

---

# Chapter 13 Memory: Monitoring Usage and Tuning - Notes

---

## 13.2 Introduction

---

Over time, systems more demanding of memory resources while RAM prices decreased and performance improved. Yet, bottlenecks in overall system performance often memory-related. CPUs and I/O subsystem can be waiting for data to be retrieved from/written to memory. Many tools for monitoring, debugging, tuning system's behavior with regard to its memory.

## 13.3 Learning Objectives:

---

- List the primary (inter-related) considerations and tasks involved in memory tuning.
- Use entries in `/proc/sys/vm` and decipher `/proc/meminfo`.
- Use **vmstat** to display information about memory, paging, I/O, processor activity, and processes' memory consumption.
- Understand how the **OOM-killer** decides when to take action and selects which processes should be exterminated to open up some memory.

## 13.4 Memory Tuning Considerations

---

Tuning memory sub-system can be complex process. Have to take note that memory usage and I/O throughput are intrinsically related, since in most cases most memory being used to cache contents of files on disk.

Thus, changing memory parameters -> large effect on I/O performance. Changing I/O parameters -> equally large converse effect on virtual memory sub-system.

When tweaking parameters in `/proc/sys/vm`, usual best practice to adjust one thing at a time and look for effects. Primary (inter-related) tasks:

- Controlling flush parameters; ie, how many pages allowed to be **dirty** and how often they flushed out to disk
- Controlling swap behavior; ie, how much pages that reflect file contents allowed to remain in memory, as opposed to those that need to be swapped out as they have no other backing store
- Controlling how much memory **overcommitment** is allowed, since many programs never need full amount of memory they request, particularly because of **copy on write (COW)** techniques

Memory tuning often subtle. What works in one system situation or load may be far from optimal in other circumstances.

## 13.5 Memory Monitoring Tools

---

Some important basic tools for monitoring and tuning memory in Linux:

### Memory Monitoring Utilities

Utility	Purpose	Package
<b>free</b>	Brief summary of memory usage	<b>procps</b>
<b>vmstat</b>	Detailed virtual memory statistics and block I/O, dynamically updated	<b>procps</b>

Simplest tool to use is **free**:

```
c7:/tmp> free -m
              total    used    free   shared  buff/cache   available
Mem:      15895  2946   6085    3580     6863       8984
Swap:      20023     0    20023
```

## 13.6 /proc/sys/vm

`/proc/sys/vm` directory contains many tunable knobs to control **Virtual Memory** system. Exactly what appears in directory depends somewhat on kernel version. Almost all entries writable (by **root**).

Note: values can be changed either by directly writing to entry, or using **sysctl** utility. Values can be set at boot time by modifying `/etc/sysctl.conf`.

Can find full documentation for `/proc/sys/vm` directory in kernel source (or kernel documentation package on distribution) usually under `Documentation/sysctl/vm.txt`.

### `proc/sys/vm` Entries

Entry	Purpose
<code>admin_reserve_kbytes</code>	Amount of free memory reserved for privileged users
<code>block_dump</code>	Enables block I/O debugging
<code>compact_memory</code>	Turns on or off <b>memory compaction</b> (essentially defragmentation) when configured into the kernel
<code>dirty_background_bytes</code>	<b>Dirty</b> memory threshold that triggers writing uncommitted pages to disk
<code>dirty_background_ratio</code>	Percentage of total pages at which kernel will start writing dirty data out to disk
<code>dirty_bytes</code>	The amount of dirty memory a process needs to initiate writing on its own
<code>dirty_expire_centisecs</code>	When dirty data is old enough to be written out in hundredths of a second
<code>dirty_ratio</code>	Percentage of pages at which a process writing will start writing out dirty data on its own
<code>dirty_writeback_centisecs</code>	Interval in which periodic writeback daemons wake up to flush. If set to zero, no automatic periodic writeback
<code>drop_caches</code>	Echo 1 to free page cache, 2 to free dentry and inode caches, 3 to free all. note only <b>clean</b> cached pages are dropped; do <b>sync</b> first to flush dirty pages
<code>extfrag_threshold</code>	Controls when kernel should compact memory
<code>hugepages_treat_as_movable</code>	Used to toggle how huge pages are treated
<code>hugetlb_shm_group</code>	Sets a group ID that can be used for System V huge pages
<code>laptop_mode</code>	Can control a number of features to save power on laptops
<code>legacy_va_layout</code>	Use old layout (2.4 kernel) for how memory mappings are displayed

<code>lowmem_reserve_ratio</code>	Controls how much low memory is reserved for pages that can only be there; ie, pages which can go in high memory instead will do so. Only important on 32-bit systems with high memory
<code>max_map_count</code>	Maximum number of memory mapped areas a process may have. Default is 64 K
<code>min_free_kbytes</code>	Minimum free memory that must be reserved in each zone
<code>mmap_min_addr</code>	How much address space a user process cannot memory map. Used for security purposes, to avoid bugs where accidental kernel null dereferences can overwrite the first pages used in an application
<code>nr_hugepages</code>	Minimum size of hugepage pool
<code>nr_pdflush_hugepages</code>	Maximum size of the hugepage pool = <code>nr_hugepages*nr_overcommit_hugepages</code>
<code>nr_pdflush_threads</code>	Current number of <b>pdflush</b> threads; not writeable
<code>oom_dump_tasks</code>	If enabled, dump information produced when <b>oom-killer</b> cuts in
<code>oom-kill-allocating_task</code>	If set, <b>oom-killer</b> kills task that triggered out-of-memory situation, rather than trying to select best one
<code>overcommit_kbytes</code>	One can set either <code>overcommit_ratio</code> or this entry, not both
<code>overcommit_memory</code>	If 0, kernel estimates how much free memory is left when allocations are made. If 1, permits all allocations until memory actually does run out. If 2, prevents any overcommitment
<code>overcommit_ratio</code>	If <code>overcommit_memory</code> = 2 memory commitment can reach swap plus this percentage of <b>RAM</b>
<code>page-cluster</code>	Number of pages that can be written to swap at once, as a power of two. Default is 3 (which means 8 pages)
<code>panic_on_oom</code>	Enable system to crash on an out of memory situation
<code>percpu_pagelist_fraction</code>	Fraction of pages allocated for each cpu in each zone for hot_pluggable CPU machines
<code>scan_unevictable_pages</code>	If written to, system will scan and try to move pages to try and make them reclaimable
<code>stat_interval</code>	How often vm statistics are updated (default 1 second) by <b>vmstat</b>
<code>swappiness</code>	How aggressively should the kernel swap
<code>user_reserve_kbytes</code>	If <code>overcommit_memory</code> is set to 2 this sets how low the user can draw memory resources
<code>vfs_cache_pressure</code>	How aggressively the kernel should reclaim memory used for inode and dentry cache. Default is 100; if 0 this memory is never reclaimed due to memory pressure

## 13.7 vmstat

**vmstat**: multi-purpose tool that displays information about memory, paging, I/O, processor activity and processes. Has many options. General form of command:

```
$ vmstat [options] [delay] [count]
```

If **delay** given in seconds, report repeated at interval count times. If **count** not given, **vmstat** will keep reporting statistics

forever, until killed by signal, such as `Ctrl-C`.

If no other arguments given, can see what **vmstat** displays, where first line shows averages since last reboot, while succeeding lines show activity during specified interval.

```
$ vmstat 2 4
```

vmstat

Fields shown are:

#### vmstat Fields

Field	Subfield	Meaning
Processes	r	Number of processes waiting to be scheduled in
Processes	b	Number of processes in uninterruptible sleep
memory	swpd	Virtual memory used (KB)
memory	free	Free (idle) memory (KB)
memory	buff	Buffer memory (KB)
memory	cache	Cached memory (KB)
sw ap	si	Memory swapped in (KB)
sw ap	so	Memory swapped out (KB)
I/O	bi	Blocks read from devices (block/sec)
I/O	bo	Blocks written to devices (block/sec)
system	in	Interrupts/second
system	cs	Context switches/second
CPU	us	CPU time running user code (percentage)
CPU	sy	CPU time running kernel (system) code (percentage)
CPU	id	CPU time idle (percentage)
CPU	wa	Time waiting for I/O (percentage)
CPU	st	Time "stolen" from virtual machine (percentage)

If option `-s m` given, memory statistics will be in MB instead of KB.

With `-a` option, **vmstat** displays information about **active** and **inactive** memory, where **active** memory pages are those which have been recently used. May be **clean** (disk contents are up to date) or **dirty** (need to be flushed to disk eventually). By contrast, **inactive memory** pages have not been recently used and are more likely to be clean and are released sooner under memory pressure:

```
$ vmstat -a 2 4
```

Memory can move back and forth between active and inactive lists, as they get newly referenced, or go a long time between uses.

```
vmstat
```

To get table of memory statistics and certain event counters, use `-s` option:

```
vmstats
```

To get table of disk statistics, use `-d` option:

```
vmstatd
```

#### vmstat Disk Fields

Field	Subfield	Meaning
reads	total	Total reads completed successfully
reads	merged	Grouped reads (resulting in one I/O)
reads	ms	Milliseconds spend reading
writes	total	total writes completed successfully
writes	merged	Grouped writes (resulting in one I/O)
writes	ms	Milliseconds spent writing
I/O	cur	I/O in progress
I/O	sec	seconds spent for I/O

If want to just get some quick statistics on only one partition, use `-p` option:

```
vmstatp
```

## 13.8 /proc/meminfo

As noted earlier, relatively lengthy summary of memory statistics in `/proc/meminfo`:

```
procmeminfo
```

Worthwhile to go through listing and understand most of the entries:

`/proc/meminfo` **Entries**

Entry	Meaning
<code>MemTotal</code>	Total usable RAM (physical minus some kernel reserved memory)
<code>MemFree</code>	Free memory in both low and high zones
<code>Buffers</code>	Memory used for temporary block I/O storage
<code>Cached</code>	Page cache memory, mostly for file I/O

SwapCached	Memory that was swapped back in but is still in the swap file
Active	Recently used memory, not to be reclaimed first
Inactive	Memory not recently used, more eligible for reclamation
Active (anon)	Active memory for <b>anonymous</b> pages
Inactive (anon)	Inactive memory for <b>anonymous</b> pages
Active (file)	Active memory for <b>file</b> -backed pages
Inactive (file)	Inactive memory for <b>file</b> -backed pages
Unevictable	Pages which can not be swapped out of memory or released
Mlocked	Pages which are locked in memory
SwapTotal	Total swap space available
SwapFree	Swap space not being used
Dirty	Memory which needs to be written back to disk
Writeback	Memory actively being written back to disk
AnonPages	Non-file back pages in cache
Mapped	Memory mapped pages, such as libraries
Shmem	Pages used for shared memory
Slab	Memory used in slabs
SReclaimable	Cached memory in slabs that can be reclaimed
SUnreclaim	Memory in slabs that can't be reclaimed
KernelStack	Memory used in kernel stack
PageTables	Memory being used by page table structures
Bounce	Memory used for block device bounce buffers
WritebackTmp	Memory used by <b>FUSE</b> filesystems for writeback buffers
CommitLimit	Total memory available to be used, including overcommitment
Committed_AS	Total memory presently allocated, whether or not it is used
VmallocTotal	Total memory available in kernel for <b>vmalloc</b> allocations
VmallocUsed	Memory actually used by <b>vmalloc</b> allocations
VmallocChunk	Largest possible contiguous <b>vmalloc</b> area
HugePages_Total	Total size of the huge page pool
HugePages_Free	Huge pages that are not yet allocated
HugePage_Rsvd	Huge pages that have been reserved, but not yet used
HugePages_Surp	Huge pages that are surplus, used for overcommitment

Note: exact entries seen may depend on exact kernel version being run.

## 13.9 OOM Killer

Simplest way to deal with memory pressure -> permit memory allocations to succeed as long as free memory is available, fail when all memory exhausted.

Second simplest way -> use **swap** space on disk to push some resident memory out of core. In this case, total available memory (at least in theory) is actual RAM plus size of **swap** space. Hard part of this is to figure out which pages of memory to swap out when pressure demands. In this approach, once swap space filled, requests for new memory must fail.

Linux, however, goes one better. Permits system to overcommit memory, so that it can grant memory requests that exceed size of RAM plus **swap**. Might seem foolhardy, but many (if not most) processes do not use all requested memory.

Example 1: program that allocates 1 MB buffer, and then uses only few pages of memory. Example 2: every time child process forked, receives copy of entire memory space of parent. Because Linux uses COW (copy on write) technique, unless one of the processes modifies memory, no actual copy needs to be made. However, kernel has to assume that copy might need to be done.

Thus, kernel permits overcommitment of memory, but only for pages dedicated to user processes. Pages used within kernel not swappable, and always allocated at request time.

Can modify, and even turn off this overcommitment by setting value of `/proc/sys/vm/overcommit_memory`:

- 0: (default) Permit overcommitment, but refuse obvious overcommits, and give root users somewhat more memory allocation than normal users
- 1: All memory requests are allowed to overcommit
- 2: Turn off overcommitment. Memory requests will fail when the total memory commit reaches the size of the **swap** space plus a configurable percentage (50 by default) of RAM. This factor is modified changing `/proc/sys/vm/overcommit_ratio`.

If available memory exhausted, Linux invokes **OOM-killer** (**Out Of Memory**) to decide which process(es) to exterminate to open up memory.

No precise science, algorithm must be heuristic, cannot satisfy everyone. In minds of many developers, purpose of OOM-killer to permit graceful shutdown, rather than be part of normal operations.

An amusing take on this by Andries Brouwer (<https://lwn.net/Articles/104185/>):

An aircraft company discovered that it was cheaper to fly its planes with less fuel on board. The planes would be lighter and use less fuel and money was saved. On rare occasions however the amount of fuel was insufficient, and the plane would crash. This problem was solved by the engineers of the company by the development of a special OOF (out-of-fuel) mechanism. In emergency cases a passenger was selected and thrown out of the plane. (When necessary, the procedure was repeated.) A large body of theory was developed and many publications were devoted to the problem of properly selecting the victim to be ejected. Should the victim be chosen at random? Or should one choose the heaviest person? Or the oldest? Should passengers pay in order not to be ejected, so that the victim would be the poorest on board? And if for example the heaviest person was chosen, should there be a special exception in case that was the pilot? Should first class passengers be exempted? Now that the OOF mechanism existed, it would be activated every now and then, and eject passengers even when there was no fuel shortage. The engineers are still studying precisely how

this malfunction is caused.

In order to make decisions of who gets sacrificed to keep system alive, value called **badness** computed (can be read from `/proc/[pid]/oom_score`) for each process on system and order of killing determined by this value.

Two entries in same directory can be used to promote/demote likelihood of extermination. Value of `oom_adj` : number of bits points should be adjusted by. Normal users can only increase badness. Decrease (negative value for `oom_adj`) can only be specified by superuser. Value of `oom_adj_score` directly adjusts point value. Note: use of `oom_adj` deprecated.

##

[Back to top](#)

---

[Previous Chapter](#) - [Table of Contents](#) - [Next Chapter](#)