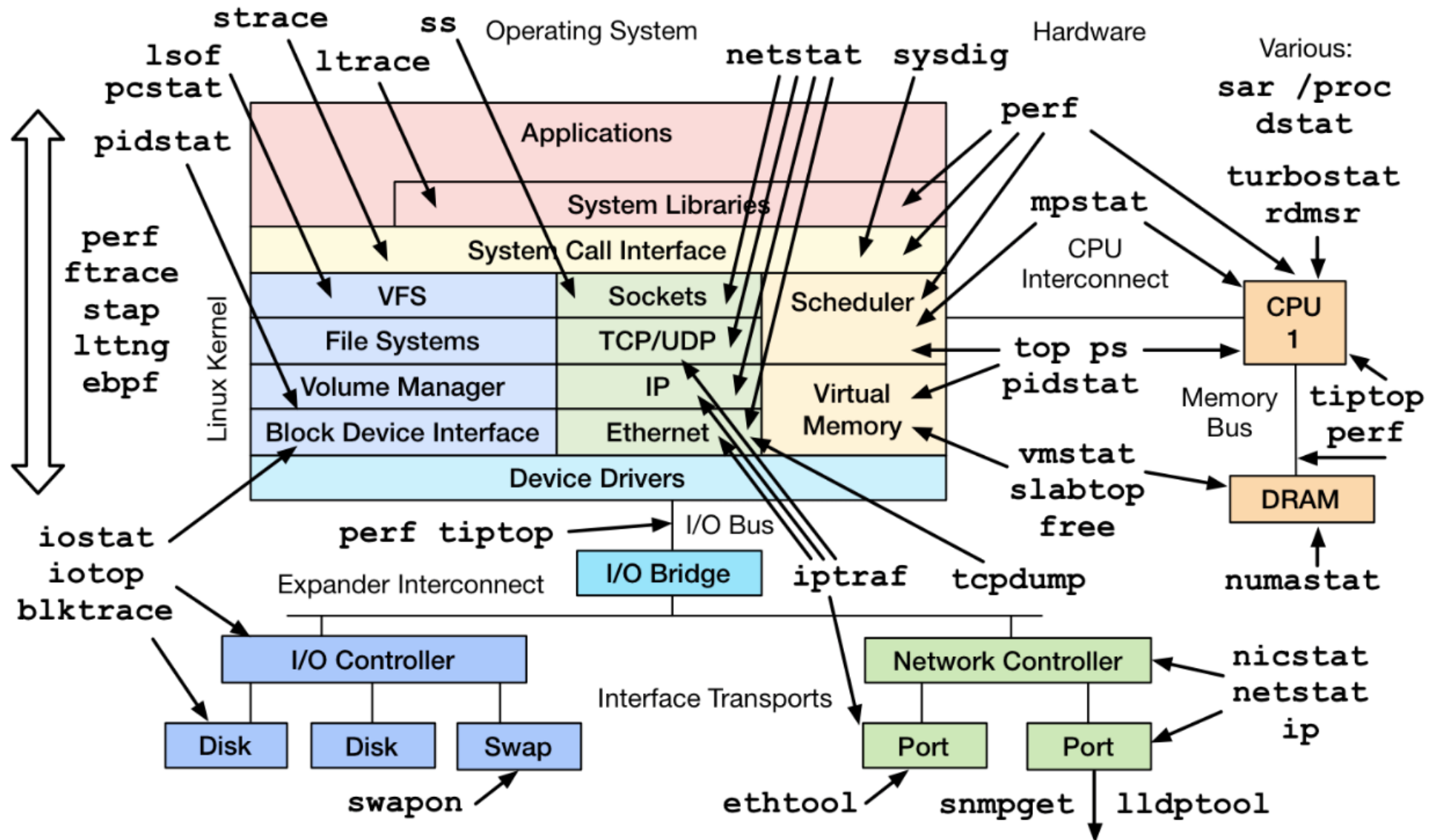
- Scheduler algorithm
 - echo deadline > /sys/block/<dev>/queue/scheduler
- IO request queue
 - echo 1024 > /sys/block/<dev>/queue/nr_requests

Network tuning

- Interrupts balance
 - `/proc/irq/IRQ/smp_affinity`
- Backlogs
 - `net.core.netdev_max_backlog`
 - `net.core.somaxconn`
 - `net.ipv4.tcp_max_syn_backlog`
- Reduce TCP overhead
 - `net.ipv4.tcp_sack`
 - `net.ipv4.tcp_fack`
- Reduce connection overhead
 - `net.ipv4.tcp_fin_timeout`
 - `net.ipv4.tcp_tw_reuse`
- Enable auto-tuning
 - `net.ipv4.tcp_moderate_rcvbuf`
 - `net.ipv4.tcp_window_scaling`

***The tool isn't important – it's important to
have a way to measure everything***

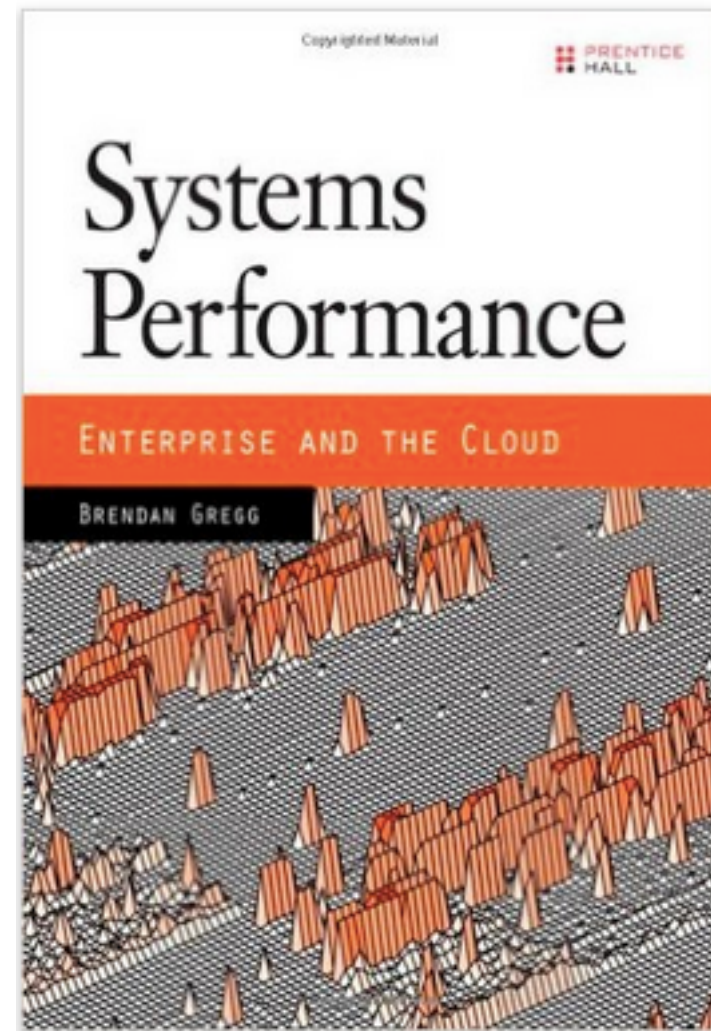
-- Brendan Gregg



Architects look at thousands of buildings during their training, and study critiques of those buildings written by masters. In contrast, most software developers only ever get to know a handful of large programs well — usually programs they wrote themselves — and never study the great programs of history. As a result, they repeat one another's mistakes rather than building on one another's successes.

--The Architecture of Open Source Applications

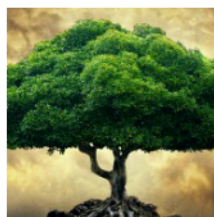
推荐两本书



reference

- <http://tutorials.jenkov.com/software-architecture/computer-architecture.html>
- <http://exadat.co.uk/2015/01/29/cpus-memory-storage-and-database-engines-the-shape-of-things-to-come/>
- <http://mechanical-sympathy.blogspot.com/2013/02/cpu-cache-flushing-fallacy.html>
- <http://duartes.org/gustavo/blog/post/what-your-computer-does-while-you-wait/>
- <https://software.intel.com/en-us/articles/detecting-memory-bandwidth-saturation-in-threaded-applications>
- <https://software.intel.com/en-us/articles/optimizing-applications-for-numa>
- https://www.thomas-krenn.com/en/wiki/Linux_Storage_Stack_Diagram
- http://www.mimuw.edu.pl/~lichota/09-10/Optymalizacja-open-source/Materialy/10%20-%20Dysk/gelato_ICE06apr_blktrace_brunelle_hp.pdf
- <http://codecapsule.com/2014/02/12/coding-for-ssds-part-2-architecture-of-an-ssd-and-benchmarking/>
- <http://codecapsule.com/2014/02/12/coding-for-ssds-part-3-pages-blocks-and-the-flash-translation-layer/>

- <http://cd-docdb.fnal.gov/0019/001968/001/Linux-Pkt-Recv-Performance-Analysis-Final.pdf>
- <http://www.slideshare.net/hisaki/x86-hardware-for-packet-processing>
- <http://balodeamit.blogspot.com/2013/10/receive-side-scaling-and-receive-packet.html>
- <https://hpi.de/plattner/research/tools-methods-for-enterprise-systems-design-and-engineering.html>
- <http://sdepl.ucsd.edu/cgi-bin/yman2html?m=tcp&s=7>
- <http://zh.wikipedia.org/wiki/Trie>
- http://en.wikipedia.org/wiki/Perfect_hash_function
- <http://goog-perftools.sourceforge.net/doc/tcmalloc.html>
- <https://software.intel.com/en-us/articles/avoiding-and-identifying-false-sharing-among-threads>
- <https://www.certificationkits.com/cisco-certification/cisco-ccnp-tshoot-642-832-exam-study-center/cisco-ccnp-tshoot-troubleshooting-networks/>
- <http://brendangregg.com/books.html>
- <http://blog.tanelpoder.com/2013/02/21/peeking-into-linux-kernel-land-using-proc-filesystem-for-quickndirty-troubleshooting/>
- http://goog-perftools.sourceforge.net/doc/heap_profiler.html
- <http://clang.llvm.org/docs/AddressSanitizer.html>
- <http://brendangregg.com/linuxperf.html>



淘宝千石

扫描上面的二维码，关注我吧