

MM101: Introduction to Memory Management

Christopher Lameter <cl@linux.com>

@qant



Overview

- Memory and processes
- Real/Virtual memory and Paging
- Machine configuration
- Processes use of memory
- Overcommit
- Knobs
- Processor cache use





Pages and physical page frame numbers

- Division of memory into "pages"
 - 1-N bytes become split at page size boundaries and become
 M = N/page-size

pages

- We can then refer to memory by the Page Frame Number (PFN) and an offset into the page.
- Common size is 4k (Intel legacy issues)
- The MMU "creates" virtual addresses.





Basics of "paging"

Paging Henny Henny Virtual Page Number Offset Prime Offset

- Processes have virtual memory
- -> PFN
- Page Tables
- Faults
 - Major
 - Minor
- Virtual vs physical



Process Memory

- Virtual memory maps to physical memory
- A view of memory distinct
 - for each process.
- Pages shared
- Access control
- Copy on Write





Swap, Zero pages etc.

- Swap page
- Zero page
- Read data behavior
- Write data behavior

l see a number zero!



Anonymous vs file backed pages



Kernel Basic memory information

/proc/meminfo

/sys/devices/system/ has lots of more detailed information on hardware (processors and memory)

Commands: numactl --hardware free, top, dmesg

MemTotal: MemFree: MemAvailable: 30823580 kB Buffers: Cached: SwapCached: Active: Inactive: Active(anon): Inactive(anon): Active(file): Inactive(file): Unevictable: Mlocked: SwapTotal: SwapFree: Dirty: Writeback:

31798552 kB 25949124 kB 220988 kB 4679188 kB 0 kB2803000 kB 2336992 kB 240776 kB 6432 kB 2562224 kB 2330560 kB 0 kB 0 kB 2097148 kB 2097148 kB 48 kB 0 kB

AnonPages: Mapped: Shmem: Slab: SReclaimable: SUnreclaim: KernelStack: PageTables:

239716 kB 195596 kB 7396 kB 550628 kB 443040 kB 107588 kB 6840 kB 11176 kB



Inspecting a processes use of memory

/proc/<pid>/status /proc/<pid>/*maps

(there are other files in /proc/<pid>/* with more information about the processes)

Commands: **ps, top**

Name: sshd VmPeak: VmSize: VmLck: VmPin: VmHWM: VmRSS: RssAnon: RssFile: RssShmem:

65772 kB 65772 kB $0 \, \text{kB}$ $0 \, \text{kB}$ 6008 kB 6008 kB 1216 kB 4792 kB $0 \, \text{kB}$



V

V

V

V

V

1332 kB
132 kB
492 kB
8076 kB
168 kB
0 kB



User limit (ulimit)

- Max memory size
- Virtual memory
 - Stack size
- and lots of other controls.



Overcommit configuration

- Virtual memory use vs physical
- overcommit_kbytes overcommit_memory
 - 0 overcommit. Guess if mem is available.
 - 1 Overcommit. Never say there is no memory
 - 2 Only allocate according to the ratio

```
overcommit_ratio
```







Important VM control knobs

Found in **/proc/sys/vm**

More descriptions of these knobs in Kernel source code.

linux/Documentation/admin-guide

admin reserve kbytes dirty writeback centisecs min free kbytes numa zonelist order stat refresh block dump drop caches min slab ratio oom dump tasks swappiness compact memory extfrag threshold min unmapped ratio oom kill allocating task user reserve kbytes compact unevictable allowed hugetlb shm group mmap min addr overcommit kbytes vfs cache pressure dirty background bytes laptop mode mmap rnd bits overcommit memory watermark scale factor dirty_background_ratio legacy va layout mmap rnd compat bits overcommit ratio zone reclaim mode dirty bytes lowmem reserve ratio nr hugepages page-cluster dirty expire centisecs max map count nr hugepages mempolicy panic on oom dirty ratio memory failure early kill nr overcommit hugepages percpu pagelist fraction dirtytime_expire_seconds memory_failure_recovery numa_stat stat_interval



• The online Kernel Administrators Guide:

https://www.kernel.org/doc/html/v4.14/admin-guide/index. html

- Kernel.org has wikis and documentation (www.kernel.org)
- Consult the manpages (especially for system calls and coding)

"Simple" Memory Access

- **UMA** (Uniform Memory Access)
- Any access to memory has the same characteristics (performance and latency)
- The vast major of systems have only UMA.
- But there is always the processor cache hierarchy
 - The CPU is fast, memory is slow
 - Caches exist to avoid accesses to r memory
- Aliasing
- Coloring
- Cache Miss
- Trashing



CPU Cache Access Latencies in Clock Cycles



NUMA Memory

- Memory with different access characteristics
- Memory Affinities depending on where a process was started
- Control NUMA allocs with memory policies
- System Partitioning using Cpusets and Containers
- Manual memory *migration*
- Automatic memory migration



- Typical memory is handled in chunks of base page size (Intel 4k, IBM PowerX 64K, ARM 64K)
- Systems support larger memory chunks of memory called Huge pages (Intel 2M)
- Must be pre configured on boot in order to guarantee that they are available
- Required often for I/O bottlenecks on Intel.
- 4TB requires 1 billion descriptors with 4K pages. Most of this is needed to compensate for architectural problems on Intel. Intel processors have difficulties using modern SSDs and high
 Default
 Page 24
 Page 24
- Large contiguous segments (I/O performance)
- Fragmentation issues
- Uses files on a special file system that must be explicitly requested by mmap operations from special files.



Default	HugePages pool
72 GB RAM	24 GB RAM HugePages
4K 4K<	4K 4K 4K 4K 4K 2 MB 4K 4K 4K 4K 4K 4K 4K 4K 2 MB 4K 4K 2 MB 2 MB 2 MB



With

For questions and feedback please reach out to me at - cl@linux.com

http://gentwo.org/christoph



THE LINUX FOUNDATION OPEN SOURCE SUMMIT NORTH AMERICA