

Linux Kernel : Memory Management

(c) 2020 kaiwanTECH

Linux Kernel : Memory Management

Virtual Memory

- **Introduction**
- **The How**
- **The Why**

Linux Kernel : Memory Management

Why VM?

- **Memory is virtualized**

- Programmer does not worry about phy RAM, availability, etc

- **Information Isolation**

- Each process run's in it's own private VAS

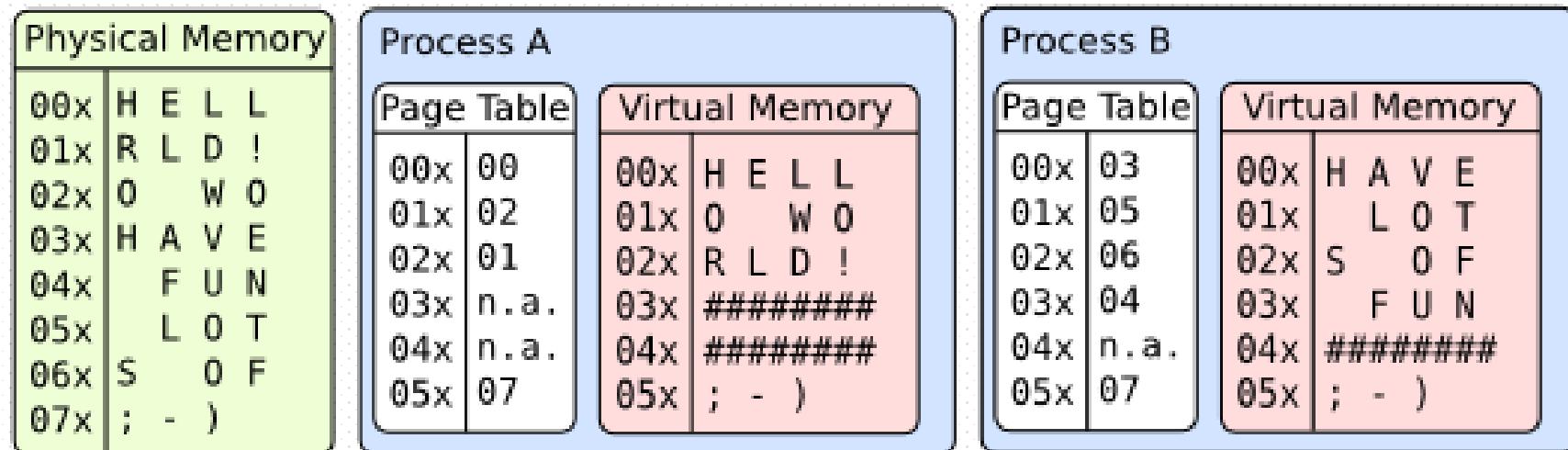
- **Fault Isolation**

- One process crashing will not cause others to be affected

- Multithreaded (firefox) vs multiprocess (chrome)

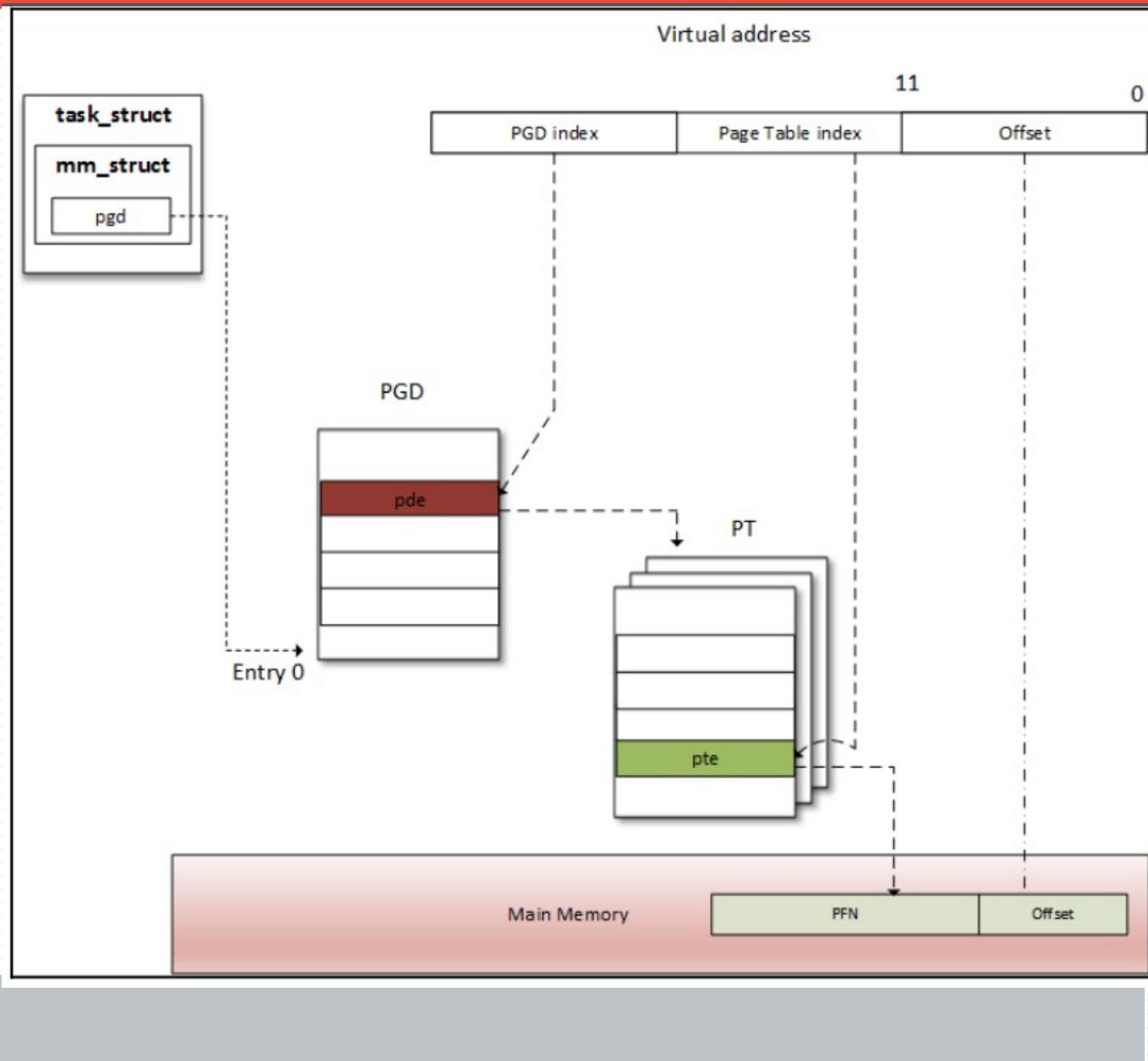
Linux Kernel : Memory Management

Paging Tables - a *very simplistic and conceptual view*



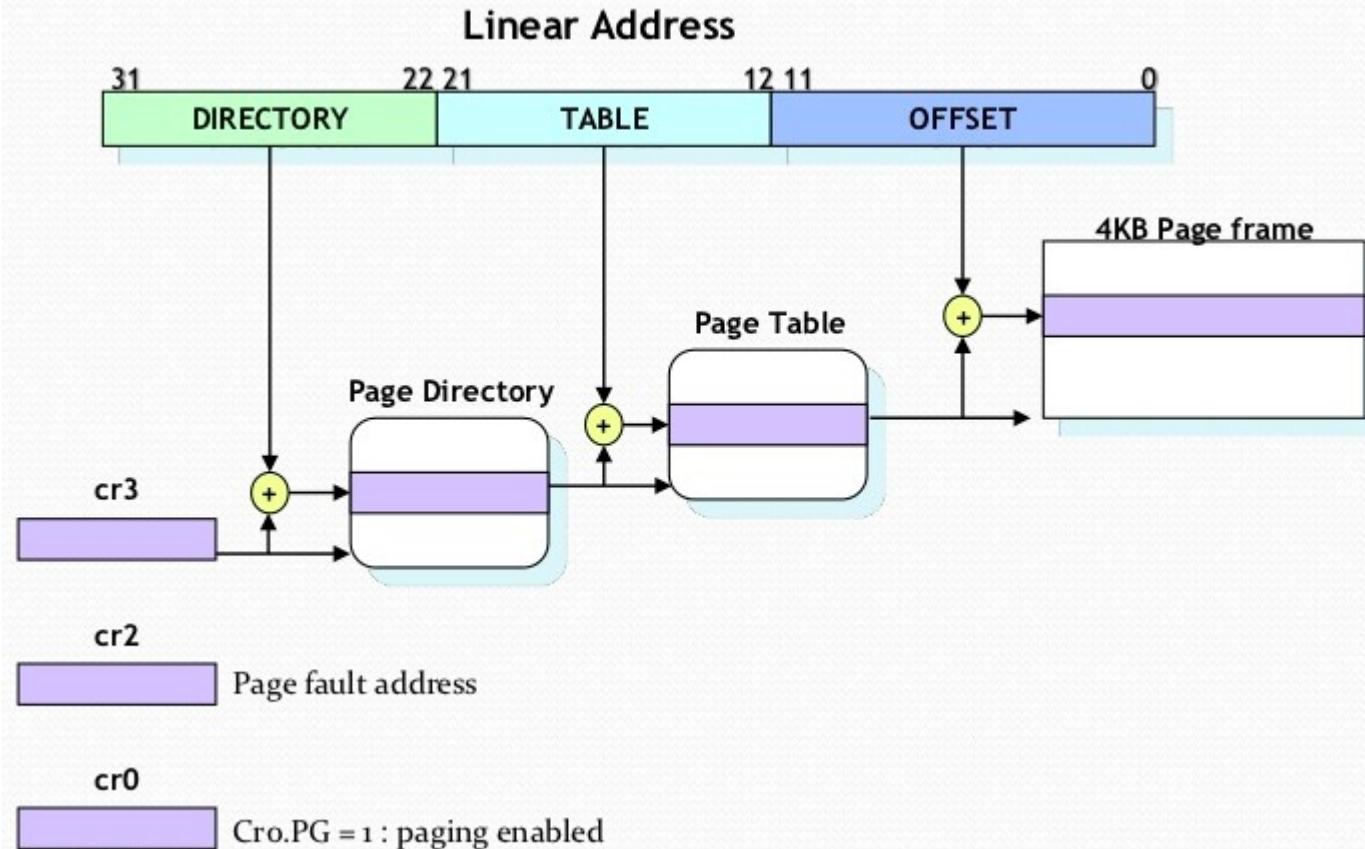
Source

Page Table Layout on the x86_64



Source: Linux Device Driver Development, John Madieu, Packt, 2017

Paging on the x86_32 : 2-level



Source

Paging on the x86_64 : 4-level

Source: ULK3

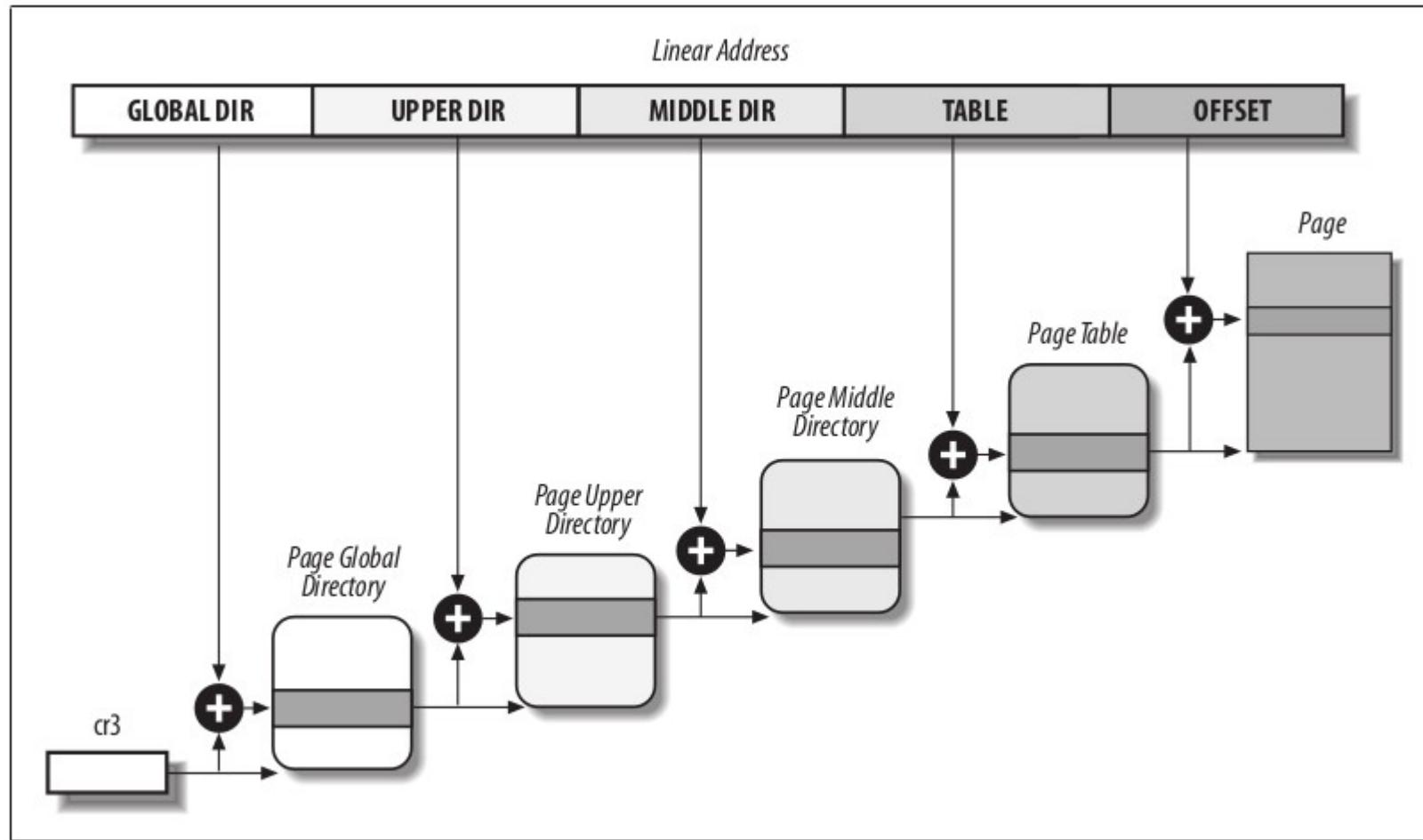


Figure 2-12. The Linux paging model

Paging on the x86_64

4-level Paging on x86_64

Article: Page Table Layout

6 65 55555 433 32 22 11 0

[Unused:16][PGD:9][PUD:9][PMD:9][PTE:9][offset:12]

11

Eq. (virtual) address:

00000000 00000000 00000000 00000000 00000000 10110001 01110000 00010000

00 00 00 00 00 . . . b1 70 10

0xb17010

And so we can now visualise what a given example address (actually represents):

6 5 4 3 2 1

43210987654321098765432109876543210987654321098765432109876543210987654321

[RESERVED][PGD][PUD] [PMD][PTE][OFFSET]

```
|           |           | -PAGE_SHIFT  
|           |           |-----PMD_SHIFT  
|           |-----PUD_SHIFT  
|-----PGDIR SHIFT
```

But how does this relate to a 'sparse' layout of pages, what do these offsets actually reference?

Paging on the x86_64

4-level Paging on x86_64

Article: Page Table Layout

The best way of showing how this fits together is diagrammatically, so taking the example address from above:

[RESERVED][PGD][PUD][PMD][PTE][OFFSET]

PGD offset = 0000000000 = 0

PUD offset = 000000000 = 0

PMD offset = 000000101 ≡ 5

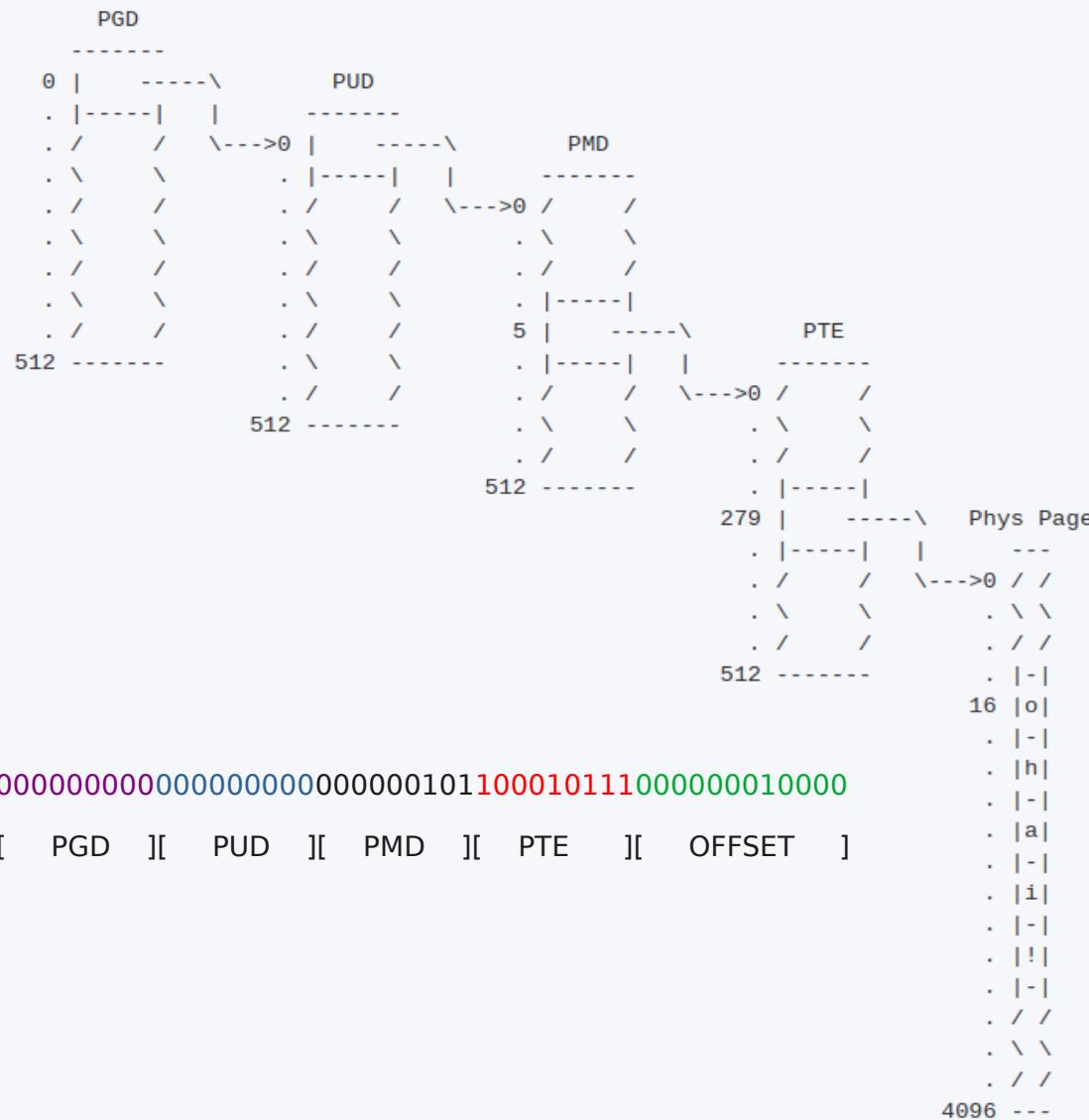
PTE offset = 100010111 = 279

phy offset = 000000010000 = 16

Page Table Layout on the x86_64

4-level Paging on x86_64

Article: Page Table Layout



Address Translation : ARM64

Address Translation on the ARMv8

Source (pg 12-3)

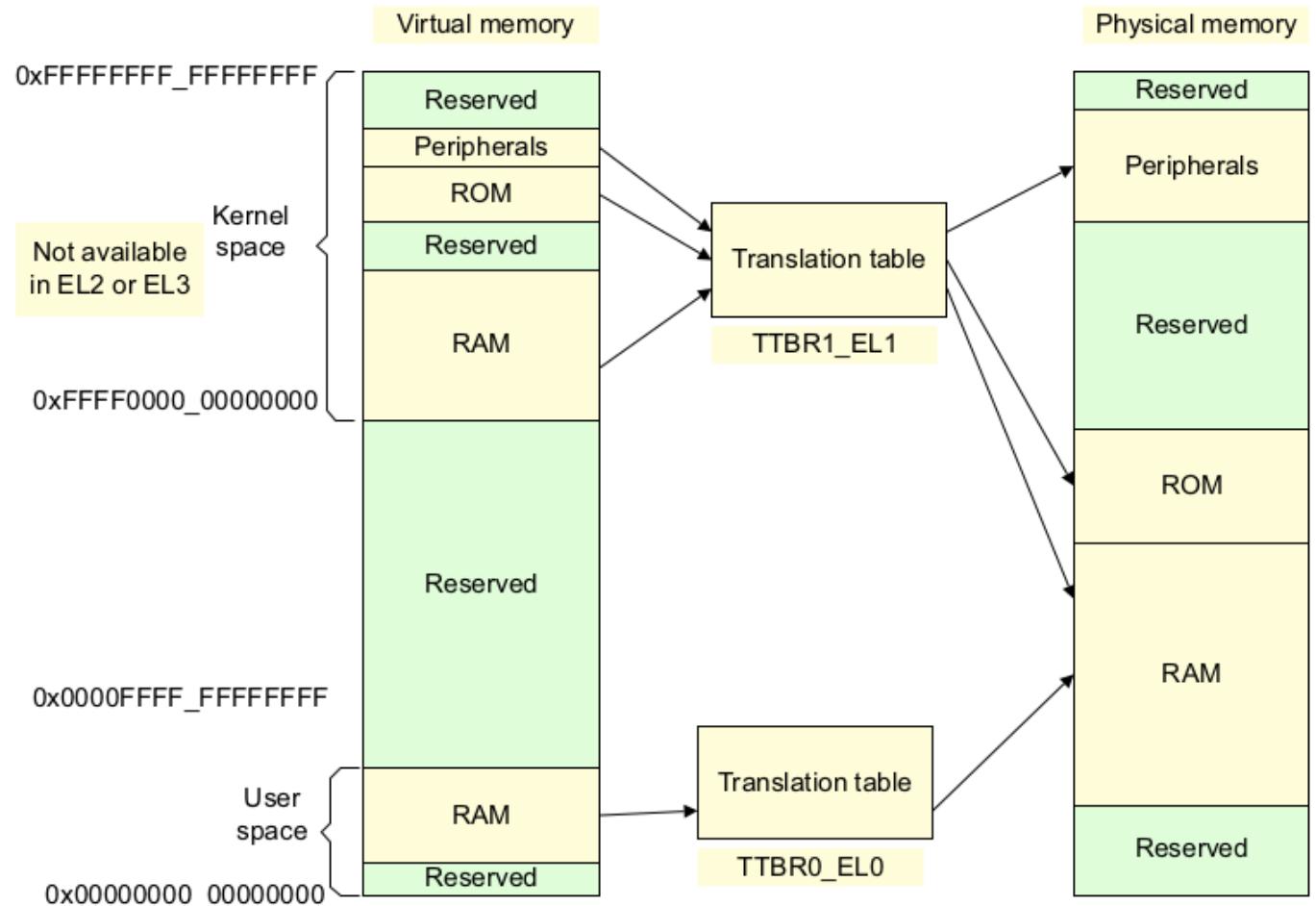
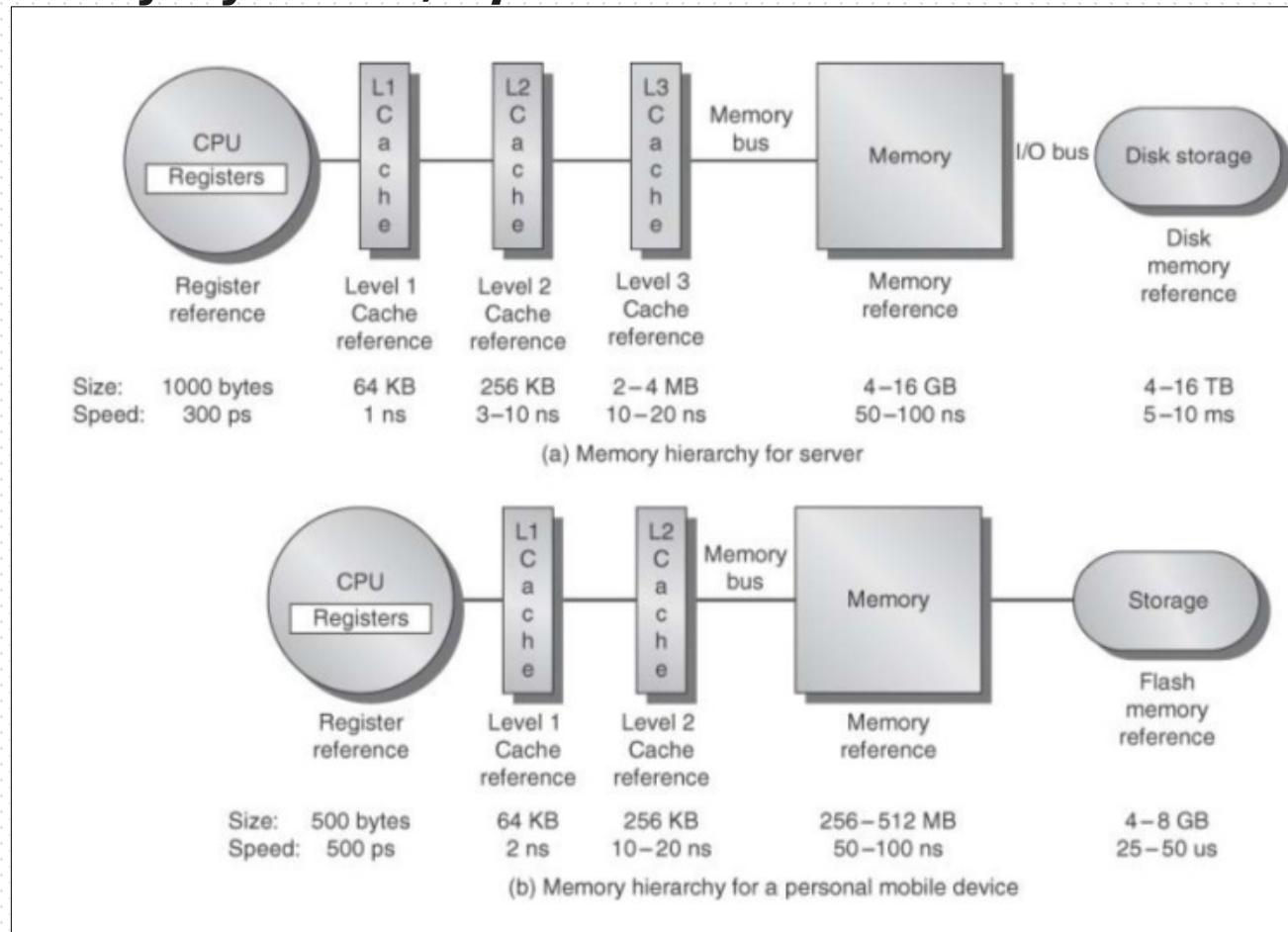


Figure 12-3 Address translation using translation tables

Linux Kernel : Memory Management

Memory Pyramid ; Speed and Size

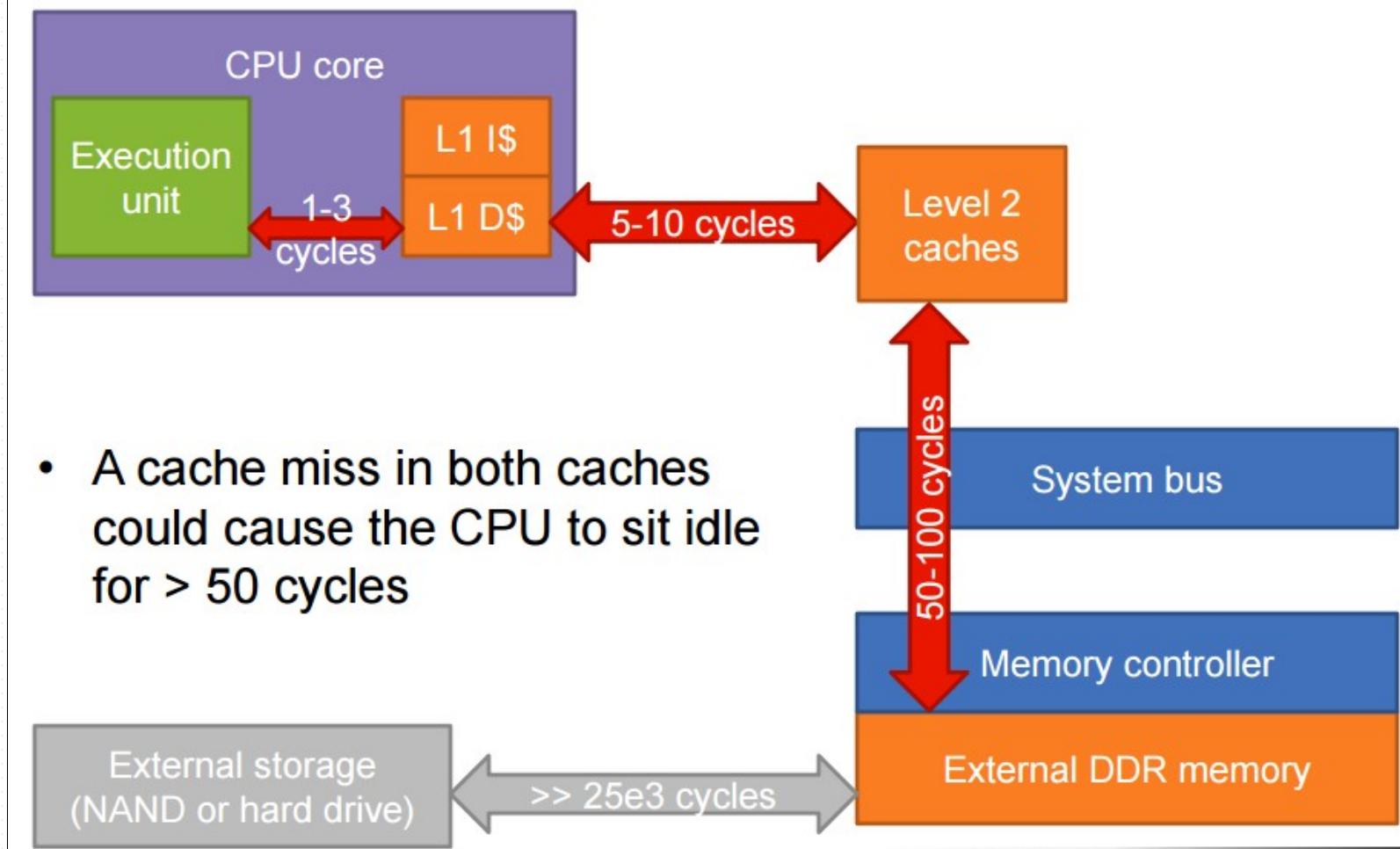


*Computer Architecture,
A Quantitative Approach,
5th Ed by Hennessy &
Patterson.*

Linux Kernel : Memory Management

CPU Caches

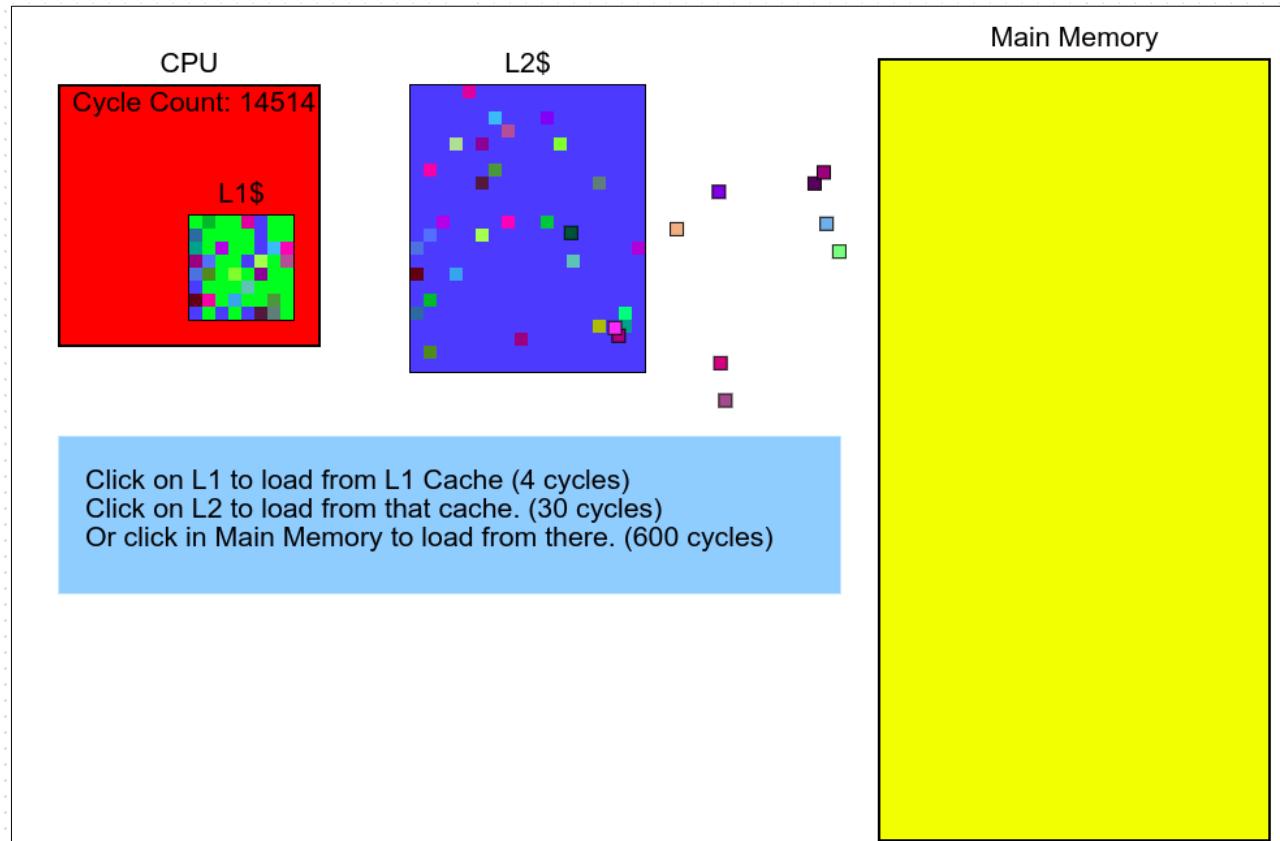
Memory latency



Linux Kernel : Memory Management

CPU Caches

**MUST SEE and TRY
(it's interactive) !!!**



<http://www.overbyte.com.au/misc/Lesson3/CacheFun.html>

Linux Kernel : Memory Management

CPU Caching [below, with : perf stat find arch/ -type f]

Performance counter stats for 'find arch/ -type f':

```
110.05 msec task-clock          # 0.005 CPUs utilized
  1,987  context-switches       # 0.018 M/sec
    23   cpu-migrations        # 0.209 K/sec
   277   page-faults            # 0.003 M/sec
13,93,53,680   cycles           # 1.266 GHz
10,06,09,295   instructions     # 0.72 insn per cycle
  2,00,78,349   branches         # 182.452 M/sec
    6,56,553   branch-misses     # 3.27% of all branches
21.608155978 seconds time elapsed ← first time (hardly any hits!)
 0.064850000 seconds user
 0.057220000 seconds sys
```

[...]

Performance counter stats for 'find arch/ -type f':

```
43.61 msec task-clock          # 0.213 CPUs utilized
   18   context-switches       # 0.413 K/sec
    0   cpu-migrations        # 0.000 K/sec
   277   page-faults            # 0.006 M/sec
8,70,28,961   cycles           # 1.996 GHz
7,12,75,161   instructions     # 0.82 insn per cycle
 1,40,37,741   branches         # 321.894 M/sec
    3,24,005   branch-misses     # 2.31% of all branches
0.204390622 seconds time elapsed ← second time, immediately; majority cache hits!
 0.028115000 seconds user
 0.016066000 seconds sys
```

Linux Kernel : Memory Management

CPU Caching : measuring LLC (Last Level Cache) hits/misses

1. With perf !

```
sudo perf top -e cache-misses -e cache-references -a -ns  
pid,cpu,comm
```

2. With the modern approach : eBPF !

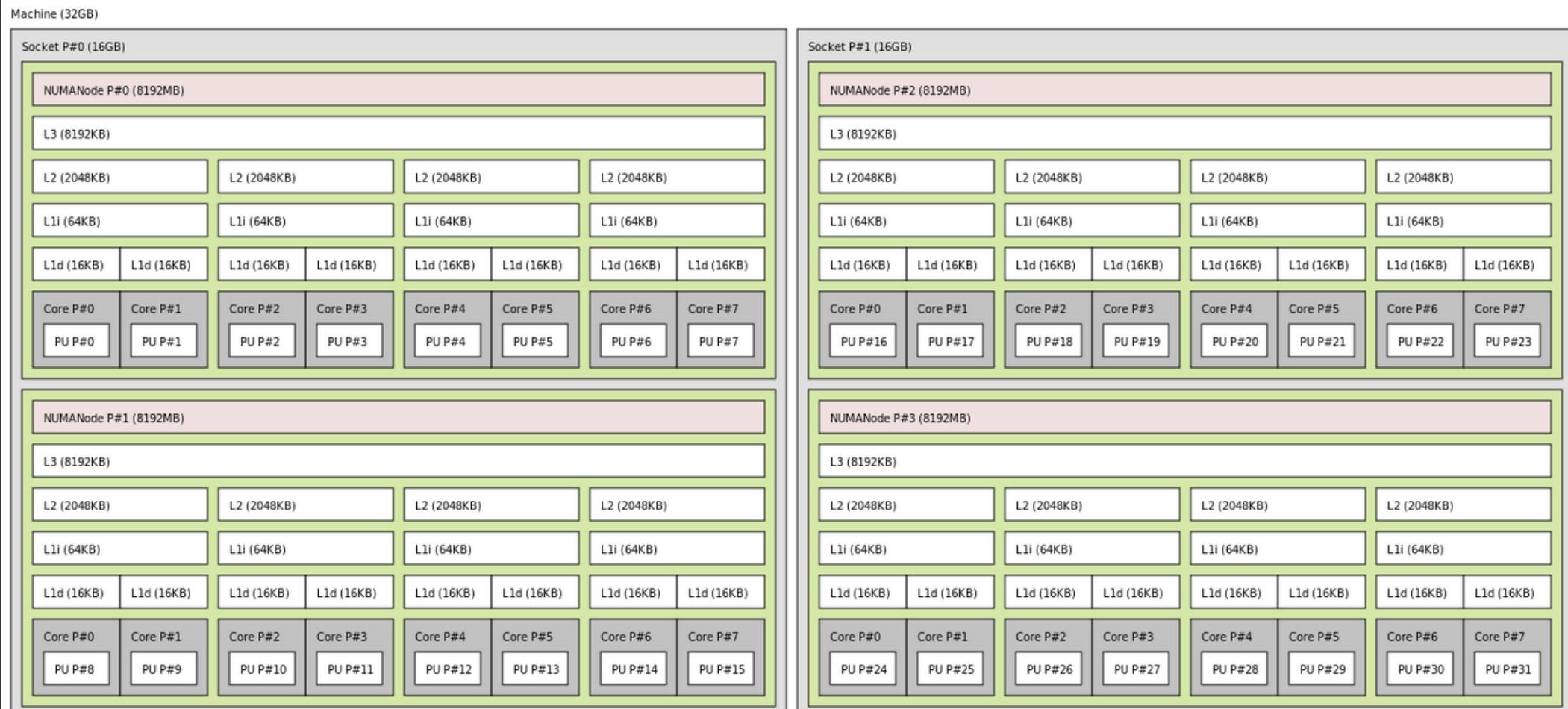
```
$ sudo llcstat-bpfcc
```

Running for 10 seconds or hit Ctrl-C to end.

[^] CPID	NAME	CPU	REFERENCE	MISS	HIT%
3828	Chrome_ChildIOT	3	8300	31400	0.00%
223043	chrome	2	6900	29700	0.00%
129894	nspr-3	3	7500	6100	18.67%
0	swapper/0	0	1432900	322500	77.49%
283	jbd2/sda6-8	0	2500	6200	0.00%
[...]					

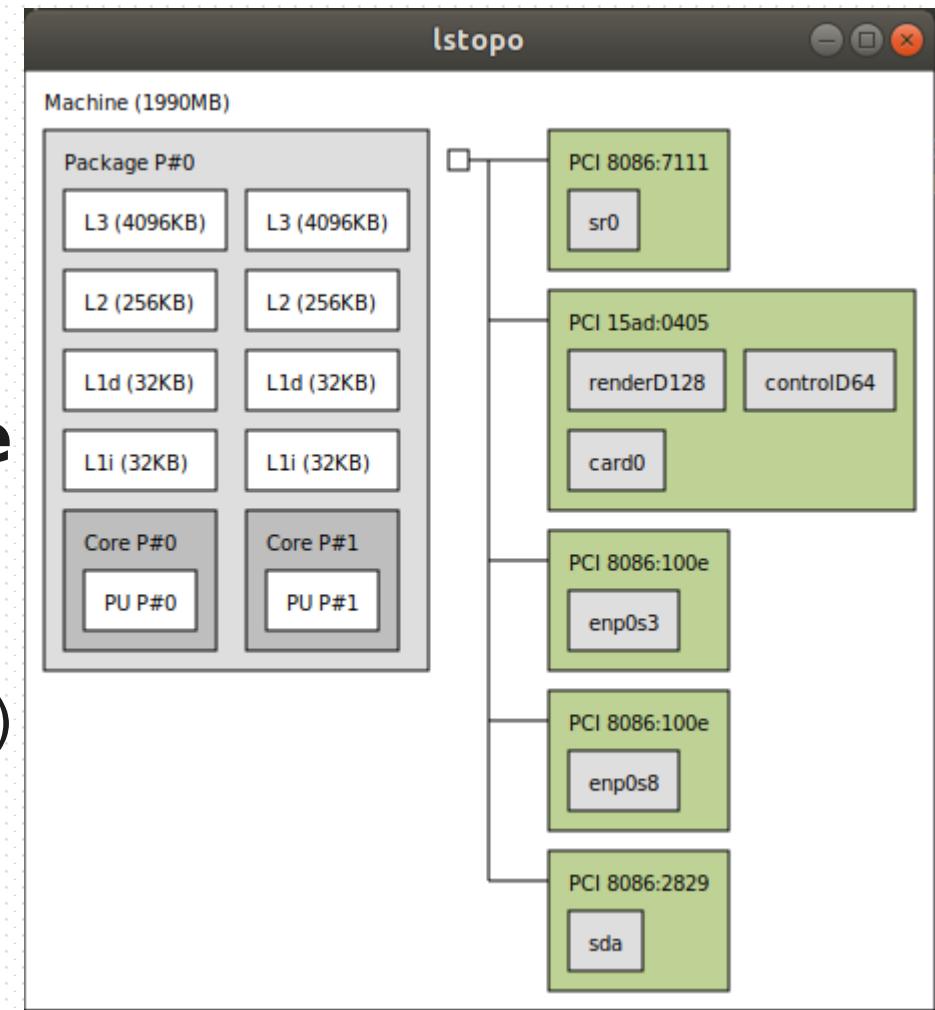
Linux Kernel : Memory Management

An example (pic) from the [Wikipedia page on CPU Cache](#):



Linux Kernel : Memory Management

- Can see any system's CPU/cache 'topology' in a fantastic graphical manner
- How? Install and run the *Istopo(1)* utility
 - Ubuntu: sudo apt install hwloc
 - (also use the CLI hwloc-* utils!)
 - Istopo



Linux Kernel : Memory Management

Software developers would do well to be aware how programming in a cache-aware manner can greatly optimize code:

- keep all important structure members ('hotspots') together and at the top
- ensure the structure memory start is CPU cacheline-aligned
- don't let a member "fall off" a cacheline (use padding if required; compiler can help)
- an **LLC (Last Level Cache) Miss** is expensive!
- use profilers (perf and **eBPF** tools (*bcc*)) to measure

[Good Article:

[How L1 and L2 CPU caches work, and why they're an essential part of modern chips, Joel Hruska, Feb 2016](#)

Also, careful of cache coherence issues; see these Wikipedia animations.

<< *Lookup 1LinuxMM.pdf* for details on:

Hardware Paging: n-Level Paging, Linux's 4 (or now 5)-level arch-independent paging model, IA-32, ARM-32, x86_64 hardware paging, etc

>>

Linux Kernel : Memory Management

Virtual Address Space Splitting

- Recall that Linux is a *monolithic OS*
- The reality is that every process has its own private VAS, and all of them share a “common mapping” for the upper region - the kernel segment!
- The *typical* “VM split”

On a 32-bit x86 (IA-32) Linux system
3 GB : 1 GB :: user-space : kernel-space

On a 32-bit ARM Linux system
2 GB : 2 GB :: user-space : kernel-space
[sometimes 3 GB : 1 GB]



Linux Kernel : Memory Management

Common powers of 2 : kilobyte onwards

Name (Symbol)	Value		Name (Symbol)	Value
<u>kilobyte</u> (kB)	10^3	2^{10}	<u>kibibyte</u> (KiB)	2^{10}
<u>megabyte</u> (MB)	10^6	2^{20}	<u>mebibyte</u> (MiB)	
gigabyte (GB)	10^9	2^{30}	<u>gibibyte</u> (GiB)	2^{30}
<u>terabyte</u> (TB)	10^{12}	2^{40}	<u>tebibyte</u> (TiB)	2^{40}
petabyte (PB)	10^{15}	2^{50}	<u>pebibyte</u> (PiB)	2^{50}
<u>exabyte</u> (EB)	10^{18}	2^{60}	<u>exbibyte</u> (EiB)	2^{60}
<u>zettabyte</u> (ZB)	10^{21}	2^{70}	<u>zebibyte</u> (ZiB)	2^{70}
<u>yottabyte</u> (YB)	10^{24}	2^{80}	<u>yobibyte</u> (YiB)	2^{80}

Linux Kernel : Memory Management

PAGE_OFFSET

*Splitting point is called **PAGE_OFFSET***

- location from which physical RAM is direct-mapped (1:1) into kernel VAS

On an IA-32:

```
$ zcat /proc/config.gz |grep -i VMSPLIT  
CONFIG_VMSPLIT_3G=y  
# CONFIG_VMSPLIT_2G is not set  
# CONFIG_VMSPLIT_1G is not set  
$
```

Tip: Numbers to recognize : PAGE_OFFSET

32-bit

0xc000 0000 = 3 GB ; 0x8000 0000 = 2 GB

64-bit

0x0000 0100 0000 0000 = 1 TB (2^{40}) ;

0x0000 0200 0000 0000 = 2 TB (2^{41}) << MIPS >>

0x0000 8000 0000 0000 = 128 TB (2^{48}) << x86_64 4-level >>

Linux Kernel : Memory Management

64-bit VAS

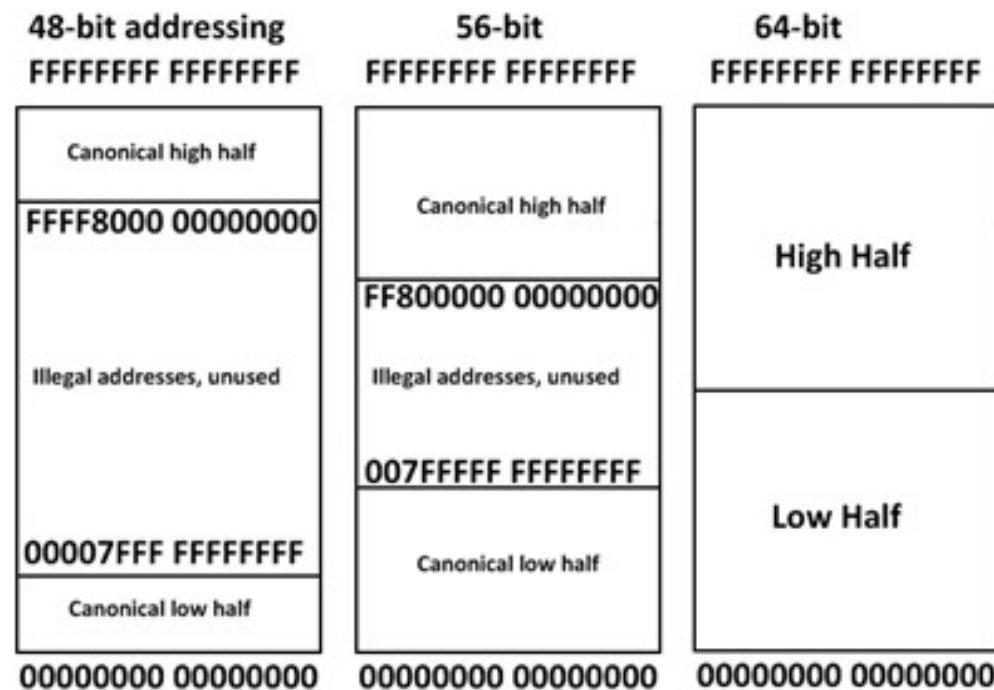


Figure 2 - Memory Addressing

Pic: Professional Linux Kernel Architecture, Mauerer

Linux Kernel : Memory Management

Arch-specific VM Splits

With standard 4 KB Page size

Arch	N-Level	Addr Bits	VM "Split"	Userspace		Kernel-space	
				Start vaddr	End vaddr	Start vaddr	End vaddr
IA-32	2	32	3 GB : 1 GB	0x0	0xbffff fffff	0xc000 0000	0xfffff fffff
ARM	2	32	2 GB : 2 GB	0x0	0x7fff fffff	0x8000 0000	0xfffff fffff
x86_64	4	48	128 TB : 128 TB	0x0	0x0000 7ffff ffff fffff	0xfffff 8000 0000 0000	0xfffff fffff ffff fffff
	5*	56	64 PB : 64 PB	0x0	0x00ff fffff ffff fffff	0xff00 0000 0000 0000	0xfffff fffff ffff fffff
Aarch64	3	39	512 GB : 512 GB	0x0	0x0000 007f ffff fffff	0xfffff ff800 0000 000	0xfffff fffff ffff fffff
	4	48	256 TB : 256 TB	0x0	0x0000 fffff ffff fffff	0xfffff 0000 0000 0000	0xfffff fffff ffff fffff

* >= 4.14 Linux

IMP Update! Please see the newly updated (as of Linux 5.0) x86_64 process memory layout doc here:
https://www.kernel.org/doc/Documentation/x86/x86_64/mm.txt

Linux Kernel : Memory Management

03 Jan 2018 : Meltdown and Spectre hardware errors

- Root cause: CPU speculative code/data fetches
- Security concerns
- “Due to a processor-level bug (to do with speculative code execution), the traditional approach of keeping kernel paging tables and thus kernel virtual address space (VAS) within the process - anchored in the upper region - is now susceptible to attack! And thus needs to be changed. Linux kernel devs have been working furiously at it...”
- Mitigation: **KPTI : kernel page table isolation**

```
$ uname -a
Linux ... 4.14.11-300.fc27.x86_64 #1 SMP Wed Jan 3 13:52:28 UTC 2018 x86_64 x86_64 x86_64 GNU/Linux
$ dmesg |grep "isolation"
Kernel/User page tables isolation: enabled
$ dmesg |grep -i spectre
kern :info : [Sat Mar  3 15:16:00 2018] Spectre V2 : Mitigation: Full generic retpoline
kern :info : [Sat Mar  3 15:16:00 2018] Spectre V2 : Spectre v2 mitigation: Filling RSB on context switch
kern :info : [Sat Mar  3 15:16:00 2018] Spectre V2 : Spectre v2 mitigation: Enabling Indirect Branch Prediction Barrier
$
```

<< See 2LinuxMM doc for more details >>

Linux Kernel : Memory Management

Meltdown/Spectre Effects

- Traditionally, Linux would keep the kernel paging tables as part of *every process* - *in its “upper region” - the ‘canonical higher half’*
- With Meltdown/Spectre, that decision has been reversed!
- Performance impact: resulted in (much) higher context-switching times for user ↔ kernel system call code paths
 - Already visible in ‘cloud’ virtualization workloads, etc
- [Details: “... In current kernels, each process has a single PGD; one of the first steps taken in the KPTI patch series is to create a second PGD. The original remains in use **when the kernel** is running; it maps the full address space. The **second** is made active (at the end of the patch series) **when the process is running in user space**. It points to the same directory hierarchy for pages belonging to the process itself, but the portion describing kernel space (which sits at the high end of the virtual address space) **is mostly absent**. ...” : *LWN, 20Dec2017, Corbet*
- New! 14 May 2019: *Intel and other hardware vendors disclose a similar class of hardware-related (side channel) security vulnerabilities, collectively called MDS (Microarchitectural Data Sampling); link*

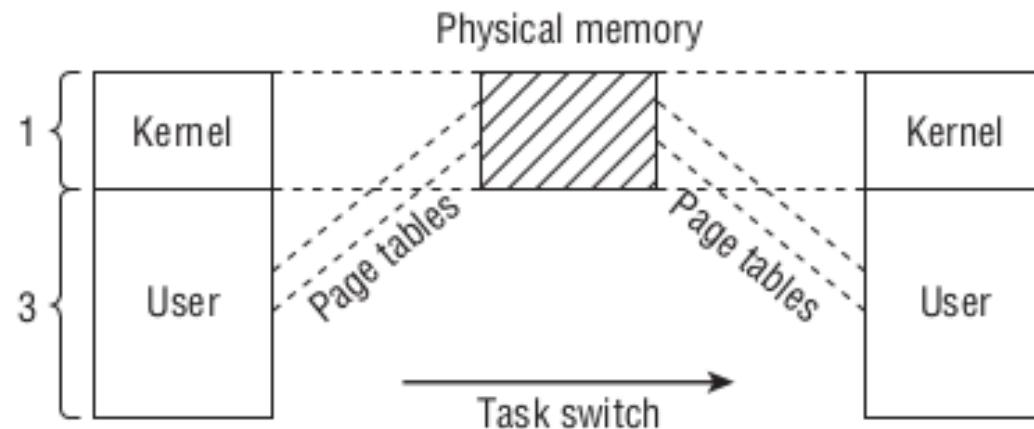
Linux Kernel : Memory Management

Meltdown/Spectre Effects

- Mitigating the detrimental slowdowns due to the page table isolation (KPTI) solution
 - *Use THP (Transparent Huge Pages)*
 - *Use PCID (on x86_64); Process Context Identifiers; helps reduce TLB shootdowns significantly. Available “properly” from Linux 4.14*
A useful article: [PCID is now a critical performance/security feature on x86 \[ARM\[64\] equivalent is the Context ID register; see this: ‘ARM Context ID Register & Process Context Switch’\]](#)
 - *Identify and reduce*
 - *large system call usage*
 - *interrupt sources*
 - [29Oct2019]
[Running on Intel? If you want security, disable hyper-threading, says Linux kernel maintainer; Speculative execution bugs will be with us for a very long time](#)

Linux Kernel : Memory Management

- Physical – platform – RAM is “direct-mapped” into the kernel segment at boot
- The OS is the resource manager; it manages RAM, does *not* hoard it!
- Gives it to userspace on demand (via demand paging)



Pic: Professional Linux Kernel Architecture, Mauerer

Figure 3-14: Connection between virtual and physical address space on IA-32 processors.

Linux Kernel : Memory Management

THP – Transparent Huge Pages

A quick summarization

- 2.6.38 onwards
- Kernel address space is always THP-enabled, **reducing TLB pressure** (for eg.: with typical 4Kb page size, to translate (va to pa) 2 MB space, one would require $(2 * 1024\text{Kb} / 4\text{Kb}) = 512$ TLB entries/translations.
With THP enabled, each kernel page is 2 MB, thus requiring 1 TLB entry only!)
- Transparent to applications*
- (Userland) Pages are swappable (split into 4k pages & swapped)
- Some CPU overhead (khugepaged checks continuously for smaller 4K pages that can be merged together to create a 2MB huge page)
- Currently works only on anonymous memory regions (planned to add support for page cache and tmpfs pages)

Linux Kernel : Memory Management

Show and run the vm_user.c app

Running on a 32-bit IA-32 system (and OS) with a 3 GB : 1 GB VM Split

```
$ getconf -a |grep LONG_BIT
LONG_BIT                      32
$ ./vm_user
wordsize : 32
&main = 0x0804847d &g = 0x0804a028 &u = 0804a030 &heapptr = 0x082a5008 &loc = 0xbfa6e5d8
$
```

Running on an ARM-32 system (and OS) with a 2 GB : 2 GB VM Split

(Qemu-based ARM Versatile Express platform)

```
$ ./vm_user
32-bit: &main = 0x00008578 &g = 0x00010854 &u = 0x0001085c &heapptr = 0x00011008 &loc = 0x7ebf4d08
$
```

Running on a 64-bit (x86_64) system (and OS) [128 TB : 128 TB VM split]

```
$ getconf -a |grep LONG_BIT
LONG_BIT                      64
$ ./vm_user
wordsize : 64
&main = 0x0000557de98f3165 &g = 0x0000557de98f6010 &u = 0x0000557de98f6018 &heapptr = 0x0000557deb690260 &loc =
0x00007ffcf811766c
$
```

Running on an ARM64 (Raspberry Pi 3) [256 TB : 256 TB VM split]

```
$ ./vm_user
wordsize : 64
&main = 0x0000aaaac50028a4 &g = 0x0000aaaac5013010 &u = 0x0000aaaac5013018 &heapptr = 0x0000aaaac8ec0260 &loc =
0x0000ffffe6e3c1fc
$
```

Linux Kernel : Memory Management

Physical Memory Organization

[N]UMA, zones

Buddy Sys Alloc

Slab Cache

Content

Kernel dynamic allocation APIs

BSA

kmalloc

vmalloc

Custom slab

Kernel segment - diagrams + dmesg @boot k layout

Page Cache -diagram, blog article [mmap]

OOM killer

Demand paging

VM overcommit

Page fault handling basics

Glibc malloc behavior

User VAS mapping - proc

Linux Kernel : Memory Management

Physical Memory Organization

SMP - Symmetric Multi Processing

One OS, one RAM, multiple CPUs (cores)

AMP - Asymmetric Multi Processing

> 1 OS, one RAM, multiple CPUs (cores)

UMA - Uniform Memory Access

RAM treated as a uniform hardware resource by the OS; does not matter which CPU a thread is running upon, memory looks the same

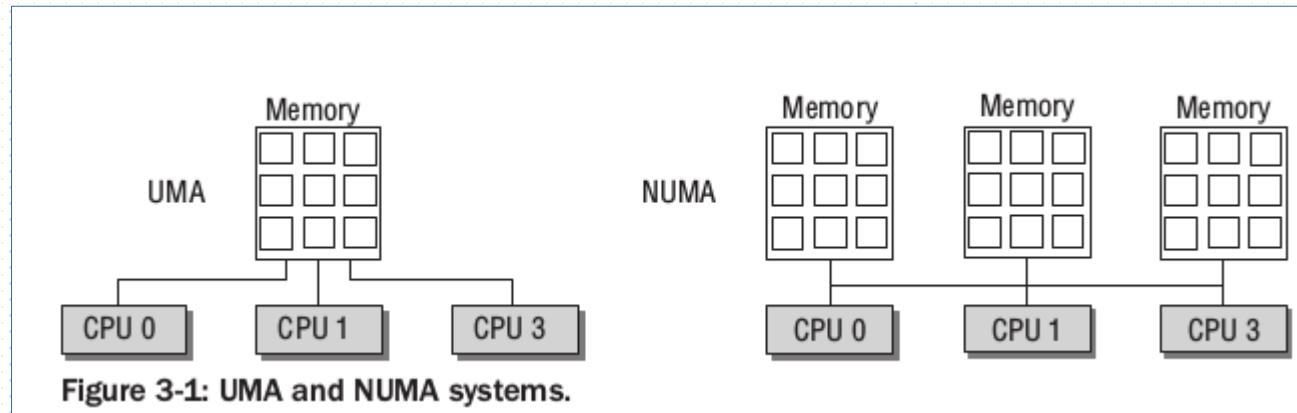
NUMA - Non-Uniform Memory Access

RAM “banks” treated as distinct hardware resources by the OS – called ‘nodes’; does matter on which CPU a thread is running upon, should use RAM from the “closest” node.

Linux (Windows, UNIXes) is NUMA-aware.

Linux Kernel : Memory Management

[N]UMA



Above pic: "Professional Linux Kernel Architecture", W Mauerer, Wrox Press

[NUMA on Wikipedia](#)

NUMA machines are always multiprocessor; each CPU has local RAM available to it to support very fast access; also, all RAM is linked to all CPUs via a bus.

(Above) [Image Source](#)

Linux Kernel : Memory Management

NUMA system, 4 nodes

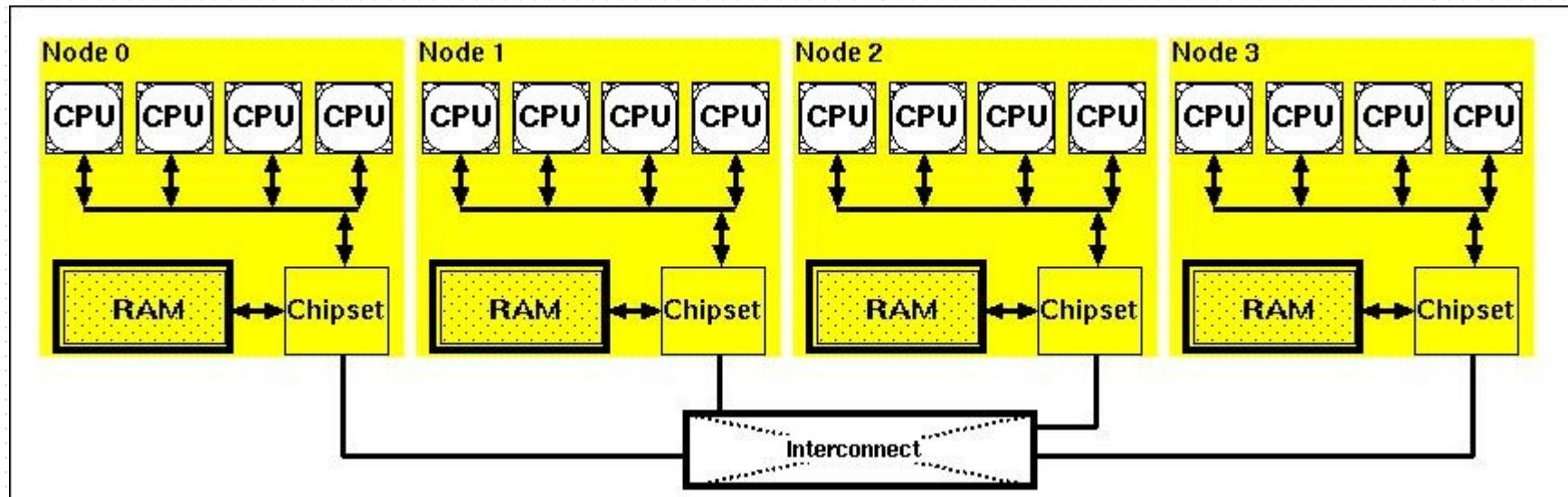
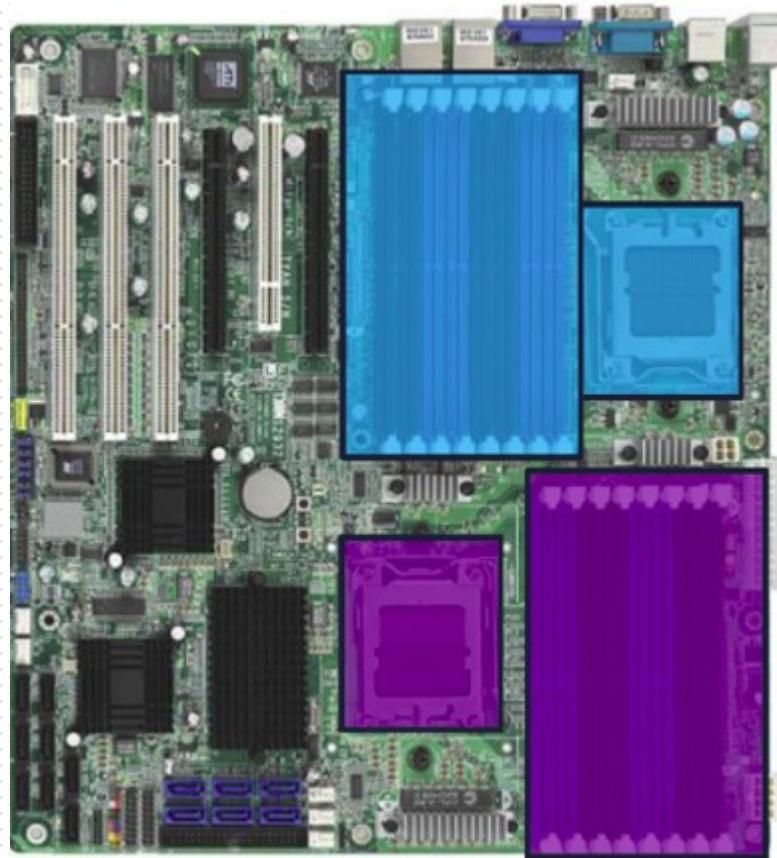


Image Source

Linux Kernel : Memory Management

NUMA server



Linux Kernel : Memory Management

Actual NUMA servers



NUMA-Q, IBM



Integrity SuperdomeX, HP

Linux Kernel : Memory Management

Linux

RAM is divided into *nodes*

One node for each CPU core that has *local RAM* available to it

A node is represented by the data structure *pg_data_t*

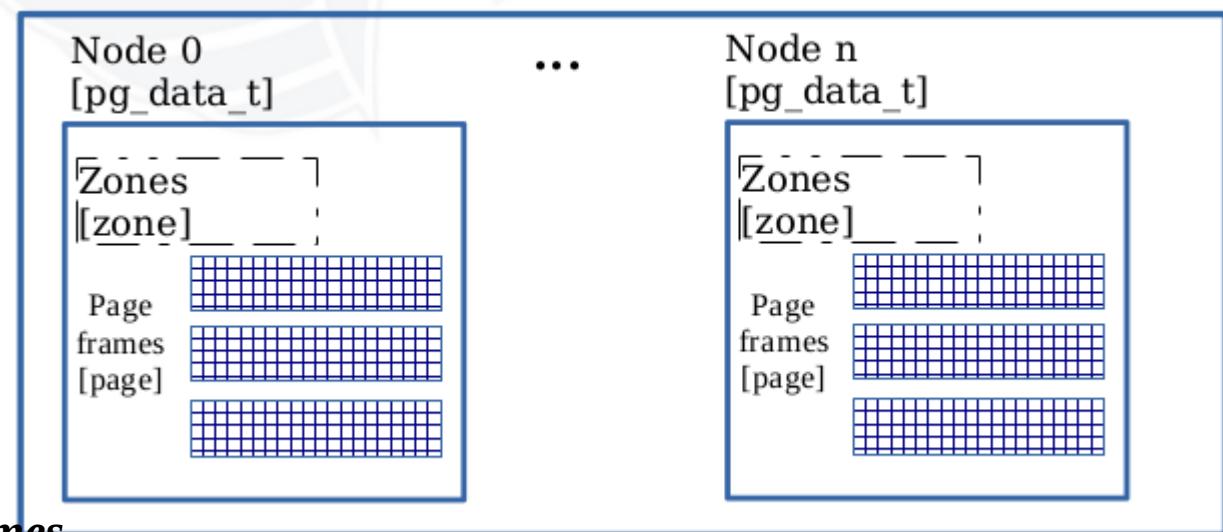
**Each node is further split into
(at least one, upto four) *zones***

`ZONE_DMA`

`ZONE_DMA32`

`ZONE_NORMAL`

`ZONE_HIGHMEM`



Each zone consists of *page frames*

A page frame is represented (and managed) by the data structure *page*

Linux Kernel : Memory Management

Emulating a NUMA box with Qemu

```
$ qemu-system-x86_64 --enable-kvm  
-kernel <...>/4.4.21-x86-tags/arch/x86/boot/bzImage  
-drive file=images/wheezy.img,if=virtio,format=raw  
-append root=/dev/vda -m 1024 -device virtio-serial-pci  
-smp cores=4 -numa node,cpus=0-1,mem=512M  
-numa node,cpus=2-3,mem=512M  
-monitor stdio
```

...

```
(qemu) info numa  
2 nodes  
node 0 cpus: 0 1  
node 0 size: 512 MB  
node 1 cpus: 2 3  
node 1 size: 512 MB
```

Linux Kernel : Memory Management

Zones

- Often enable a workaround over a hardware or software issue.
Eg.
 - ZONE_DMA : 0 – 16 MB on old Intel with the ISA bus
 - ZONE_HIGHMEM : what if the 32-bit platform has more RAM than there is kernel address space? (eg. an embedded system with 3 GB RAM on a machine with a 3:1 VM split)
 - ... and so on
- The ‘lowmem’ region – direct mapped platform RAM – is often placed as ZONE_NORMAL (but not always)
- Kernel inits zones at boot

Linux Kernel : Memory Management

False Sharing

[“Avoiding and Identifying False Sharing Among Threads”, Intel](#)

In symmetric multiprocessor (SMP) systems, each processor has a local cache. The memory system must guarantee cache coherence. False sharing occurs when threads on different processors modify variables that reside on the same cache line. This invalidates the cache line and forces an update, which hurts performance. This article covers methods to detect and correct false sharing.

Prevented by ensuring the variables cannot “live” in the same CPU cacheline

Linux Kernel : Memory Management

An example of false-sharing corrected

Each zone is represented by struct zone, which is defined in the

File : [3.10.24]: include/linux/mmzone.h

```
struct zone {
[...]
    spinlock_t          lock;
--snip--
    struct free_area   free_area[MAX_ORDER];
--snip--
    ZONE_PADDING(_pad1_)

    /* Fields commonly accessed by the page reclaim scanner */
    spinlock_t          lru_lock;
<<
File : [3.10.24] : include/linux/mmzone.h
90 /*
91 * zone->lock and zone->lru_lock are two of the hottest locks in the kernel.
92 * So add a wild amount of padding here to ensure that they fall into separate
93 * cachelines. There are very few zone structures in the machine, so space
94 * consumption is not a concern here.
95 */
```

>>

...

Linux Kernel : Memory Management

Linux Kernel Memory Allocation

- Layered Approach

Slab Cache/Allocator

**Buddy System Allocator (BSA) /
aka 'Page Allocator'**

Linux Kernel : Memory Management

Buddy System Allocator (BSA)

The kernel's memory alloc/de-alloc “engine”

The key data structure is the ‘free list’:

- An array of pointers to doubly linked circular lists
- 0 to MAX_ORDER-1 (MAX_ORDER is usually 10); thus 11 entries
- each index is called the ‘order’ of the list
 - holds a chain of free memory blocks
 - each memory block
 - is of size 2^{order}
 - is guaranteed to be *physically contiguous*
- *allocations are always rounded power-of-2 pages*
 - exception: when using the *alloc_pages_exact* API

Linux Kernel : Memory Management

The BSA (Buddy System Allocator) free list

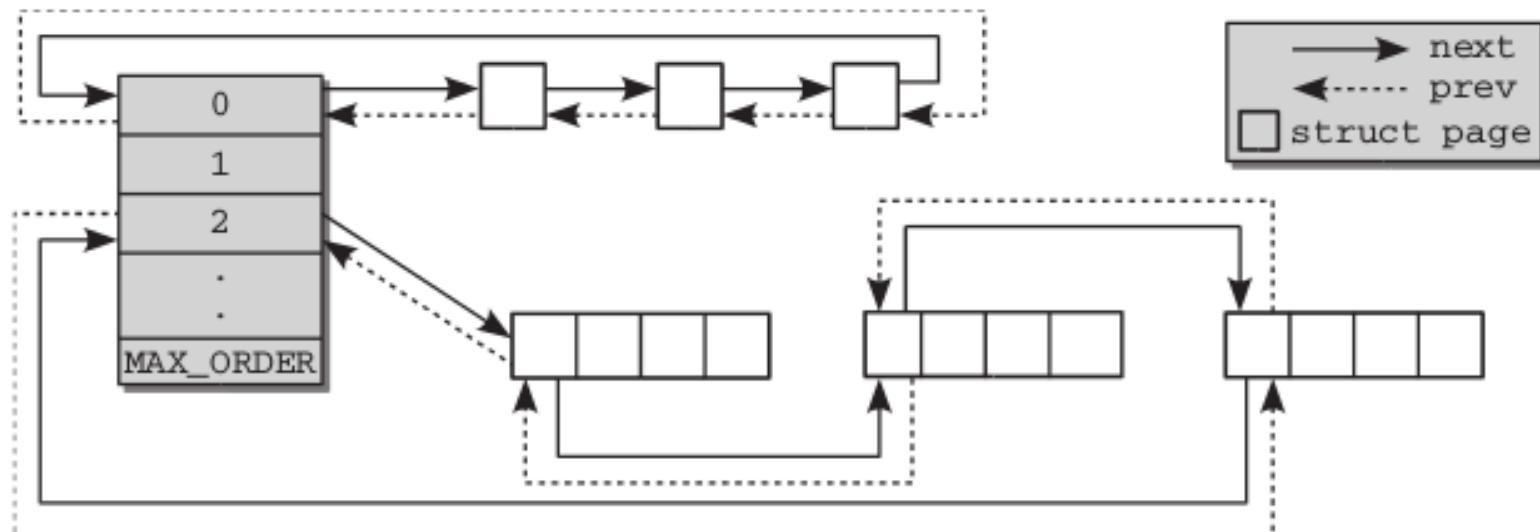


Figure 3-22: Linking blocks in the buddy system.

Linux Kernel : Memory Management

On a system with 1 GB RAM:

```
$ cat /proc/buddyinfo
Node 0, zone      DMA      0      0      0      1      2      1      1      0      1      1      1      3
Node 0, zone    DMA32     49     32     11      5      4     41     15     13     16     15     90
$
```

On a system with 16 GB RAM:

```
$ cat /proc/buddyinfo
Node 0, zone      DMA      2      3      6      2      3      1      0      0      1      1      1      3
Node 0, zone    DMA32   14617   3480   1303    465    388    97    63    23    11    3      0
Node 0, zone   Normal   1208   1909    303    609    552   180    49    10     3     0      0
$
```

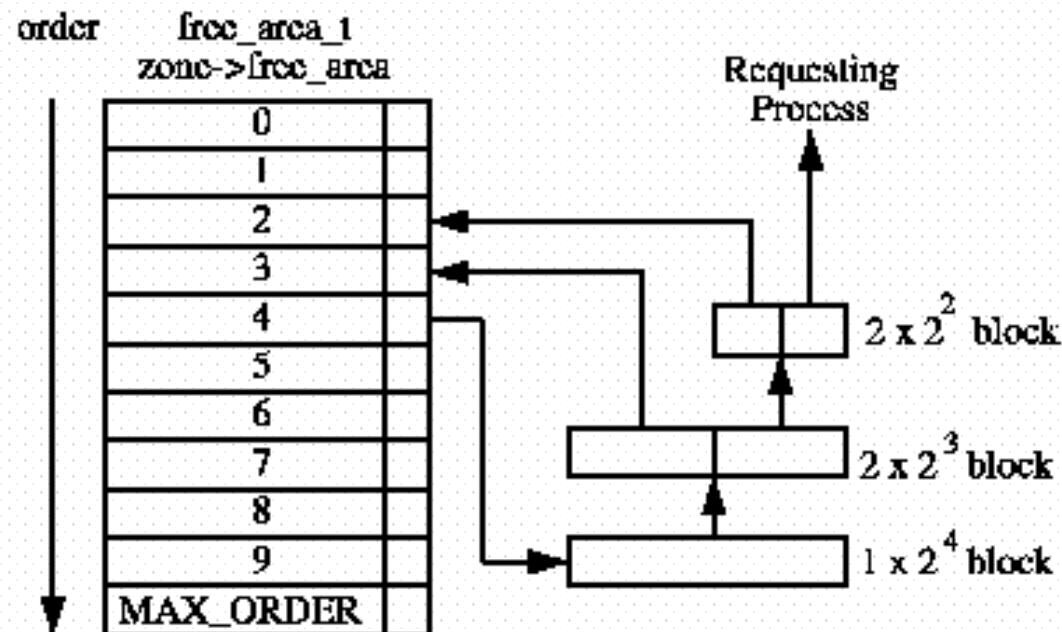
Order 0

Order 10

Linux Kernel : Memory Management

BSA: alloc of 4 KB (2^2) chunk by splitting:

- a 2^4 (16 KB) chunk into 2 chunks of 2^3 (8K) each
- further splitting the 2^3 (8 KB) chunk into 2 chunks of 2^2 (4 KB) each
- the left-over 2^3 chunk goes onto order 3
- the left-over 2^2 chunk goes onto order 2



Linux Kernel : Memory Management

BSA: one freelist per node per zone (in each node)

- Makes the (de)alloc NUMA-aware
- Actually, from 2.6.24, it's more complex: the addition of *per-migration-type* freelists means there are upto six freelists per node per zone *in each node)!

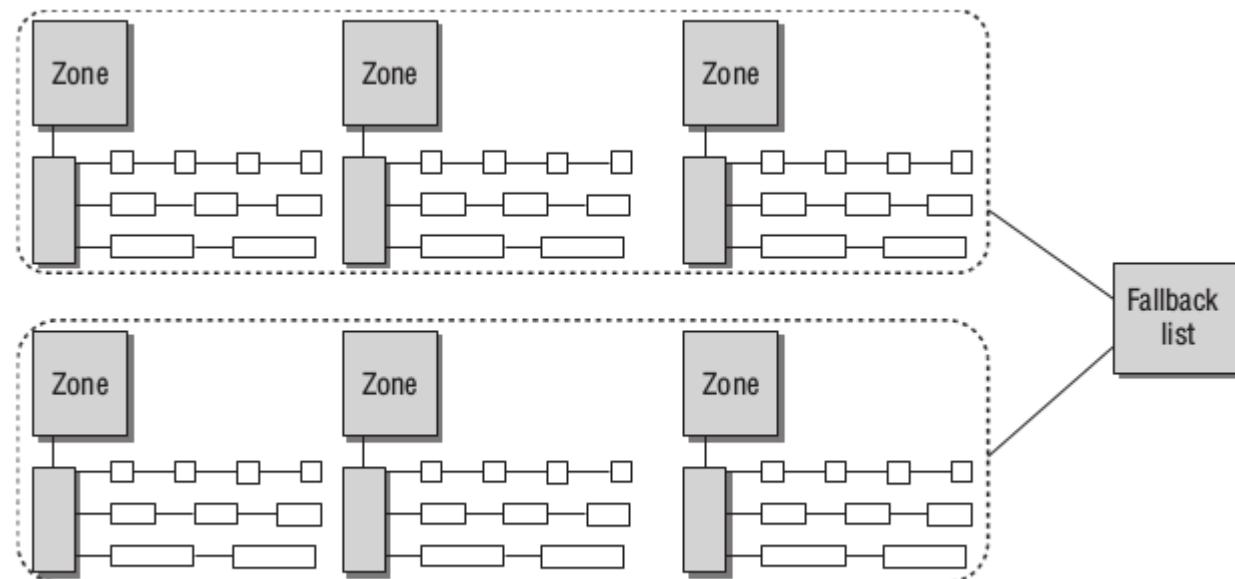


Figure 3-23: Relationship between buddy system and memory zones/nodes.

Linux Kernel : Memory Management

BSA / Page Allocator

Pros

- very fast
- physically contiguous page-aligned memory chunks
- actually helps defragment RAM by *merging* 'buddy' blocks
 - buddy block: block of the same size and physically contiguous

Cons

- primary issue: internal fragmentation / **wastage!**
 - granularity is a page; so asking for 128 bytes gets you 4096

Linux Kernel : Memory Management

Slab Cache / Slab Allocator

- The BSA's major wastage problem is addressed by the slab allocator
- Slab : two Big Ideas
 - caches 'objects' - commonly used kernel data structures
 - caches fragments of RAM

Look it up!

```
sudo vmstat -m
```

Linux Kernel : Memory Management

Check out the available slabs for the kmalloc:

	Num	Total	Size	Pages
\$ sudo vmstat -m grep "kmalloc\-\-"				
kmalloc-8192	52	52	8192	4
kmalloc-4096	109	136	4096	8
kmalloc-2048	267	272	2048	8
kmalloc-1024	860	1040	1024	8
kmalloc-512	547	672	512	8
kmalloc-256	1198	1728	256	16
kmalloc-192	1470	1470	192	21
kmalloc-128	1280	1280	128	32
kmalloc-96	1890	1890	96	42
kmalloc-64	3136	3136	64	64
kmalloc-32	3456	3456	32	128
kmalloc-16	2816	2816	16	256
kmalloc-8	3584	3584	8	512
\$				

num: Number of currently active objects

total: Total number of available objects

size: Size of each object

pages: Number of pages with at least one active object

Linux Kernel : Memory Management

Kernel Dynamic Memory Alloc APIs

- BSA (Buddy System Allocator)
- Slab
 - builtin APIs
 - custom slab cache

Linux Kernel : Memory Management

Kernel Dynamic Memory Alloc APIs

BSA - 'low-level' APIs

Most important of them is:

```
unsigned long __get_free_page(gfp_mask);
unsigned long __get_free_pages(gfp_t gfp_mask, unsigned int order);
unsigned long get_zeroed_page(gfp_t gfp_mask);
void *alloc_pages_exact(size_t size, gfp_t gfp_mask);
```

Free with

```
free_page(unsigned long addr);
void free_pages(unsigned long addr, unsigned int order);
void free_pages_exact(void *virt, size_t size);
```

Linux Kernel : Memory Management

Kernel Dynamic Memory Alloc APIs

A ‘golden rule’ :: ‘**cannot sleep in atomic context**’; hence, cannot call any API that results in a call to *schedule()* in atomic context (any kind of interrupt code is an atomic context: ISR top half, bottom halves (tasklets, softirqs))

GFP (get_free_page) Mask (gfp_t):

- a bitmask of several possible values
- most are used internally
- ‘module authors’ (driver devs) use these:

GFP_KERNEL : use when it’s safe to sleep - process-context *only*, no locks held; might cause the process context to block

GFP_ATOMIC : use when one *cannot* sleep; interrupt-context; high priority; will succeed or fail immediately (no blocking)

_GFP_ZERO : zeroes out the memory region

Linux Kernel : Memory Management

Kernel Dynamic Memory Alloc APIs

Which Flag to Use When

<i>Situation</i>	<i>Solution</i>
Process context, can sleep	Use GFP_KERNEL
Process context, cannot sleep	Use GFP_ATOMIC, or perform your allocations with GFP_KERNEL at an earlier or later point when you can sleep
Interrupt handler	Use GFP_ATOMIC
Softirq	Use GFP_ATOMIC
Tasklet	Use GFP_ATOMIC

Linux Kernel : Memory Management

Example *Low-level BSA*

drivers/hv/connection.c

```
...
vmbus_connection.int_page =
    (void *)__get_free_pages(GFP_KERNEL|__GFP_ZERO, 0);
    if (vmbus_connection.int_page == NULL) {
        ret = -ENOMEM;
        goto cleanup;
    }
...
free_pages((unsigned long)vmbus_connection.int_page, 0);
vmbus_connection.int_page = NULL;
...
```

Linux Kernel : Memory Management

Slab Allocator / Slab Cache APIs

- Layered above the BSA
- Gets *it's* memory from the BSA

APIs ::

kmalloc() or kzalloc()

- Vast majority of kernel allocations

- Meant for the usage case where the amount of memory required is less than a page; this is the common case!!

```
#include <linux/slab.h>
void *kmalloc(size_t size, gfp_t flags);
void *kzalloc(size_t size, gfp_t flags);
```

Free the memory with:

```
void kfree(const void *objp); // for both k[m|z]alloc()
void kzfree(const void *objp); // zero the mem & free it
                                (security)
```

Linux Kernel : Memory Management

Example *Slab allocator*

drivers/hv/connection.c

```
...
msginfo = kzalloc(sizeof(*msginfo) +
                  sizeof(struct vmbus_channel_initiate_contact),
                  GFP_KERNEL);
if (msginfo == NULL) {
    ret = -ENOMEM;
    goto cleanup;
}
...
kfree(msginfo);
```

An unwritten rule: GFP_KERNEL allocations for low-order
(\leq PAGE_ALLOC_COSTLY_ORDER (=3)) never fail (implies, \leq 8 pages)!

Linux Kernel : Memory Management

Runtime Usage of kmalloc

One can actually “track” these functions via the powerful **Kprobes** kernel framework! Below, output seen while using a **Jprobe and a Kretprobe on __kmalloc** :

[...]

```
jp_kmalloc: 40 callbacks suppressed
__kmalloc(96)
PRINT_CTX:: [000] systemd-journal :215  | d..0
__kmalloc(24)
PRINT_CTX:: [000][kworker/u2:2]:3529  | d..0
__kmalloc(152)
ret_handler: 40 callbacks suppressed
= 0xfffff95e9b8a3db40
PRINT_CTX:: [000][kthreadd]:2  | d..0
__kmalloc(32)
= 0xfffff95e9b8a457e0
```

[...]

Note- Jprobes removed from 4.15 onward

Linux Kernel : Memory Management

k[m|z]alloc Upper Limit on a Single Allocation

- From 2.6.22, the upper limit on x86, ARM, x86_64, etc is **4 MB**
- Technically, a function of the CPU page size and the value of MAX_ORDER
 - with a page size of 4 KB and MAX_ORDER of 10, the upper limit is 4 MB

A detailed article describing the same:

[*kmalloc and vmalloc : Linux kernel memory allocation API Limits*](#)

Linux Kernel : Memory Management

kmalloc Pros

- Uses the kernel's identity-mapped RAM (obtained and mapped at boot)
 - usually in ZONE_NORMAL
 - called the '**lowmem**' region
- it's thus **very fast** (pre-mapped, no page table setup required)
- returned memory chunks are guaranteed to be
 - *physically* contiguous
 - on a hardware cacheline boundary

kmalloc Limitations

- there is an upper limit
- security: by default, freed memory is not cleared, which could result in info-leakage scenarios. Can build the kernel with CONFIG_PAGE_POISONING (performance impact)

Linux Kernel : Memory Management

The *ksize* API, given a pointer to slab-allocated memory (in effect, memory allocated via the *kmalloc*), returns the actual number of bytes allocated.

```
size_t ksize(const void *objp);
```

It could be more than what was requested. E.g.

```
size_t actual=0;  
void *ptr = kmalloc(160, GFP_KERNEL); // ask for 160 bytes  
actual = ksize(ptr); // will get 192 bytes (actual ← 192)
```

Demo: the mm/allocsize_tests/ksz_test kernel module

Linux Kernel : Memory Management

vmalloc

- When you require more memory than kmalloc can provide (typically 4 MB)
- vmalloc provides a *virtually contiguous* memory region
- max size depends on the hardware platform and the kernel config; the kernel's VMALLOC region size - a global pool for all vmalloc's
- only for software use (f.e. not ok for DMA transfers); *only* process context (not in interrupt code); might cause process context to block
- slower (page table setup, demand paging based memory, unlike the kmalloc)

```
void *vmalloc(unsigned long size);
void *vzalloc(unsigned long size);
void vfree(const void *addr);
```

Linux Kernel : Memory Management

vmalloc: example code

fs/cramfs/uncompress.c

```
...
int cramfs_uncompress_init(void)           << see next slide >>
{
    if (!initialized++) {
        stream.workspace = vmalloc(zlib_inflate_workspacesize());
        if (!stream.workspace) {
            initialized = 0;
            return -ENOMEM;
        }
    }
...
}
```

arch/x86/kvm/cpuid.c

```
...
r = -ENOMEM;
cpuid_entries = vmalloc(sizeof(struct kvm_cpuid_entry) *
                        cpuid->nent);
if (!cpuid_entries)
    goto out;
...
}
```

Linux Kernel : Memory Management

Examine current vmalloc allocations via procfs (on an ARM)

```
ARM # cat /proc/vmallocinfo
```

```
0xa0800000-0xa0802000      8192 of_iomap+0x40/0x48 phys=1e001000 ioremap
0xa0802000-0xa0804000      8192 of_iomap+0x40/0x48 phys=1e000000 ioremap
0xa0804000-0xa0806000      8192 l2x0_of_init+0x68/0x234 phys=1e00a000 ioremap
0xa0806000-0xa0808000      8192 of_iomap+0x40/0x48 phys=10001000 ioremap
```

```
...
```

```
0xa088b000-0xa088d000      8192 devm_ioremap+0x48/0x7c phys=10016000 ioremap
0xa088d000-0xa0899000      49152 cramfs_uncompress_init+0x38/0x74 pages=11 vmalloc
0xa0899000-0xa08dc000      274432 jffs2_zlib_init+0x20/0x80 pages=66 vmalloc
```

```
...
```

```
#
```

Also:

When unsure, can use the **kvmalloc()**; will invoke kmalloc() as a first preference, and fallback to the **vmalloc()**:

```
void *kvmalloc(size_t size, gfp_t flags);
```

Linux Kernel : Memory Management

TIP: to get memory with particular protections

mm/vmalloc.c

```
...
 * For tight control over page level allocator and protection flags
 * use __vmalloc() instead.
...
```

```
void * __vmalloc(unsigned long size, gfp_t gfp_mask, pgprot_t prot)
```

arch/x86/include/asm/pgtable_types.h defines these protection symbols (and more):

```
...
PAGE_KERNEL_RO
PAGE_KERNEL_RX
PAGE_KERNEL_EXEC
PAGE_KERNEL_NOCACHE
PAGE_KERNEL_LARGE
PAGE_KERNEL_LARGE_EXEC
PAGE_KERNEL_IO
PAGE_KERNEL_IO_NOCACHE
...
```

Linux Kernel : Memory Management

Custom Slab Cache

- If a data structure in your kernel code is very often allocated and deallocated, it's perhaps a good candidate for a ***custom slab cache***
- Slab layer exposes APIs to allow custom cache creation and management
- **APIs**

```
struct kmem_cache *kmem_cache_create(const char *, size_t, size_t,  
                                     unsigned long,  
                                     void (*)(void *));  
void kmem_cache_destroy(struct kmem_cache *);
```

```
void *kmem_cache_[z]alloc(struct kmem_cache *cachep, gfp_t flags);  
void kmem_cache_free(struct kmem_cache *cachep, void *objp);
```

<< Custom Slab Cache code example >>

Linux Kernel : Memory Management

Custom Slab Cache

Misc

```
unsigned int kmem_cache_size(struct kmem_cache *s);
```

Cache Shrinker Interfaces

```
int kmem_cache_shrink(struct kmem_cache *cachep);
extern int register_shrinker(struct shrinker *);
extern void unregister_shrinker(struct shrinker *);
```

Linux Kernel : Memory Management

Managed memory for drivers with the devres APIs (1)

- Devres : *device resources* manager
- A framework developed for the device model, enabling auto-freeing of memory allocated via the ‘devres’ APIs
- The freeing occurs on driver detach
- *Requires* the struct device pointer though

<i>In place of</i>	Devres API
kmalloc / kzalloc	devm_kmalloc / devm_kzalloc
dma_alloc_coherent	dmam_alloc_coherent

Linux Kernel : Memory Management

Managed memory for drivers with the `devres` APIs (2)

- DON'T try and replace all allocations with the 'managed' variant APIs
- They are best suited to device driver `probe()` methods, ensuring correct initialization, and if it fails, auto-freeing of resources ensuring correct de-init
- Ref:
“... Please note that managed resources (whether it's memory or some other resource) are meant to be used in code responsible for the probing the device. They are generally a wrong choice for the code used for opening the device, as the device can be closed without being disconnected from the system. Closing the device requires freeing the resources manually, which defeats the purpose of managed resources.

The memory allocated with `kzalloc()` should be freed with `kfree()`. The memory allocated with `devm_kzalloc()` is freed automatically. It can be freed with `devm_kfree()`, but it's usually a sign that the managed memory allocation is not a good fit for the task. ...”

- Many more 'managed' interfaces available; see the kernel documentation here:
<https://www.kernel.org/doc/Documentation/driver-model/devres.txt>

Linux Kernel : Memory Management

Dynamic Memory Allocation APIs to Use

Required Memory	Method to Use	API(s)
Physically contiguous perfectly rounded power-of-2 pages (of size between 1 page and 4 MB)	BSA / page allocator	<code>alloc_page[s]</code> , <code>__get_free_page[s]</code> , <code>get_zeroed_page</code> , <code>alloc_pages_exact</code>
Common case: physically contiguous memory of size less than 1 page	Slab Allocator (cache)	<code>kmalloc</code> , <code>kzalloc</code> (<code>kzalloc()</code> preferred)
Physically contiguous memory (not a power-of-2) of size between 1 page and kmalloc max slab (typically 8 KB) or higher (upto 4 MB typically)	Slab Allocator (cache)	<code>kmalloc</code> , <code>kzalloc</code> : but check actual size allocated with <code>ksize()</code> and accordingly optimize
Virtually contiguous memory of large size (> 4 MB typically)	Indirect via BSA	<code>vmalloc</code> or CMA
Custom data structures	Custom slab caches	<code>kmem_cache_*</code> APIs

Linux Kernel : Memory Management

Kernel Segment / VAS

- Layout of the kernel virtual address space
- CPU-dependant
- <see “Kernel Memory Layout on ARM Linux” :
<http://www.arm.linux.org.uk/developer/memory.txt>

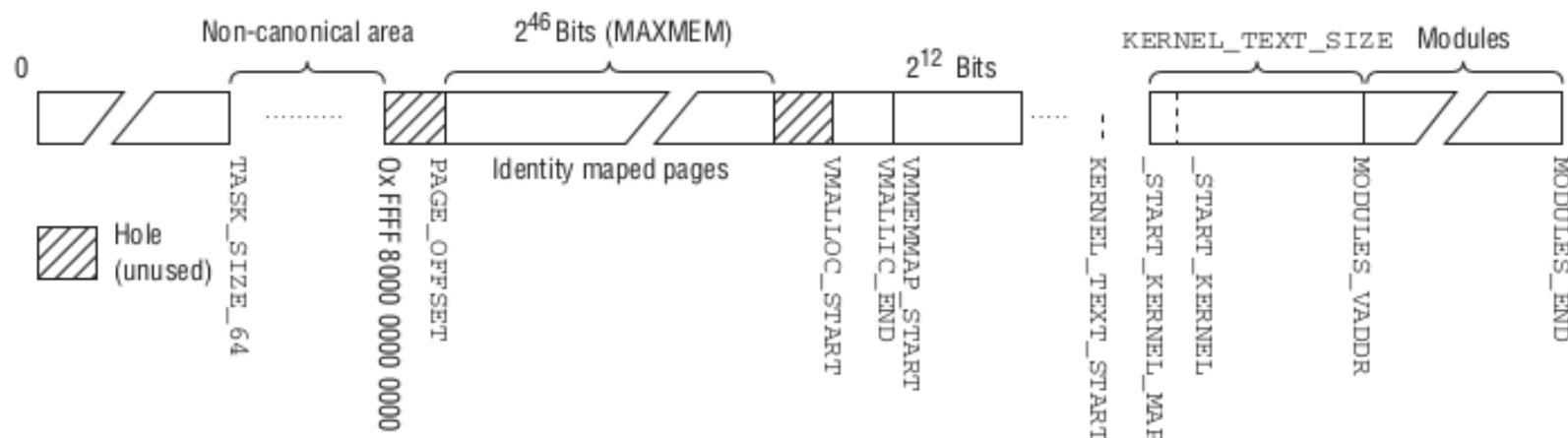
Eg.: On an ARM-32 Linux (Versatile Express CA-9 platform with 512 MB RAM):

```
$ dmesg
[...]
Virtual kernel memory layout:
  vector    : 0xfffff0000 - 0xfffff1000   (  4 kB)
  fixmap    : 0xfffc00000 - 0xfff00000   (3072 kB)
  vmalloc   : 0xa0800000 - 0xff800000   (1520 MB)
  lowmem   : 0x80000000 - 0xa0000000   ( 512 MB)
  modules   : 0x7f000000 - 0x80000000   ( 16 MB)
  .text     : 0x80008000 - 0x80800000   (8160 kB)
  .init     : 0x80b00000 - 0x80c00000   (1024 kB)
  .data     : 0x80c00000 - 0x80c664cc   ( 410 kB)
  .bss     : 0x80c68000 - 0x814724b8   (8234 kB)
```

Linux Kernel : Memory Management

Kernel Segment on the x86_64

```
#define __PAGE_OFFSET __AC(0xffff810000000000, UL)
#define PAGE_OFFSET __PAGE_OFFSET
#define MAXMEM __AC(0x3fffffff, UL)
```



Source: PLKA

Linux Kernel : Memory Management

Kernel Segment on a Yocto-emulated Aarch64 with 512 MB RAM

```
$ dmesg
[...]
[ 0.000000] Virtual kernel memory layout:
[ 0.000000]     modules   : 0xffffffff8000000000 - 0xffffffff8008000000 ( 128 MB)
[ 0.000000]     vmalloc   : 0xffffffff8008000000 - 0xffffffffbebffff0000 ( 250 GB)
[ 0.000000]         .text    : 0xffffffff8008080000 - 0xffffffff80086c0000 ( 6400 KB)
[ 0.000000]         .rodata   : 0xffffffff80086c0000 - 0xffffffff8008860000 ( 1664 KB)
[ 0.000000]         .init    : 0xffffffff8008860000 - 0xffffffff8008900000 ( 640 KB)
[ 0.000000]         .data    : 0xffffffff8008900000 - 0xffffffff80089bd800 ( 758 KB)
[ 0.000000]         .bss    : 0xffffffff80089bd800 - 0xffffffff8008a1969c ( 368 KB)
[ 0.000000]     fixed    : 0xffffffffbefef7fd000 - 0xffffffffbefec00000 ( 4108 KB)
[ 0.000000]     PCI I/O  : 0xffffffffbefee00000 - 0xffffffffbeffe00000 ( 16 MB)
[ 0.000000]     vmemmap  : 0xffffffffbf00000000 - 0xfffffffffc000000000 ( 4 GB maximum)
[ 0.000000]             : 0xffffffffbf00000000 - 0xffffffffbf008000000 ( 8 MB actual)
[ 0.000000]     memory   : 0xfffffffffc000000000 - 0xfffffffffc020000000 ( 512 MB)
[ 0.000000]                                     << 'lowmem' region - platform RAM >>
[...]
```

Linux Kernel : Memory Management

Kernel Segment on a Raspberry Pi 3 B+ Aarch64 with 1 GB RAM
running 64-bit Ubuntu 18.04.2 LTS

```
$ dmesg
[...]
Virtual kernel memory layout:
modules : 0xfffff000000000000 - 0xfffff000008000000 ( 128 MB)
vmalloc : 0xfffff000008000000 - 0xfffff7dffbf00000 (129022 GB)
  .text : 0x (ptrval) - 0x (ptrval) ( 11776 KB)
  .rodata : 0x (ptrval) - 0x (ptrval) ( 4096 KB)
  .init : 0x (ptrval) - 0x (ptrval) ( 6144 KB)
  .data : 0x (ptrval) - 0x (ptrval) ( 1343 KB)
  .bss : 0x (ptrval) - 0x (ptrval) ( 1003 KB)
fixed : 0xfffff7dfffe7fb000 - 0xfffff7dfffec00000 ( 4116 KB)
PCI I/O : 0xfffff7dffffe00000 - 0xfffff7dfffffe00000 ( 16 MB)
vmemmap : 0xfffff7e00000000000 - 0xfffff800000000000 ( 2048 GB maximum)
          0xfffff7e4d36000000 - 0xfffff7e4d36ed0000 ( 14 MB actual)
memory : 0xfffff934d80000000 - 0xfffff934dbb400000 ( 948 MB)
          << 'lowmem' region - platform RAM >>
[...]
```

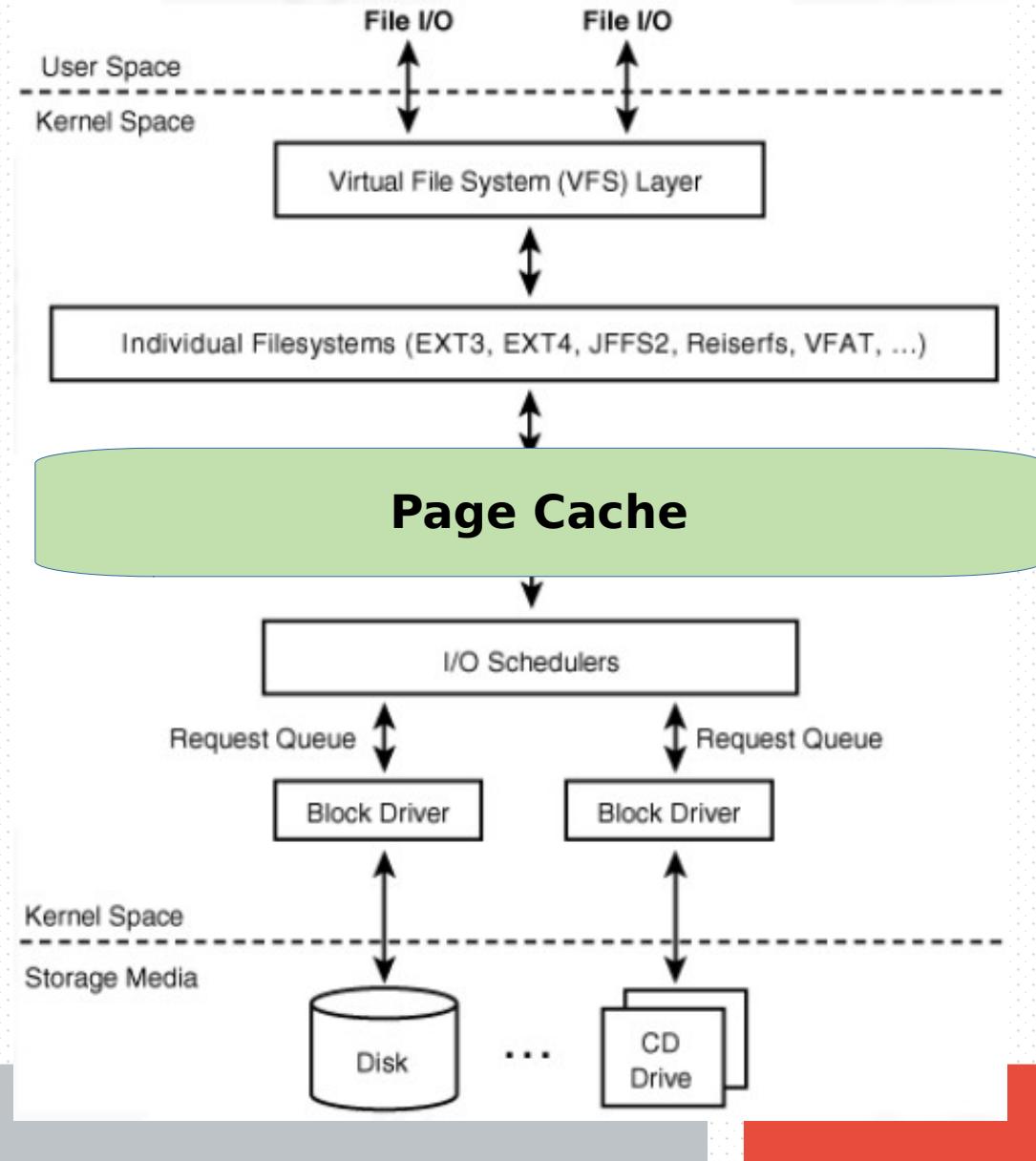
Linux Kernel : Memory Management

The Page Cache

Fundamental property of a cache:

(a) the memory medium of the memory that is being cached is *much* slower to access than the cache

(b) principle of locality of reference: spatial (in space) and temporal (in time)



Linux Kernel : Memory Management

Page Cache

Bypass the page cache by using O_DIRECT in open(2)

Clean - empty out - the page cache:

```
sync && echo 1 > /proc/sys/vm/drop_caches
```

```
# free -h
```

	total	used	free	shared	buff/cache	available
Mem:	15G	9.1G	290M	975M	6.2G	5.3G
Swap:	7.6G	0B	7.6G			
Total:	23G	9.1G	7.9G			

```
# sync
```

```
# echo 1 > /proc/sys/vm/drop_caches
```

```
# free -h
```

	total	used	free	shared	buff/cache	available
Mem:	15G	9.1G	4.8G	975M	1.7G	5.3G
Swap:	7.6G	0B	7.6G			
Total:	23G	9.1G	12G			

Linux Kernel : Memory Management

Page Cache

The article

["Page Cache , the Affair Between Memory and Files"](#), Gustav Duarte, clearly demonstrates why using memory-mapped IO is far superior to the “usual” approach – using the read/write system calls (or library equivalents) in a loop – especially for the use-case where the amount of IO to be performed is quite large.

Linux Kernel : Memory Management

The OOM (Out Of Memory) Killer and Demand Paging

- Where's the VM?
 - In the 'memory pyramid'
 - Registers, CPU Caches, RAM, Swap, (+ newer: nvdimm's)
- What if - worst case scenario - *all of these are completely full!?*
- *Then the OOM Killer jumps in, and*
 - Kills the process (and descendants) with highest memory usage
 - Uses heuristics to determine the target process

Linux Kernel : Memory Management

The OOM (Out Of Memory) Killer and Demand Paging

Then the OOM Killer jumps in, and

- Kills the process (and descendants) with highest memory usage
- Uses heuristics to determine the target process
- OOM Score (in /proc/<pid>/oom_score)
- Range: [0-1000]
 - 0 : the process is not using any memory available to it
 - 1000 : the process is using 100% of memory available to it
- oom_score_adj
 - **Net score = oom_score + oom_score_adj**
 - So: oom_score_adj = -1000 => *never kill this*
 - oom_score_adj = +1000 => *always kill this*
- *Details: see the man page on proc(5) - search for oom_**

Linux Kernel : Memory Management

The OOM (Out Of Memory) Killer and Demand Paging

A script to display the OOM ‘score’ and ‘adjustment’ values:

```
# Src: https://dev.to/rampage/surviving-the-linux-oom-killer-2ki9
printf 'PID\tOOM Score\tOOM Adj\tCommand\n'
while read -r pid comm
do
    [ -f /proc/$pid/oom_score ] && [ $(cat /proc/$pid/oom_score) != 0 ] &&
        printf '%d\t%d\t%d\t%s\n' "$pid" "$(cat /proc/$pid/oom_score)" "$(cat /proc/$pid/oom_score_adj)" "$comm"
done < <(ps -e -o pid= -o comm=) | sort -k 2nr
```

```
$ ./show_oom_score
```

PID	OOM Score	OOM Adj	Command
23584	93	0	VirtualBoxVM
3176	49	0	Web Content
9540	39	0	Web Content
3906	37	0	Web Content
3115	28	0	firefox
3428	23	0	Web Content
[...]			

```
$
```

Linux Kernel : Memory Management

The OOM (Out Of Memory) Killer and Demand Paging

- Can one deliberately invoke the OOM killer?
 - Easy!

```
# echo f > /proc/sysrq-trigger
```

- “Manually”: write a “crazy allocator” ‘C’ program
 - *Check out the **oom_killer_try.c** program now*

Linux Kernel : Memory Management

Demand paging

VM overcommit

Page fault handling basics

Glibc malloc behavior

User VAS mapping - proc

Linux Kernel : Memory Management