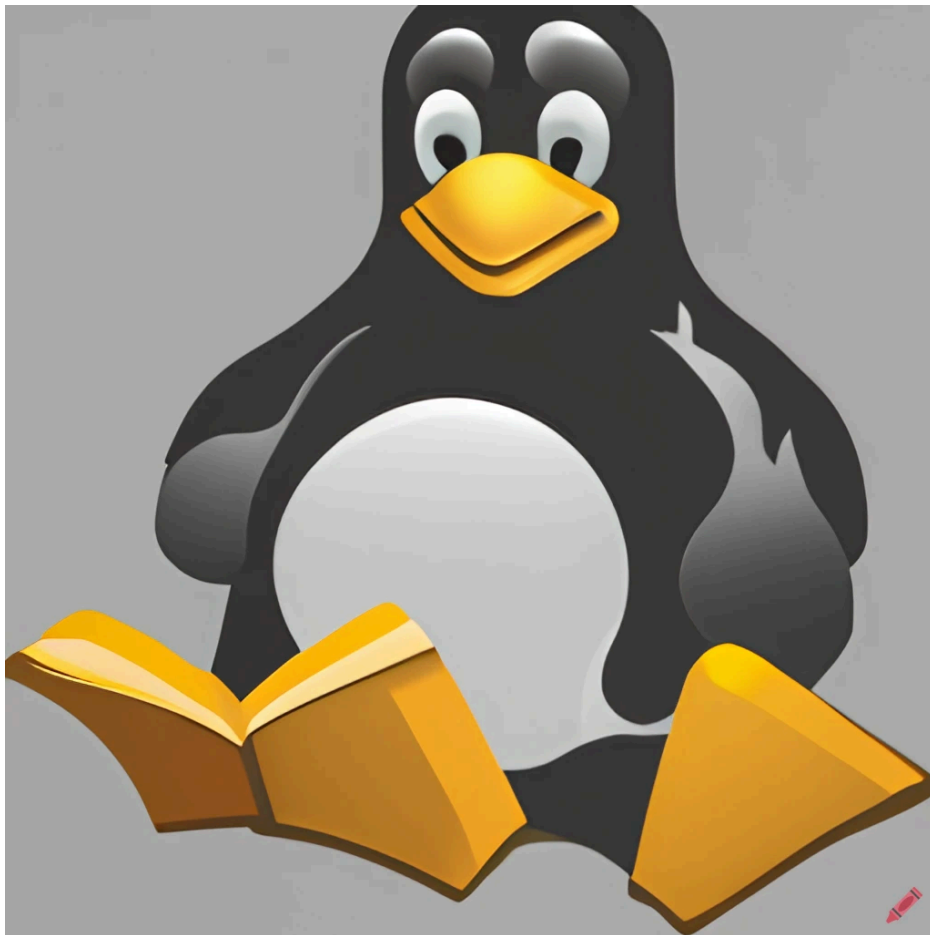


The Linux Process Journey

Version 9.0

June-2024

By Dr. Shlomi Boutnaru



Created using [Craiyon AI Image Generator](#)

Table of Contents

Table of Contents.....	2
Introduction.....	6
swapper (PID 0).....	7
init (PID 1).....	8
Kernel Threads.....	9
kthreadd (PID 2).....	10
migration.....	12
charger_manager.....	14
idle_inject.....	15
kworker (Kernel Thread Worker).....	16
kdevtmpfs.....	17
cpuhp (CPU Hotplug).....	18
khungtaskd (Kernel Hang Task Daemon).....	18
kswapd.....	20
kcompactd.....	21
md (Multiple Device Driver).....	23
mld (Multicast Listener Discovery).....	25
ksmd (Kernel Same Page Merging).....	26
ttm_swap.....	27
watchdogd (Watchdog Daemon).....	28
zswap-shrink.....	30
khugepaged (Kernel Huge Pages Daemon).....	31
krfcommd (Kernel Radio Frequency Communication Daemon).....	32
ksgxd (Kernel Software Guard eXTensions Daemon).....	33
jbd2 (Journal Block Device 2).....	34
netns (Network Namespace).....	35
oom_reaper (Out-of-Memory Reaper).....	36
kpsmoused (Kernel PS/2 Mouse Daemon).....	37
Slub_flushwq (SLUB Flush Work Queue).....	38
pgdatinit.....	39
kblockd (Kernel Block Daemon).....	40
writeback.....	41
kdamond (Data Access MONitor).....	42
kintegrityd (Kernel Integrity Daemon).....	43
kthrotld (Kernel Throttling Daemon).....	44
scsi_eh (Small Computer System Interface Error Handling).....	45
blkcg_punt_bio.....	46
napi (New API).....	47

kauditd (Kernel Audit Daemon)	48
tpm_dev_wq (Trusted Platform Module Device Work Queue)	49
ipv6_addrconf (IPv6 Address Auto Configuration)	50
mm_percpu_wq (Per-CPU Memory Work Queue)	52
inet_frag_wq (IP Fragmentation Work Queue)	53
kstrp (Stream Parser)	54
devfreq_wq	55
dmcrypt_write (Device Mapper for Transparent Encryption/Decryption)	56
ModemManager (Modem Management Daemon)	57
kerneloops	58
xargs (Extended Arguments)	59
cpp (The C Preprocessor)	60
ntpd (Network Time Protocol Daemon)	61
gold (The GNU ELF Linker)	62
kmod (Linux Kernel Module Handling)	63
cfg80211 (Wireless Configuration)	64
kdmflush (Kernel Device Mapper Flush)	65
nvme-wq (Non-Volatile Memory Express Work Queue)	66
] (Checking File Types and Comparing Values)	67
rcub (Read-Copy Update Boost)	68
dmesg (Print/Control the Kernel Ring Buffer)	69
login (Begin Session on The System)	70
su (Substitute\Switch User)	71
kacpid (Kernel Advanced Configuration and Power Interface Daemon)	72
thermald (Thermal Daemon)	73
dhcpd (Dynamic Host Configuration Protocol Daemon)	74

Introduction

When starting to learn OS internals I believe that we must understand the default processes executing (roles, tasks, etc). Because of that I have decided to write a series of short writeups named "Process ID Card" (aimed at providing the OS vocabulary).

Overall, I wanted to create something that will improve the overall knowledge of Linux in writeups that can be read in 1-3 mins. I hope you are going to enjoy the ride.

In order to create the list of processes I want to explain, I have installed a clean Ubuntu 22.10 VM (Desktop version) and executed ps (as can be seen in the following image - not all the output was included).

```
UID          PID     PPID  C  STIME TTY          TIME CMD
root         1       0  1  07:15 ?           00:00:05 /lib/systemd/systemd splash --system --deserialize 26
root         2       0  0  07:15 ?           00:00:00 [kthreadd]
root         3       2  0  07:15 ?           00:00:00 [rcu_gp]
root         4       2  0  07:15 ?           00:00:00 [rcu_par_gp]
root         5       2  0  07:15 ?           00:00:00 [kworker/0:0-events]
root         6       2  0  07:15 ?           00:00:00 [kworker/0:0H-events_highpri]
root         9       2  0  07:15 ?           00:00:00 [mm_percpu_wq]
root        10       2  0  07:15 ?           00:00:00 [rcu_tasks_rude_]
root        11       2  0  07:15 ?           00:00:00 [rcu_tasks_trace]
root        12       2  0  07:15 ?           00:00:00 [ksoftirqd/0]
```

Probably the best way to do it is to go over the processes by the order of their PID value. The first one I want to talk about is the one we can't see on the list, that is PID 0 (we can see it is the PPID for PID 1 and PID 2 - on them in the next posts).

Lastly, you can follow me on twitter - @boutnaru (<https://twitter.com/boutnaru>). Also, you can read my other writeups on medium - <https://medium.com/@boutnaru>. Lastly, You can find my free eBooks at <https://TheLearningJourneyEbooks.com>.

Lets GO!!!!!!

swapper (PID 0)

Historically, old Unix systems used swapping and not demand paging. So, swapper was responsible for the “Swap Process” - moving all pages of a specific process from/to memory/backing store (including related process’ kernel data structures). In the case of Linux PID 0 was used as the “idle process”, simply does not do anything (like nops). It was there so Linux will always have something that a CPU can execute (for cases that a CPU can’t be stopped to save power). By the way, the idle syscall is not supported since kernel 2.3.13 (for more info check out “man 2 idle”). So what is the current purpose of swapper today? helping with pageout ? cache flushes? idling? buffer zeroing? I promise we will answer it in more detail while going through the other processes and explaining the relationship between them.

But how can you believe that swapper (PID 0) even exists? if you can’t see it using ps. I am going to use “bpftrace” for demonstrating that (if you don’t know about bpftrace, I strongly encourage you to read about it). In the demo I am going to trace the kernel function “hrtimer_wakeup” which is responsible for waking up a process and move it to the set of runnable processes. During the trace I am going to print the pid of the calling process (BTW, in the kernel everything is called a task - more on that in future posts) and the executable name (the comm field of the task_struct [/include/linux/sched.h]). Here is the command: `sudo bpftrace -e 'kfunc:hrtimer_wakeup { printf("%s:%d\n",curtask->comm,curtask->pid); }'`.



```
Attaching 1 probe...
swapper/0:0
swapper/2:0
swapper/0:0
swapper/2:0
swapper/2:0
swapper/0:0
swapper/2:0
swapper/0:0
swapper/2:0
swapper/0:0
```

From the output we can see we have 3 instances of swapper: swapper/0, swapper/1 and swapper/2 all of them with PID 0. The reason we have three is because my VM has 3 virtual CPUs and there is a swapper process for each one of them - see the output of the command in the following image.

init (PID 1)

After explaining about PID 0, now we are going to talk about PID 1. Mostly known as “init”. init is the first Linux user-mode process created, which runs until the system shuts down. init manages the services (called demons under Linux, more on them in a future post). Also, if we check the process tree of a Linux machine we will find that the root of the tree is init.

There are multiple implementations for init, each of them provide different advantages among them are: SysVinit, launched, systemd, runit, upstart, busybox-init and OpenRC (those are examples only and not a full list). Thus, based on the implementation specific configuration files are read (such as /etc/inittab - SysVinit), different command/tools to manage demons (such as service - SysVinit and systemctl - systemd), and different scripts/profiles might be executed during the boot process (runlevels of SysVinit vs targets in systemd).

The creation of init is done by the kernel function “rest_init”¹. In the code we can see the call to “user_mode_thread” which spawns init, later in the function there is a call to “kernel_thread” which creates PID 2 (more information about it in the upcoming pages ;-).

Now we will go over a couple of fun facts about init. First, in case a parent process exits before all of its children process, init adopts those child processes. Second, only the signals which have been explicitly installed with a handler can be sent to init. Thus, sending “kill -9 1” won’t do anything in most distributions (try it and see nothing happens). Remember that different init implementations handle signals in different ways.

Because they are multiple init implementations (as we stated before) we can determine the one installed in the following manner. We can perform “ls -l /sbin/init”. If it is not a symlink it is probably SysVinit, else if it points to “/lib/systemd/systemd” than systemd is in use (and of course they are other symlinks to the other implementation - you can read about it in the documentation of each init implementation). As you can see in the attached screenshot Ubuntu 22.10 uses systemd.

¹ <https://elixir.bootlin.com/linux/v6.1.8/source/init/main.c#L683>

Kernel Threads

Before we will go over kthreadd I have decided to write a short post about kernel threads (due to the fact kthreadd is a kernel thread). We will go over some characteristics of kernel threads. First, kernel threads always execute in Kernel mode and never in User mode. Thus, kernel threads have basically all privileges and have no userspace address associated with them.

Second, both user mode process and kernel threads are represented by a task_struct inside the Linux kernel. As with all other user tasks, kernel threads are also part of the OS scheduling flow and can be executed on any CPU (there are cases in which there is a specific kernel thread for each CPU, we have seen it with swapper in the first post). Third, all kernel threads are descendants of kthreadd - Why is that? We will explain it in the next post focused on kthreadd.

Lastly, let's investigate kernel threads using /proc and see the difference in information retrieved from a regular user process (aka user task). There are multiple file entries in "/proc/pid" that contain information in case of a user mode process but are empty in case of a kernel thread, such as: "maps", "environ", "auxv", "cmdline" (I suggest reading "man proc" to get more info about them). Also, the fd and fdinfo directories are empty and the link "exe" does not point to any executable. In the attached screenshot we can see some of the difference between PID 1 [example of a regular user mode process] and PID 2 [example for a kernel thread]. BTW, the screenshot below was taken from an online/browser based Linux implementation called JSLinux - <https://bellard.org/jslinux>.

```
localhost:/# uname -a
Linux localhost 4.12.0-rc6-g48ec1f0-dirty #21 Fri Aug 4 21:02:28 CEST 2017 i586
Linux
localhost:/# cat /etc/issue
Welcome to Alpine Linux 3.12
Kernel \r on an \m (\l)

localhost:/# ls -l /proc/1/exe
lrwxrwxrwx  1 root  root          0 Aug 11 23:17 /proc/1/exe -> /bin/busybox
localhost:/# ls -l /proc/2/exe
ls: /proc/2/exe: cannot read link: No such file or directory
lrwxrwxrwx  1 root  root          0 Aug 11 23:16 /proc/2/exe
localhost:/# cat /proc/1/environ
HOME=/TERM=linuxTZ=UTC+07:00localhost:/#
localhost:/# cat /proc/2/environ
```

kthreadd (PID 2)

After explaining about PID 1, now we are going to talk about PID 2. Basically, kthreadd is the “kernel thread daemon”. Creation of a new kernel thread is done using kthreadd (We will go over the entire flow). Thus, the PPID of all kernel threads is 2 (checkout ps to verify this). As explained in the post about PID 1 (init) the creation of “kthreadd” is done by the kernel function “rest_init”². There is a call to the function “kernel_thread” (after the creation of init).

Basically, the kernel uses “kernel threads” (kthreads from now on) in order to run background operations. Thus, it is not surprising that multiple kernel subsystems are leveraging kthreads in order to execute async operations and/or periodic operations. In summary, the goal of kthreadd is to make available an interface in which the kernel can dynamically spawn new kthreads when needed.

Overall, kthreadd continuously runs (infinite loop³) and checks “kthread_create_list” for new kthreads to be created. In order to create a kthread the function “kthread_create”⁴ is used, which is a helper macro for “kthread_create_on_node”⁵. We can also call “kthread_run”⁶ could also be used, it is just a wrapper for “kthread_create”. The arguments passed to the creating function includes: the function to run in the thread, args to the function and a name.

While going over the source code we have seen that “kthread_create” calls “kthread_create_on_node”, which instantiates a “kthread_create_info” structure (based on the args of the function). After that, that structure is queued at the tail of “kthread_create_list” and “kthreadd” is awakened (and it waits until the kthread is created, this is done by “__kthread_create_on_node”⁷). What “kthreadd” does is to call “create_thread” based on the information queued. “create_thread” calls “kernel_thread”, which then calls “kernel_clone”. “kernel_clone” executes “copy_process”, which creates a new process as a copy of an old one - the caller needs to kick-off the created process (or thread in our case). By the way, the flow of creating a new task (recall every process/thread under Linux is called task and represented by “struct task_struct”) from user mode also gets to “copy_process”.

For the sake of simplicity, I have created a flow graph which showcases the flow of creating a kthread, not all the calls are there, only those I thought are important enough. Also, in both cases of macros/functions I used the verb “calls”. The diagram appears at the end of the post. Let me know if it is clear enough or do you think I should change something.

² <https://elixir.bootlin.com/linux/v6.1.8/source/init/main.c#L683>

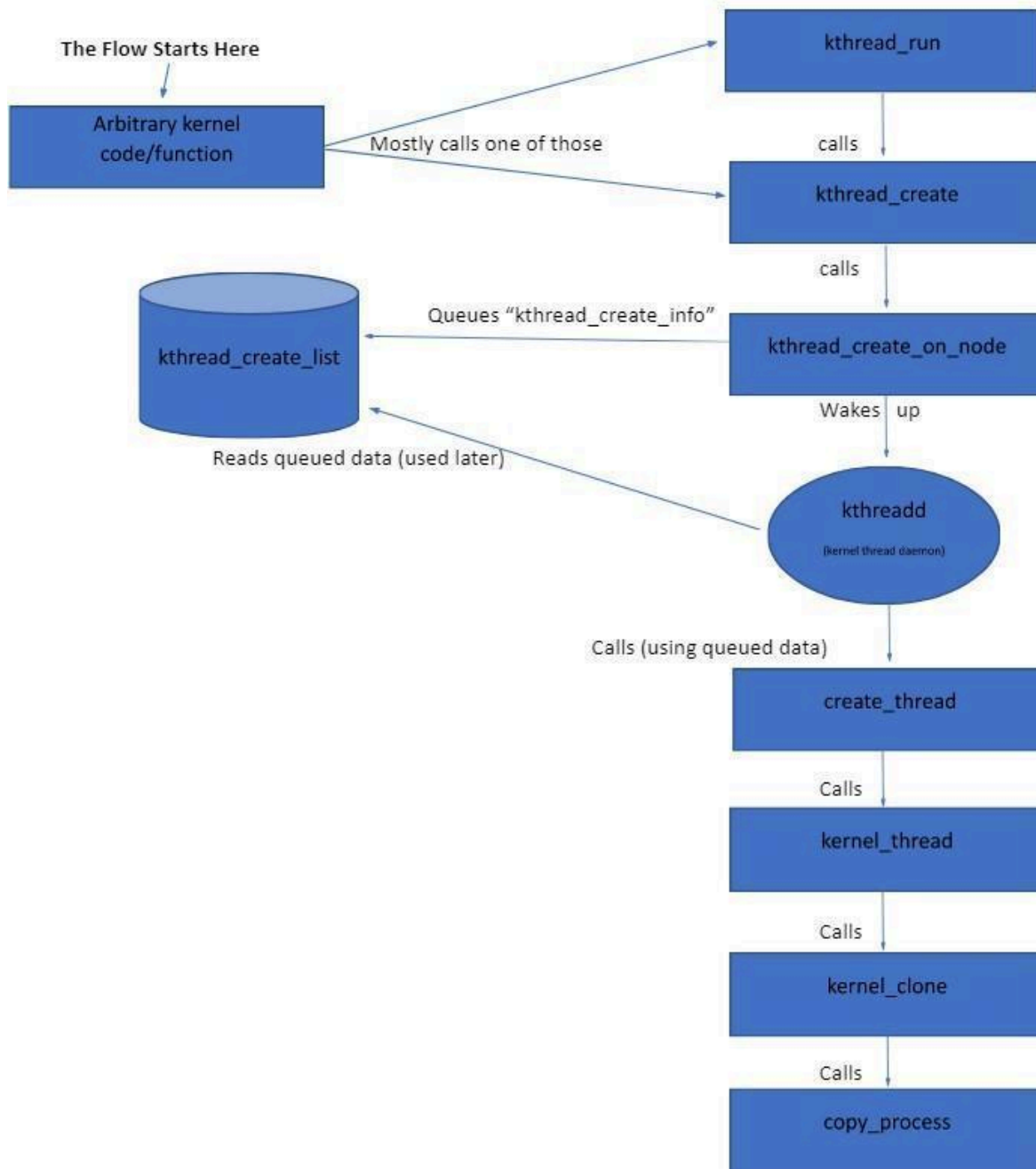
³ <https://elixir.bootlin.com/linux/v6.1.12/source/kernel/kthread.c#L731>

⁴ <https://elixir.bootlin.com/linux/v6.1.12/source/include/linux/kthread.h#L27>

⁵ <https://elixir.bootlin.com/linux/v6.1.12/source/kernel/kthread.c#L503>

⁶ <https://elixir.bootlin.com/linux/v6.1.12/source/include/linux/kthread.h#L51>

⁷ <https://elixir.bootlin.com/linux/v6.1.12/source/kernel/kthread.c#L414>



migration

One of the goals of an operating system is to handle and balance resources across the hardware of the compute entity. In order to do that, Linux has a kernel thread named “migration” which has an instance on every vCPU. By the way, the naming format is “migration/N” where N is the id of the vCPU.

By default threads are not constrained to a vCPU and can be migrated between them in the next call to “schedule()” (which calls the main scheduler function, which is “__scheduler()”⁸). It is done mainly in case the scheduler identifies an unbalanced across the runqueues (the queue in which processes which are in ready/runnable state are waiting to use the processor) of the vCPUs.

It is important to state that we can influence this flow by setting the affinity of a thread (for more read “man 2 sched_setaffinity”. We will talk about that in a future post). There could be performance, cache and other impacts for doing that (but that is also a topic for a different writeup).

I have created a small demo which shows the working of “migration”. For that I have created a VM running Ubuntu 22.04 with 3 vCPUs. In order to trace the usage of “move_queue_task” I have used bpftrace with the following command: **sudo bpftrace -e 'kfunc:move_queued_task { printf("%s moved %s to %d CPU\n",curtask->comm,args->p->comm,args->new_cpu); }'**. The output of the command is shown below. The one-liner prints: the name of the task calling “move_queued_task”, the name of the task which is moved and id the vCPU which the task is moved to.

```
Attaching 1 probe...
migration/2 moved sudo to 1 CPU
migration/1 moved dpkg to 2 CPU
migration/1 moved apt to 0 CPU
migration/1 moved update-motd-upd to 0 CPU
migration/1 moved (snap) to 0 CPU
migration/2 moved friendly-recv to 0 CPU
migration/2 moved lvm2-activation to 0 CPU
migration/0 moved (direxec) to 2 CPU
migration/2 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 2 CPU
migration/2 moved (direxec) to 1 CPU
migration/2 moved (direxec) to 0 CPU
migration/0 moved udiskd to 1 CPU
migration/2 moved bash to 1 CPU
migration/2 moved bash to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/2 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
```

⁸ <https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L6544>

In summary, what the kernel thread “migration” does is to move threads from highly loaded vCPUs to others which are less crowded (by inserting them to a different run-queue). A function which is used by “migration” in order to move a task to a new run-queue is “move_queued_task” (<https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L2325>).

charger_manager

The “charger_manager” kernel thread is created by a freezable workqueue⁹. Freezable workqueues are basically frozen when the system is moved to a suspend state¹⁰. Based on the kernel source code “charger_manager” is responsible for monitoring the health (like temperature monitoring) of the battery and controlling the charger while the system is suspended to memory¹¹. The “Charger Manager” kernel module is written by MyungJoo Ham¹².

Moreover, the kernel documentation states that the “Charger Manager” also helps in giving an aggregated view to user-space in case there are multiple chargers for a battery. In case they are multiple batteries with different chargers on a system, that system would need multiple instances of “Charger Manager”¹³.

On my Ubuntu VM (22.04.1 LTS) this kernel module is not compiled as a separate “*.ko” file. It is compiled into the kernel itself (builtin), as you can see in the output of “modinfo” in the screenshot below.

```
Troller $ modinfo charger_manager
name:          charger_manager
filename:      (builtin)
license:      GPL
file:         drivers/power/supply/charger-manager
description:  Charger Manager
author:       MyungJoo Ham <myungjoo.ham@samsung.com>
Troller $ █
```

⁹ <https://elixir.bootlin.com/linux/latest/source/drivers/power/supply/charger-manager.c#L1749>

¹⁰ <https://lwn.net/Articles/403891/>

¹¹ <https://elixir.bootlin.com/linux/latest/source/drivers/power/supply/charger-manager.c>

¹² <https://elixir.bootlin.com/linux/latest/source/drivers/power/supply/charger-manager.c#L1768>

¹³ <https://www.kernel.org/doc/html/v5.3/power/charger-manager.html>

idle_inject

On our plate this time we are going to talk about the kernel thread “idle_inject”, which was merged to the kernel in about 2009. The goal of “idle_inject” is forcing idle time on a CPU in order to avoid overheating.

If we think about it, “idle_inject” adds latency, thus it should be considered only if CPUFreq (CPU Frequency scaling) is not supported. Due to the fact the majority of modern CPUs are capable of running a different clock frequency and voltage configuration we can use CPUFreq in order to avoid overheating.

Overall, there is one “idle_inject” kernel thread per processor (with the name pattern “idle_inject/N”, where N is the id of the processor) - as shown in the screenshot below. Also, all of them are created at init time.

The “idle_inject” kernel threads will call “idle_inject_fn()”->”play_idle_precise()” to inject a specified amount of idle time. After all of the kernel threads are woken up, the OS sets a timer for the next cycle. When the timer interrupt handler wakes the threads for all processors based on a defined “cpu-mask” (affected by idle injection). By the way, when I set a kprobe on “idle_inject_fn()” for 3 hours on my VM it was never called ;-)

```
Troller# ps -eo user,comm,pid,ppid | grep idle_inject
root      idle_inject/0      16      2
root      idle_inject/1      19      2
root      idle_inject/2      25      2
Troller# █
```

kworker (Kernel Thread Worker)

A kworker is a kernel thread that performs processing as part of the kernel, especially in the case of interrupts, timers, I/O, etc. It is based on workqueues which are async execution mechanisms, that execute in “process context” (I will post on workqueus in more details separately, for now it is all that you need to know).

Overall, there are a couple of kworkers running on a Linux machine. The naming pattern of kworkers includes: the number of the core on which it is executed, the id of the thread and can contain also string that hints what the kworker does (check the output of ‘ps -ef | grep kworker’).

```
6          2  0 07:15 ?          00:00:00 [kworker/0:0H-events_highpri]
82         2  0 07:15 ?          00:00:02 [kworker/0:1H-kblockd]
113        2  0 07:15 ?          00:00:00 [kworker/u3:0]
46277      2  0 11:11 ?          00:00:00 [kworker/u2:1-events_unbound]
46547      2  0 11:20 ?          00:00:01 [kworker/0:1-events]
46624      2  0 11:23 ?          00:00:00 [kworker/u2:2-kcryptd/253:0]
46867      2  0 11:28 ?          00:00:00 [kworker/0:0-inet_frag_wq]
47091      2  0 11:33 ?          00:00:00 [kworker/u2:0-events_unbound]
47299      2  0 11:36 ?          00:00:00 [kworker/0:2-events]
```

The big question is - “How do we know what each kworker is doing?”. It’s a great question, the way in which we are going to answer it is by using ftrace (function tracing inside the kernel - I suggest reading more about that - <https://www.kernel.org/doc/Documentation/trace/ftrace.txt>). The command we are going to use are:

```
echo workqueue:workqueue_queue_work > /sys/kernel/debug/tracing/set_event
cat /sys/kernel/debug/tracing/trace_pipe > /tmp/trace.log
```

The first one enables the tracing regarding workqueus. The second reads the tracing data and saves it to a file. We can also run “cat /sys/kernel/debug/tracing/trace_pipe | grep kworker” and change the grep filter to a specific kworker process. In the trace we will see the function name that each kworker thread is going to execute.

```
kworker/u2:2-46624 [000] d... 17855.481276: workqueue_queue_work: work struct=00000000da1e6721 function=flush_to_idisc
workqueue=events_unbound req_cpu=8192 cpu=4294967295
kworker/u2:1-48183 [000] d... 17855.525798: workqueue_queue_work: work struct=00000000be96cc25 function=ata_sff_pio_ta
< workqueue=ata_sff req_cpu=8192 cpu=0
kworker/u2:1-48183 [000] d... 17856.038232: workqueue_queue_work: work struct=000000001e1ee94f function=kcryptd_crypt
dm_crypt] workqueue=kcryptd/253:0 req_cpu=8192 cpu=4294967295
kworker/u2:1-48183 [000] d... 17857.542509: workqueue_queue_work: work struct=00000000be96cc25 function=ata_sff_pio_ta
< workqueue=ata_sff req_cpu=8192 cpu=0
kworker/u2:1-48183 [000] d... 17859.558293: workqueue_queue_work: work struct=00000000be96cc25 function=ata_sff_pio_ta
< workqueue=ata_sff req_cpu=8192 cpu=0
kworker/u2:1-48183 [000] d... 17860.134032: workqueue_queue_work: work struct=000000001e1ee94f function=kcryptd_crypt
dm_crypt] workqueue=kcryptd/253:0 req_cpu=8192 cpu=4294967295
kworker/u2:1-48183 [000] d... 17860.134074: workqueue_queue_work: work struct=00000000e0b6b12c function=kcryptd_crypt
dm_crypt] workqueue=kcryptd/253:0 req_cpu=8192 cpu=4294967295
```

kdevtmpfs

“kdevtmpfs” is a kernel thread which was created using the “kthread_run” function¹⁴. “kdevtmpfs” creates a devtmpfs which is a tmpfs-based filesystem (/dev). The filesystem is created during bootup of the system, before any driver code is registered. In case a driver-core requests a device node it will result in a node added to this filesystem¹⁵.

We can see the specific line of code that is used in order to create the mounting point “/dev”¹⁶. The mountpoint is created using the function “init_mount”¹⁷. A nice fact is that it is part of “init_*” functions which are routines that mimic syscalls but don’t use file descriptors or the user address space. They are commonly used by early init code¹⁸.

Thus, we can say the “kdevtmpfs” is responsible for managing the “Linux Device Tree”. Also, by default the name created for nodes under the filesystem is based on the device name (and owned by root) - as shown in the screenshot below (taken from copy.sh based Linux). By the way, not all devices have a node in “/dev” think about network devices ;-)

```
root@localhost:/dev# mount | grep "/dev"| head -1
dev on /dev type devtmpfs (rw,nosuid,relatime,size=10240k,nr_inodes=58635,mode=755)
root@localhost:/dev# ls -lah | head -20
total 1.0K
drwxr-xr-x 11 root root    3.4K Nov  7 02:51 .
drwxrwxrwx 17 root root      0 Nov  7 02:50 ..
crw-r--r--  1 root root   10, 235 Nov  7 02:50 autofs
drwxr-xr-x  2 root root   2.5K Nov  7 02:50 char
crw-----  1 root root     5,  1 Nov  7 02:51 console
lrwxrwxrwx  1 root root     11 Nov  7 02:50 core -> /proc/kcore
drwxr-xr-x  3 root root     60 Nov  7 02:50 cpu
crw-----  1 root root  10, 125 Nov  7 02:50 cpu_dma_latency
drwxr-xr-x  2 root root     60 Nov  7 02:50 dma_heap
drwxr-xr-x  2 root root     60 Nov  7 02:51 dri
crw-----  1 root root    29,  0 Nov  7 02:51 fb0
lrwxrwxrwx  1 root root     13 Nov  7 02:50 fd -> /proc/self/fd
crw-rw-rw-  1 root root     1,  7 Nov  7 02:50 full
drwxr-xr-x  2 root root     80 Nov  7 02:50 input
crw-r--r--  1 root root     1, 11 Nov  7 02:50 kmsg
crw-r-----  1 root root     1,  1 Nov  7 02:50 mem
drwxrwxrwt  2 root root     40 Nov  7 02:50 mqueue
crw-rw-rw-  1 root root     1,  3 Nov  7 02:50 null
crw-----  1 root root  10, 144 Nov  7 02:50 nram
```

¹⁴ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L474>

¹⁵ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L3>

¹⁶ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L377>

¹⁷ <https://elixir.bootlin.com/linux/v6.2-rc1/source/fs/init.c#L16>

¹⁸ <https://elixir.bootlin.com/linux/v6.2-rc1/source/fs/init.c#L3>

cpuhp (CPU Hotplug)

This kernel thread is part of the CPU hotplug support. It enables physically removing/adding CPUs on a specific system. There is one kernel thread per vCPU, and the pattern of the thread's name is "cpuhp/N" (where N is the id of the vCPU) - as can be seen in the screenshot below. Also, today the CPU hotplug can be used to resume/suspend support for SMP (Symmetric Multiprocessing).

If we want our kernel to support CPU hotplug the CONFIG_HOTPLUG_CPU should be enabled (it's supported on a couple of architectures such as: MIPS, ARM, x86 and PowerPC). The kernel holds the current state for each CPU by leveraging "struct cpuhp_cpu_state"¹⁹.

We can configure the CPU hotplug mechanism using sysfs (/sys/devices/system/cpu). For example we can shut down and bring up a CPU by writing "0" and "1" respectively to the "online" file in the directory representing the CPU (for which we want to change the status) - checkout the screenshot below (the Linux VM I am testing on has 3 vCPUs).

In order to bring the CPU down the function "cpu_device_down"²⁰ is called. In order to bring up a CPU function "cpu_device_up"²¹ is called.

```
Troller # pwd
/sys/devices/system/cpu
Troller # ls
cpu0  cpufreq  isolated  offline  power  uevent
cpu1  cpuidle  kernel_max  online  present  vulnerabilities
cpu2  hotplug  modalias  possible  smt
Troller # echo 0 > ./cpu2/online
Troller # dmesg | tail -2
[147586.057954] kvm-clock: cpu 1, msr b7001041, secondary cpu clock
[148846.125346] smpboot: CPU 2 is now offline
Troller # echo 1 > ./cpu2/online
Troller # dmesg | tail -2
[148846.125346] smpboot: CPU 2 is now offline
[148874.835266] smpboot: Booting Node 0 Processor 2 APIC 0x2
```

¹⁹ <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L65>

²⁰ <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L1225>

²¹ <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L1439>

kswapd

The kernel thread “kswapd” is the background page-out daemon of Linux (swaps processes to disk). You can see the creation of the kernel thread in the source of the kernel - <https://elixir.bootlin.com/linux/latest/source/mm/vmscan.c#L4642>. In the code we can see that a dedicated instance of “kswapd” is created for each NUMA zone (on my Ubuntu 22.10 VM I have only “kswapd0” - as shown in the screenshot below).

Overall, the goal of the “kswapd” is to reclaim pages when memory is running low. In the old days, the “kswapd” was woken every 10 seconds but today it is only wakened by the page allocator, by calling “wakeup_kswapd”²⁴. The code of the page allocator is located at “mm/page_alloc.c”²⁵.

Basically, “kswapd” trickles out pages so the system has some free memory even if no other activity frees up anything (like by shrinking cache). Think about cases in which operations work in asynchronous contexts that cannot page things out.

The major function which is called by “kswapd” is “balance_pgdat()”²⁶. In order to see that process happening we can use the following bpftrace one-liner: “**sudo bpftrace -e 'kfunc:balance_pgdat { printf("%s:%d\n",curtask->comm,curtask->pid); }**” - You can see “kswapd0” calling it in the screenshot below. The flow of “kswapd” is based on limits, when to start shirking and “until when” to shrink (low and high limits).

```
Troller # sudo bpftrace -e 'kfunc:balance_pgdat { printf("%s:%d\n",curtask->comm,curtask->pid); }'  
Attaching 1 probe...  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97  
kswapd0:97
```

²⁴ <https://elixir.bootlin.com/linux/latest/source/mm/vmscan.c#L4555>
²⁵ https://elixir.bootlin.com/linux/latest/source/mm/page_alloc.c
²⁶ <https://elixir.bootlin.com/linux/latest/source/mm/vmscan.c#L4146>

kcompactd

When a Linux system is up and running, memory pages of different processes/tasks are scattered and thus are not physically-contiguous (even if they are contiguous in their virtual address). We can move to bigger pages size (like from 4K to 4M) but it still has its limitations like: waste of space in case of regions with small sizes and the need for multiple pages in case of large regions that can still be fragmented. Due to that, the need for memory compaction was born²⁷.

“kcompactd” is performing in the background the memory compaction flow. The goal of memory compaction is to reduce external fragmentation. This procedure is heavily dependent on page migration²⁸ to do all the heavy lifting²⁹. In order for “kcompactd” to work we should compile the kernel with “CONFIG_COMPACTON” enabled. Also, when a Linux system identifies that it is tight low in available memory the “kcompactd” won’t perform memory compaction memory³⁰.

Overall, the “kcompactd” kernel thread is created in “kcompactd_run” function³¹ which is called by “kcompactd_init”³². The function “kcompactd_init” is started by “subsys_initcall”³³, which is responsible for initializing a subsystem.

The kernel thread starts the function “static int kcompactd(void *p)”³⁴. An instance of the kernel thread is created for each node (like vCPU) on the system³⁵. The pattern of the kernel thread name is “kcompactd[IndexOfNode]” for example “kcompactd0” as we can see in the screenshot below.

“kcompactd” can be called in one of two ways: woken up or by using a timeout. It can be woken up by kswapd³⁶. Also, we can configure it using modification of the filesystem (“/proc/sys/vm/compact_memroy” for example). By the way, in the memory compaction flow of the function “compact_zone”³⁷ is executed in the context of “kcompactd”. In order to demonstrate that we can use the following one-liner using bpftrace: **sudo bpftrace -e 'kfunc:compact_zone { printf("%s:%d\n",curtask->comm,curtask->pid); }'** - The output can be seen in the screenshot below.

²⁷ <https://lwn.net/Articles/368869/>

²⁸ <https://lwn.net/Articles/157066/>

²⁹ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L5>

³⁰ <https://www.linux-magazine.com/Issues/2015/179/Kernel-News>

³¹ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2996>

³² <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3048>

³³ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3065>

³⁴ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2921>

³⁵ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3061>

³⁶ <https://www.slideshare.net/AdrianHuang/memory-compaction-in-linux-kernelpdf>

³⁷ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2289>

```
Troller # ps -ef | grep -v grep | grep kcompactd
root      37      2  0 00:15 ?                00:00:09 [kcompactd0]
Troller # ls -l /proc/sys/vm/compact_memory
--w----- 1 root root 0 Jan 14 11:54 /proc/sys/vm/compact_memory
Troller # sudo bpftrace -e 'kfunc:compact_zone { printf("%s:%d\n",curtask->comm,curtask->pid); }'
Attaching 1 probe...
kcompactd0:37
kcompactd0:37
kcompactd0:37
```

md (Multiple Device Driver)

“md” is a kernel thread which is based on a workqueue³⁸. It is responsible for managing the Linux md (multiple device) driver which is also known as the “Linux software RAID”. RAID devices are virtual devices (created from two or more real block devices). This allows multiple devices (typically disk drives or partitions thereof) to be combined into a single device to hold (for example) a single filesystem³⁹.

By using the “md” driver we can create from one/more physical devices (like disk drivers) a virtual device(s). By the use of an array of devices we can achieve redundancy, which is also known as RAID (Redundant Array of Independent Disks). For more information I suggest reading <https://man7.org/linux/man-pages/man4/md.4.html>.

Overall, “md” supports different RAID types: RAID 1 (mirroring), RAID 4, RAID 5, RAID 6 and RAID 10. For more information about RAID types I suggest reading the following link <https://www.prepressure.com/library/technology/raid>. Besides that, “md” also supports pseudo RAID technologies like: RAID 0, LINAR, MULTIPATH and FAULTY⁴⁰.

The code of “md” is included as a driver/kernel module in the source code of Linux. Thus, it can be compiled directly into the kernel or as a separate “*.ko” file. In my VM (Ubuntu 22.04) it is compiled directly into the kernel image as shown in the screenshot below.

```
Troller $ ps -ef | grep -v grep | grep "\[md]"
root          78          2  0 Dec21 ?           00:00:00 [md]
Troller $ lsmod | grep " md "
Troller $ modinfo md
name:          md_mod
filename:      (builtin)
alias:         block-major-9-*
alias:         md
description:   MD RAID framework
license:       GPL
file:          drivers/md/md-mod
parm:          start_dirty_degraded:int
parm:          create_on_open:bool
Troller $
```

³⁸ <https://elixir.bootlin.com/linux/v6.1/source/drivers/md/md.c#L9615>

³⁹ <https://linux.die.net/man/8/mdadm>

⁴⁰ https://doxfer.webmin.com/Webmin/Linux_RAID

The block devices that can be used in order to access the software RAID on Linux are in the pattern “/dev/mdN” (where N is a number [0–255])⁴¹. It can also be configured to allow access using “/dev/md/N” or “/dev/md/name”. If we want information about the current state of “md” we can query the file “/proc/mdstat” — for more information you can read <https://raid.wiki.kernel.org/index.php/Mdstat>. There is also the command line utility “mdadm” that can help with managing those devices⁴².

Lastly, the init function is declared using “subsys_initcall” (and not the “module_init”) which ensures that it will run before the device drivers that needs it (if they are using “module_init”) — <https://elixir.bootlin.com/linux/v6.1/source/drivers/md/md.c#L9947>. More information about initcalls will be included on a future writeup.

⁴¹ <https://www.oreilly.com/library/view/managing-raid-on/9780596802035/ch01s03.html>

⁴² <https://linux.die.net/man/8/mdadm>

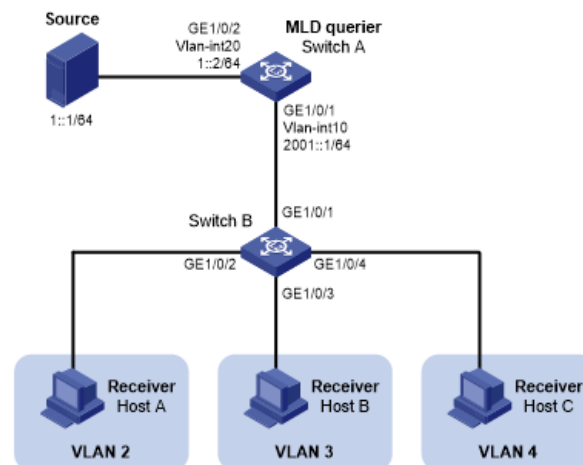
mld (Multicast Listener Discovery)

“mld” is a kernel thread which was created using a workqueue⁴³. It is the Linux implementation for the multicast listener (MLD) protocol. This protocol is used by IPv6 based routers in order to discover multicast listeners on the local network and identify which multicast addresses are of interest to those listeners. MLD is supported on different operating systems such as Windows⁴⁴ and Linux⁴⁵.

We can think about it like IGMP⁴⁶ which is used on IPv4 based networks (MLDv1 is derived from IGMPv2 and MLDv2 is similar to IGMPv3). One important difference is that MLD uses ICMPv6 message types, rather than IGMP message types⁴⁷.

Overall, MLD has three major message types: “Multicast Listener Query”, “Multicast Listener Report” and “Multicast Done”. For more information about them I suggest reading the following link⁴⁸. Also, a more detailed explanation about the different MLD operations can be found in <https://ipisco.com/lesson/mld-operations/>.

What “mld” does is to send MLD report messages⁴⁹ which are sent by an MLD host (see the diagram below⁵⁰) and processes messages⁵¹. From the source code we can see that there are definitions for structs representing both MLDv1 and MLDv2 headers.



⁴³ <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L3185>

⁴⁴ <https://learn.microsoft.com/en-us/windows/win32/winsock/igmp-and-windows-sockets>

⁴⁵ <https://lwn.net/Articles/29489/>

⁴⁶ <https://www.cloudflare.com/learning/network-layer/what-is-igmp/>

⁴⁷ <https://www.ibm.com/docs/en/zos/2.2.0?topic=protocol-multicast-listener-discovery>

⁴⁸ <https://community.cisco.com/t5/networking-knowledge-base/multicast-listener-discovery-mld/ta-p/3112082>

⁴⁹ <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L3185>

⁵⁰ https://techhub.hpe.com/eginfolib/networking/docs/switches/5130ei/5200-3944_ip-multi_cg/content/images/image33.png

⁵¹ <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L1359>

ksmd (Kernel Same Page Merging)

The kernel thread “ksm” is also known as “Kernel Same Page Merging” (and “ksmd” is ksm demon). It is used by the KVM hypervisor to share identical memory pages (supported since kernel 2.6.32) Those shared pages could be common libraries or even user data which is identical. By doing so KVM (Kernel-based Virtual Machine) can avoid memory duplication and enable more VMs to run on a single node.

In order for “ksmd” to save memory due to de-duplication we should compile the kernel with “CONFIG_KSM=y”. It is important to understand that the sharing of identical pages is done even if they are not shared by fork(). If you want to go over “ksmd” source code you can use the following link - <https://elixir.bootlin.com/linux/latest/source/mm/ksm.c>.

The way “ksmd” works is as follows. Scanning main memory for frames (“physical pages”) holding identical data and collectes the virtual memory address that they are mapped. “ksmd” leaves one of those frames and remaps each duplicate one to point to the same frame. Lastly, “ksmd” frees the other frames. All of the merge pages are marked as COW (Copy-on-Write) for cases in which one of the processes using them will want to write to the page. There is a concern that even if the memory usage is reduced the CPU usage is increased.

The kernel thread “ksmd” is created using the function kthread_run⁵². We can see from the code that the function which is the entry point of the thread is “ksm_scan_therad()” which is calling “ksm_do_scan()” which is the ksm’s scanner main worker function (it gets as input the number of pages to scan before returning). “ksmd” only merges anonymous private pages and not pagecache. Historically, the merged pages were pinned into kernel memory. Today they can be swapped like any other pages.

“ksmd” can be controlled by a sysfs interface (“/sys/kernel/mm/ksm”) - as can be seen in the screenshot below. One of the files exported by sysfs is “run” that can react to one of the following values 0/1/2. “0” means stop “ksmd” from running but keep the merged pages. “1” means run “ksmd”. “2” means stop “ksmd” from running and unmerge all currently merge pages (however leave the mergeable areas registered for next time).

```
Troller # pwd
/sys/kernel/mm/ksm
Troller # ls *
full_scans          pages_shared      pages_unshared    sleep_millisecs   stable_node_dups
max_page_sharing   pages_sharing     pages_volatile    stable_node_chains use_zero_pages
merge_across_nodes pages_to_scan     run               stable_node_chains_prune_millisecs
```

⁵² <https://elixir.bootlin.com/linux/v6.0/source/mm/ksm.c#L3188>

ttm_swap

The kernel thread “ttm_swap” is responsible for swapping GPU’s (Graphical Processing Unit) memory. Overall, TTM (Translation-Table Maps) is a memory manager that is used to accelerate devices with dedicated memory. Basically, all the resources are grouped together by objects of buffers in different sizes. TTM then handles the lifetime, the movements and the CPU mapping of those objects⁵³.

Based on the kernel documentation, each DRM (Direct Rendering Manager) driver needs a memory manager. There are two memory managers supported by DRM: TTM and GEM (Graphics Execution Manager). I am not going to talk about GEM, if you want you can start reading about in the following link - <https://docs.kernel.org/gpu/drm-internals.html>.

Moreover, “ttm_swap” is a single threaded workqueue as seen in the Linux source code⁵⁴.

Also, the man pages describe TTM as a generic memory-manager provided by the kernel, which does not provide a user-space interface (API). In case we want to use it you should checkout the interface of each driver⁵⁵.

TTM is at the end a kernel module, you can find the source code and the Makefile in the kernel source tree⁵⁶. Based on the module source code it is written by Thomas Hellstrom and Jerome Glisse⁵⁷. Also, it is described as “TTM memory manager subsystem (for DRM device)”⁵⁸. As you can see it is part of the “drivers/gpu/drm” subdirectory, which holds the code and Makefile of the drm device driver, which provides support for DRI (Direct Rendering Infrastructure) in XFree86 4.1.0+. Lastly, on my VM (Ubuntu 22.04.01) it is compiled as a separate “*.ko” file (/lib/modules/[KernelVersion]/kernel/drivers/gpu/drm/ttm.ko) - as shown in the screenshot below.

```
troller # modinfo ttm | head -15
filename:      /lib/modules/5.15.0-52-generic/kernel/drivers/gpu/drm/ttm/ttm.ko
license:      GPL and additional rights
description:   TTM memory manager subsystem (for DRM device)
author:       Thomas Hellstrom, Jerome Glisse
srcversion:   52AE33CCBE42B11150B88C3
depends:       drm
retpoline:    Y
intree:       Y
name:         ttm
vermagic:     5.15.0-52-generic SMP mod_unload modversions
sig_id:       PKCS#7
signer:       Build time autogenerated kernel key
sig_key:      49:B2:3F:66:E1:3B:8B:67:11:CE:17:63:41:27:D0:B1:28:DF:09:8C
sig_hashalgo: sha512
```

⁵³ <https://docs.kernel.org/gpu/drm-mm.html>

⁵⁴ https://elixir.bootlin.com/linux/v5.12.19/source/drivers/gpu/drm/ttm/ttm_memory.c#L424

⁵⁵ <https://www.systutorials.com/docs/linux/man/7-drm-ttm/>

⁵⁶ <https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm>

⁵⁷ https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm_module.c#L89

⁵⁸ https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm_module.c#L89

watchdogd (Watchdog Daemon)

This kernel thread “watchdogd” is used in order to let the kernel know that a serious problem has occurred so the kernel can restart the system. It is sometimes called COP (Computer Operating Properly). The way it is implemented is by opening “/dev/watchdog”, then writing at least once a minute. Every time there is a write the restart of the system is delayed.

In case of inactivity for a minute the watchdog should restart the system. Due to the fact we are not talking about a hardware watchdog the compilation of the operation depends on the state of the machine. You should know that the watchdog implementation could be software only (there are cases in which it won't restart the machine due to failure) or using a driver/module in case of hardware support⁵⁹.

If we are talking about hardware support then the watchdog module is specific for a chip or a device hardware. It is most relevant to systems that need the ability to restart themselves without any human intervention (as opposed to a PC we can reboot easily) - think about an unmanned aircraft. We need to be careful because a problem in the watchdog configuration can lead to unpredictable reboot, reboot loops and even file corruption due to hard restart⁶⁰.

The relationship between the hardware and software is as follows: the hardware is responsible to set up the timer and the software is responsible to reset the timer. When the timer gets to a specific value (configured ahead) and it is not elapsed by the software the hardware will restart the system. For an example of using hardware for this functionality you can read the following link <https://developer.toradex.com/linux-bsp/how-to/linux-features/watchdog-linux/>.

The software part is being conducted by the “watchdogd” (the software watchdog daemon) which opens “/dev/watchdog” and writes to it in order to postpone the restart of the system by the hardware - for more information you can read <https://linux.die.net/man/8/watchdog>. Examples for different watchdog drives/modules for specific chips can be found in the source tree of linux here <https://elixir.bootlin.com/linux/v6.0.11/source/drivers/watchdog>. Some examples are apple_wdt (Apple's SOC), ath79_wdt (Atheros AR71XX/AR724X/AR913X) and w83977f_wdt (Winbond W83977F I/O Chip).

We can stop the watchdog without restarting the system by closing “/dev/watchdog”. It is not possible if the kernel was compiled with “CONFIG_WATCHDOG_NOWAYOUT” enabled.

⁵⁹ <https://github.com/torvalds/linux/blob/master/Documentation/watchdog/watchdog-api.rst>

⁶⁰ <https://linuxhint.com/linux-kernel-watchdog-explained/>

Overall, in order for the watchdog to operate the kernel needs to be compiled with CONFIG_WATCHDOG=y and “/dev/watchdog” character device should be created (with major number of 10 and minor number of 130 - checkout “man mknod” if you want to create it).

Lastly, if you want to see the status of the watchdog you can use the command “wdctl”⁶¹ - As can be seen in the screenshot below⁶². For more information about the concept I suggest reading https://en.wikipedia.org/wiki/Watchdog_timer.

```
[root@ako-kaede-mirai]-(12:25am--09/06) r- -""
[r(kousekip) r` wdctl
Device: /dev/watchdog0
Identity: SP5100 TCO timer [version 0]
Timeout: 60 seconds
Pre-timeout: 0 seconds
FLAG DESCRIPTION STATUS BOOT-STATUS
KEEPALIVEPING Keep alive ping reply 1 0
MAGICCLOSE Supports magic close char 0 0
SETTIMEOUT Set timeout (in seconds) 0 0
```

⁶¹ <https://man7.org/linux/man-pages/man8/wdctl.8.html>

⁶² https://en.wikipedia.org/wiki/Watchdog_timer#/media/File:Wdctl_screenshot.png

zswap-shrink

Based on the kernel source code zswap is a backend for frontswap. Frontswap provides a “transcendent memory” interface for swap pages. In some cases we can get increased performance by saving swapped pages in RAM (or a RAM-like device) and not on disk as swap partition\swapfile⁶³. The frontends are usually implemented in the kernel while the backend is implemented as a kernel module (as we will show soon). Zswap takes pages that are in the process of being swapped out and attempts to compress and store them in a RAM-based memory pool⁶⁴.

We can say that zswap trades CPU cycles for potentially reduced swap I/O. A significant performance improvement can happen in case the reads from the swap device are much slower than the reads from the compressed cache⁶⁵. The “zswap_frontswap_store” is the function that attempts to compress and store a single page⁶⁶.

The kernel thread “zswap-shrink” is created based on a workqueue⁶⁷. On my VM (Ubuntu 22.04.1) zswap is compiled part of the kernel itself and not as a separate “*.ko” (kernel module). You can see in the screenshot below that it does not appear in the output of “lsmod” and is marked as builtin (look at the filename field) in the output of “modinfo”.

```
Troller # ps -ef | grep zswap-shrink #show the zswap-shrink kernel thread
root      128          2  0 Oct21 ?        00:00:00 [zswap-shrink]
root     169924    164567  0 20:39 pts/6    00:00:00 grep --color=auto zswap-shrink
Troller # lsmod | grep zswap #check if zswap is loaded outside the kernel
Troller # modinfo zswap #show zswap builtin
name:          zswap
filename:      (builtin)
description:   Compressed cache for swap pages
author:       Seth Jennings <sjennings@variantweb.net>
license:      GPL
file:         mm/zswap
parm:         max_pool_percent:uint
parm:         accept_threshold_percent:uint
parm:         same_filled_pages_enabled:bool
Troller # dmesg | grep zswap
[    1.071279] zswap: loaded using pool lzo/zbud
Troller # █
```

For more information like the compression used by zswap (the default one is lzo) and other parameters that can be configured for zswap I suggest reading the following link <https://wiki.archlinux.org/title/zswap>. You can also read the parameter ones “/sys/module/zswap/parameters”.

⁶³ <https://www.kernel.org/doc/html/v4.18/vm/frontswap.html>

⁶⁴ <https://elixir.bootlin.com/linux/latest/source/mm/zswap.c>

⁶⁵ <https://www.kernel.org/doc/html/v4.18/vm/zswap.html>

⁶⁶ <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1097>

⁶⁷ <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1511>

khugepaged (Kernel Huge Pages Daemon)

The kernel thread “kugepaged” is created using the “kthread_run()” function⁶⁸. It is responsible for the “Transparent Hugepage Support” (aka THP). “kugepaged” scans memory and collapses sequences of basic pages into huge pages⁶⁹.

We can manage and configure TPH using sysfs⁷⁰ or by using the syscalls “madvise”⁷¹ and “prctl”⁷². The scan of memory is done by calling “khugepaged_do_scan()”⁷³ which in turn calls “khugepaged_scan_mm_slot()”⁷⁴. In order to demonstrate that I have used the following bpftrace oneliner “**sudo bpftrace -e 'kfunc:khugepaged_scan_mm_slot{ printf("%s:%d\n",curtask->comm,curtask->pid); }**”. The output is shown in the screenshot below.

Lastly, we can also monitor the modifications made by “khugepaged” by checking the information on “/proc”. For example we can check the “AnonHugePages”/”ShmemPmdMapped”/”ShmemHugePages” in “/proc/meminfo”, which is global for the entire system. If we want information regarding a specific process/task we can use “/proc/[PID]/smaps” and count “AnonHugePages”/”FileHugeMapped” for each mapping (<https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html>).

```
Troller $ sudo bpftrace -e 'kfunc:khugepaged_scan_mm_slot{ printf("%s:%d\n",curtask->comm,curtask->pid); }'
Attaching 1 probe...
khugepaged:39
khugepaged:39
khugepaged:39
khugepaged:39
khugepaged:39
khugepaged:39
```

⁶⁸ <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2551>

⁶⁹ <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html>

⁷⁰ <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html#thp-sysfs>

⁷¹ <https://man7.org/linux/man-pages/man2/madvise.2.html>

⁷² <https://man7.org/linux/man-pages/man2/prctl.2.html>

⁷³ <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2404>

⁷⁴ <https://elixir.bootlin.com/linux/v6.1.12/source/mm/khugepaged.c#L2250>

krfcommd (Kernel Radio Frequency Communication Daemon)

“krfcommd” is a kernel which is started by executing “kthread_run()” function⁷⁵. The kernel thread executes the “rfcomm_run()” function⁷⁶. Thus, we can say that “krfcommd” is responsible for RFCOMM connections⁷⁷.

RFCOMM (Radio Frequency Communication) is a set of transport protocols on top of L2CAP which provides emulated RS-232 serial ports. It provides a simple reliable data stream (like TCP). It is used directly by many telephony related profiles as a carrier for AT commands, as well as being a transport layer for OBEX over Bluetooth⁷⁸.

Moreover, there is also an “rfcomm” cli tool in Linux. It is used to inspect and maintain RFCOMM configuration⁷⁹. For more information about RFCOMM I suggest reading <https://www.btframework.com/rfcomm.htm>. You can also go over the protocol specification⁸⁰.

Also, RFCOMM protocol supports up to 60 simultaneous connections between two Bluetooth devices. The number of connections that can be used simultaneously is implementation-specific. For the purposes of RFCOMM, a complete communication path involves two applications running on different devices (the communication endpoints) with a communication segment between them⁸¹.

Lastly, RFCOMM is implemented as a kernel module. Thus, it can be compiled directly to the kernel or separate kernel module - in the screenshot below we can see it compiled as a separate file.

```
root@localhost:~# modinfo rfcomm
filename:       /lib/modules/5.19.7-arch1-1.0/kernel/net/bluetooth/rfcomm/rfcomm.ko.zst
alias:         bt-proto-3
license:       GPL
version:       1.11
description:   Bluetooth RFCOMM ver 1.11
author:        Marcel Holtmann <marcel@holtmann.org>
srcversion:    27B7EECAEC282A1A24A7701
depends:        bluetooth
retpoline:    y
intree:        y
name:          rfcomm
vermagic:      5.19.7-arch1-1.0 SMP preempt mod_unload 686
sig_id:        PKCS#7
signer:        Build time autogenerated kernel key
sig_key:       30:9a:19:01:ba:9c:ba:d5:c0:8d:f7:a5:39:aa:c7:54:a6:c9:d8:2b
sig_hashalgo: sha512
signature:     30:64:02:30:6c:ab:da:07:56:cc:36:9d:66:06:e2:bb:98:e9:4a:50:
77:c0:37:08:0a:12:cd:5d:84:f7:2f:4a:fa:cb:5b:68:b9:c4:7b:c0:
08:1c:ec:61:33:fa:7e:a8:69:6b:fd:e7:02:30:69:c8:06:98:12:9c:
e3:b3:25:33:03:12:b1:d6:77:59:54:f5:0e:5b:d5:ff:c4:5d:d1:f1:
02:0e:16:68:2e:33:b4:97:2d:fd:be:35:1b:30:eb:17:aa:dd:01:ea:
93:0c
parm:          disable_cfc:Disable credit based flow control (bool)
parm:          channel_mtu:Default MTU for the RFCOMM channel (int)
parm:          l2cap_ertm:Use L2CAP ERTM mode for connection (bool)
```

⁷⁵ <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2215>

⁷⁶ <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2109>

⁷⁷ <https://stackoverflow.com/questions/57152408/what-is-the-internal-mechanics-of-socket-function>

⁷⁸ https://en.wikipedia.org/wiki/List_of_Bluetooth_protocols

⁷⁹ <https://linux.die.net/man/1/rfcomm>

⁸⁰ <https://www.bluetooth.com/specifications/specs/rfcomm-1-1/>

⁸¹ https://www.amd.e-technik.uni-rostock.de/ma/gol/lectures/wirlec/bluetooth_info/rfcomm.html

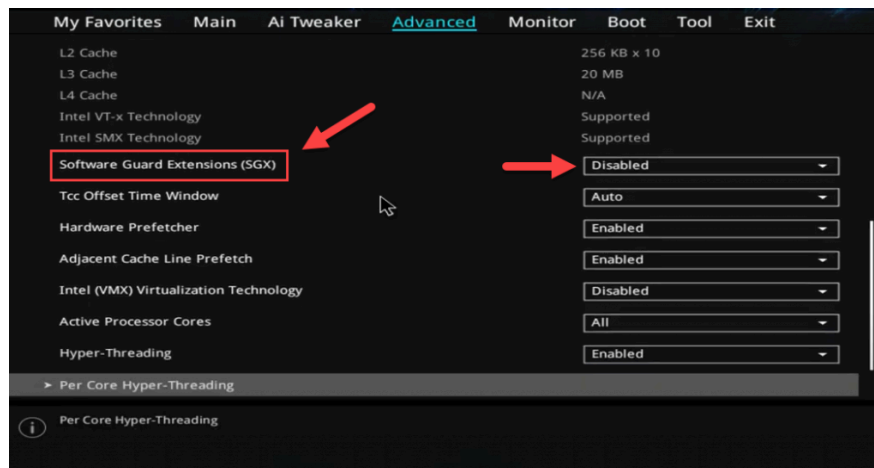
ksgxd (Kernel Software Guard eXtensions Daemon)

The kernel thread “ksgxd” is part of the Linux support for SGX (Software Guard eXtensions). Overall, SGX is a hardware security feature of Intel’s CPU that enables applications to allocate private memory regions for data and code. There is a privilege opcode “ENCLS” which allows creation of regions and “ENCLU” which is a privilege opcode that allows entering and executing code inside the regions⁸². For more information about SGX you can read my writeup about it⁸³.

“ksgxd” is a kernel which is started by executing “kthread_run()” function⁸⁴. The kernel thread executes the “ksgxd” function⁸⁵. “ksgxd” is started while SGX is initializing and at boot time it re-initializes all enclave pages. In case of over commitment “ksgxd” is also responsible for swapping enclave memory⁸⁶ like “kswapd”⁸⁷.

If you want to know if your CPU supports SGX you can use the following command: “cat /proc/cpuinfo | grep sgx” (you can also use lscpu). You can also check your UEFI (legacy BIOS) configuration to check if you - check out the screenshot below⁸⁸.

Lastly, there is a great guide for an example SGX app using a Linux VM on Azure that I encourage you to read⁸⁹. For more information about the Linux stack for SGX I suggest reading https://download.01.org/intelsgxstack/2021-12-08/Getting_Started.pdf and going over the following github repo <https://github.com/intel/linux-sgx>.



⁸² <https://docs.kernel.org/x86/sgx.html>
⁸³ <https://medium.com/@boutnaru/security-sgx-software-guard-extension-695cab7dbcb2>
⁸⁴ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L427>
⁸⁵ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L395>
⁸⁶ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L188>
⁸⁷ <https://medium.com/@boutnaru/the-linux-process-journey-kswapd-22754e783901>
⁸⁸ <https://phoenixnap.com/kb/intel-sgx>
⁸⁹ <https://tsmatz.wordpress.com/2022/05/17/confidential-computing-intel-sgx-enclave-getting-started/>

jbd2 (Journal Block Device 2)

“JBD” stands for “Journal Block Device”⁹⁰. “jbd2” is a kernel which is started by executing “kthread_run()” function⁹¹. The name of the kernel thread has the following pattern “jbd2/[DeviceName]”. The code is part of a kernel module - as you can see in the screenshot below.

Moreover, as we can see from the code it is a file system journal-writing code (part of the ext2fs journaling system). The journal is an area of reserved disk space used for logging transactional updates. The goal of “jbd2” is to schedule updates to that log⁹².

The kernel thread executes the “kjournald2()” function⁹³. This main thread function is used to manage a logging device journal. Overall, the thread has two main responsibilities: commit and checkpoint. Commit is writing all metadata buffers of the journal. Checkpoint means flushing old buffers in order to reuse an “unused section” of the log file⁹⁴.

Lastly, JBD was written by Stephen Tweedie and it is filesystem independent. There are different filesystems that are using it like ext3, ext4 and OCFS2. There are two versions: JBD created in 1998 with ext3 and JBD2 forked from JBD in 2006 with ext4⁹⁵.

```
root@localhost:~# modinfo jbd2
filename:       /lib/modules/5.19.7-arch1-1.0/kernel/fs/jbd2/jbd2.ko.zst
license:       GPL
srcversion:    7072394A13F8B3E5FCCE03C
depends:
retpoline:    Y
intree:       Y
name:         jbd2
vermagic:     5.19.7-arch1-1.0 SMP preempt mod_unload 686
sig_id:       PKCS#7
signer:       Build time autogenerated kernel key
sig_key:      30:9A:19:01:BA:9C:BA:D5:C0:8D:F7:A5:39:AA:C7:54:A6:C9:D8:2B
sig_hashalgo: sha512
signature:    30:64:02:30:0E:96:1E:1D:03:C4:F6:FD:71:26:C9:EC:8A:98:49:B8:
              91:E7:00:8A:90:43:6B:B9:D9:DD:F2:D0:64:27:8E:3B:4F:0A:CA:BD:
              3F:EC:76:4B:AD:26:79:0E:72:28:FC:C6:02:30:01:CA:42:28:FD:AA:
              D5:66:C5:16:05:2A:59:D5:BA:BE:4B:B4:DA:5E:DE:5F:1B:1B:01:06:
              7D:7B:59:12:58:D2:C5:5D:99:63:81:6B:60:D2:63:6C:0F:18:5A:26:
              9D:93
```

⁹⁰ <https://manpages.ubuntu.com/manpages/jammy/man1/pmdajbd2.1.html>

⁹¹ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L277>

⁹² <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c>

⁹³ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L169>

⁹⁴ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L152>

⁹⁵ https://en.wikipedia.org/wiki/Journaling_block_device

netns (Network Namespace)

The kernel thread “netns” is based on a single threaded workqueue⁹⁶, which is created when the network namespace is initialized (net_ns_init()). If you want to read more about “network namespaces” you can use the following link

<https://medium.com/@boutnaru/linux-namespaces-network-namespace-part-3-7f8f8e06fef3>.

Also, for a reminder you can also check out the diagram below⁹⁷.

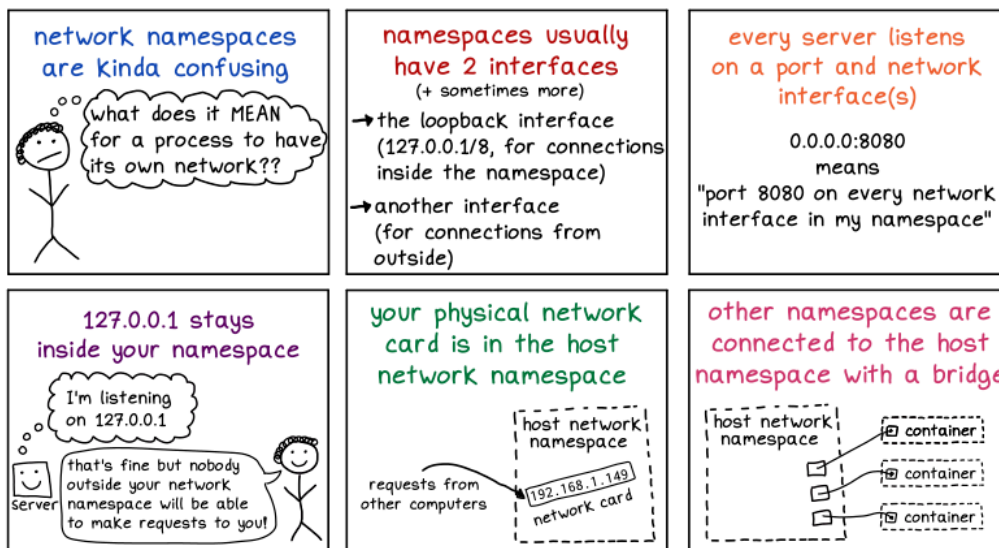
“netns” is responsible for cleaning up network namespaces. When a namespace is destroyed the kernel adds it to a cleanup list. The kernel thread “netns” goes over the list and performs the cleanup process using the “cleanup_net()” function⁹⁸.

If you want to see where all the magic happens is in “__put_net()” which queues the work on the “netns” to execute “cleanup_net()” function⁹⁹.

SULIA EVANS
@b0rk

network namespaces

18



⁹⁶ https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net_namespace.c#L1106

⁹⁷ <https://wizardzines.com/comics/network-namespaces/>

⁹⁸ https://elixir.bootlin.com/linux/v6.2.3/source/net/core/net_namespace.c#L565

⁹⁹ https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net_namespace.c#L649

oom_reaper (Out-of-Memory Reaper)

“oom_reaper” is a kernel thread which was created using the “kthread_run” function¹⁰⁰. Basically, it is the implementation of the OMM (Out-of-Memory) killer function of the Linux kernel - for more information about it I encourage you to read the following link <https://medium.com/@boutnaru/linux-out-of-memory-killer-oom-killer-bb2523da15fc>.

The function which is executed by the thread is “oom_reaper”¹⁰¹ which calls “oom_reap_task”¹⁰².

Based on the documentation the goal of the “oom_reaper” kernel thread is to try and reap the memory used by the OOM victim¹⁰³. “oom_reaper” sleeps until it is waked up¹⁰⁴ which is after OOM kills the process¹⁰⁵.

After killing the process the victim is queued so the “oom_reaper” can release the resources¹⁰⁶. You can see an example of the log created by OOM after killing a process¹⁰⁷.

```
4i7kwlazZ kernel: [ 1440] 0 1440 24585 250 0
4i7kwlazZ kernel: [ 1442] 0 1442 27085 108 0
4i7kwlazZ kernel: [ 1458] 0 1458 85095 8955 0
4i7kwlazZ kernel: Out of memory: Kill process 3223 (java) score
4i7kwlazZ kernel: Killed process 3223, UID 0, (java) total-vm:1
4i7kwlazZ yum[1458]: Updated: 7:squid-3.1.23-24.el6.x86_64
4i7kwlazZ squid[1511]: Squid Parent: child process 1513 started
```

¹⁰⁰ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L735

¹⁰¹ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L640

¹⁰² https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L609

¹⁰³ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L504

¹⁰⁴ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L680

¹⁰⁵ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L947

¹⁰⁶ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L992

¹⁰⁷ <https://blog.capdata.fr/index.php/linux-out-of-memory-killer-oom-killer-pour-un-serveur-base-de-donnees-postgresql/>

kpsmoused (Kernel PS/2 Mouse Daemon)

“kpsmoused” is a kernel thread which based on an ordered workqueue¹⁰⁸ which is allocated inside the “pmouse_init” function. “kpsmoused” is responsible for handling the input from PS/2 mouse devices.

Thus, “kpsmoused” transforms the raw data to high level event of mouse movements that be can consume from “/dev/input/mice”, “/dev/input/mouseX”, or “/dev/input/eventX”¹⁰⁹.

The kernel thread is created by the “psmouse” kernel module which is described as “PS/2 mouse driver” - as shown in the screenshot below (which was created using copy.sh). By the way, the “kpsmoused” is created as part of “/drivers/input/mouse/psmouse-base.c” since kernel 2.5.72¹¹⁰.

```
root@localhost:~# modinfo psmouse
filename:       /lib/modules/5.19.7-arch1-1.0/kernel/drivers/input/mouse/psmouse.ko.zst
license:       GPL
description:   PS/2 mouse driver
author:        UoJtech Paulik <uojtech@suse.cz>
srcversion:    900B3C75C18D7DDE1706120
alias:         serio:ty05pr*id*ex*
alias:         serio:ty01pr*id*ex*
depends:        libps2,serio
retpoline:     y
intree:        y
name:          psmouse
vermagic:      5.19.7-arch1-1.0 SMP preempt mod_unload 686
sig_id:        PKCS#7
signer:        Build time autogenerated kernel key
sig_key:       30:9a:19:01:ba:9c:ba:d5:c0:8d:f7:a5:39:aa:c7:54:a6:c9:d0:2b
sig_hashalgo: sha512
signature:     30:64:02:30:26:9e:10:64:df:8e:1f:2e:c6:2d:a8:f3:e1:53:a5:4a:
8f:01:f0:6b:f7:e7:a1:02:f4:6a:48:d0:43:fb:7b:3a:7e:0a:25:7b:
35:85:6d:cc:22:2c:b9:4a:93:2a:fc:b0:02:30:0b:da:0d:5e:c8:7d:
bd:96:ae:66:72:d4:4c:a7:64:59:e5:1d:76:0b:04:f2:20:70:93:d0:
23:ff:8b:9f:57:f5:d2:0b:9c:98:ea:37:08:d5:39:72:3b:50:87:7c:
13:c5
parm:          tpdebug:enable debugging, dumping packets to KERN_DEBUG. (bool)
parm:          recalib_delta:packets containing a delta this large will be discarded, and a recalibration may be scheduled. (int)
parm:          jumpy_delay:delay (ms) before recal after jumpiness detected (int)
parm:          spew_delay:delay (ms) before recal after packet spew detected (int)
parm:          recal_guard_time:interval (ms) during which recal will be restarted if packet received (int)
parm:          post_interrupt_delay:delay (ms) before recal after recal interrupt detected (int)
parm:          autorecal:enable recalibration in the driver (bool)
parm:          hggk_mode:default hggk mode: mouse, glidesensor or pentablet (string)
parm:          elantech_smbus:Use a secondary bus for the Elantech device. (int)
parm:          synaptics_intertouch:Use a secondary bus for the Synaptics device. (int)
parm:          proto:highest protocol extension to probe (bare, imps, exps, any). Useful for KVM switches. (proto_abbrev)
parm:          resolution:Resolution, in dpi. (uint)
parm:          rate:Report rate, in reports per second. (uint)
parm:          smartscroll:Logitech Smartscroll autorepeat, 1 = enabled (default), 0 = disabled. (bool)
parm:          a4tech_workaround:A4Tech second scroll wheel workaround, 1 = enabled, 0 = disabled (default). (bool)
parm:          resetafter:Reset device after so many bad packets (0 = never). (uint)
parm:          resync_time:How long can mouse stay idle before forcing resync (in seconds, 0 = never). (uint)
```

¹⁰⁸ <https://elixir.bootlin.com/linux/v6.2.6/source/drivers/input/mouse/psmouse-base.c#L2046>

¹⁰⁹ <https://www.kernel.org/doc/html/v5.5/input/input.html>

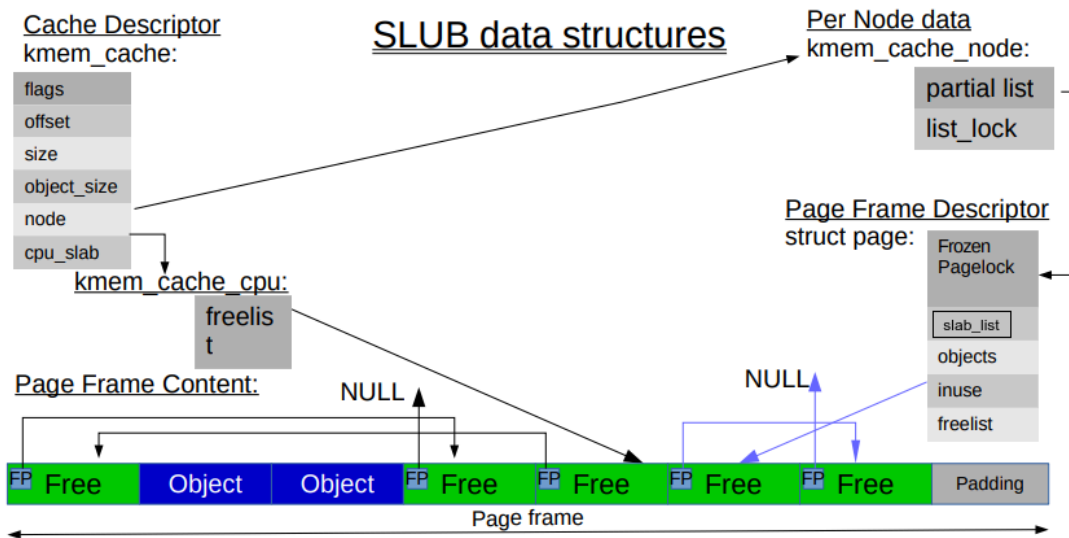
¹¹⁰ <https://elixir.bootlin.com/linux/v2.5.72/source/drivers/input/mouse/psmouse-base.c>

Slub_flushwq (SLUB Flush Work Queue)

“slub_flushwq” is a kernel thread which based on a workqueue¹¹¹ which is allocated inside the “kmem_cache_init_late” function. Based on the source code the allocation is done only if “CONFIG_SLUB_TINY” is enabled¹¹². From the documentation “CONFIG_SLUB_TINY” is for configuring SLUB allocation in order to achieve minimal memory footprint, it is not recommended for systems with more than 16 GB of RAM¹¹³. The queuing of work is done inside the “flush_all_cpus_locked” function¹¹⁴.

SLUB is also known as the “Unqueued Slab Allocator”¹¹⁵. Slab allocation is a memory management mechanism which allows efficient memory allocation of objects. It is done using reduction of fragmentation that is caused due to allocations/deallocations¹¹⁶. For more information about slab allocation I suggest reading the following link <https://hammertux.github.io/slab-allocator>.

Thus, SLUB is a slab allocator that limits the use of cache lines instead of using queued object per cpu/per node list¹¹⁷. So, it is less complicated because it does not keep queues (like for each CPU). The only queue is a linked list for all the objects in each of the slub pages¹¹⁸. The interplay between the three main data structures (kmem_cache, kmem_cache_cpu, kmem_cache_node) used by the SLUB allocator is shown in the diagram below¹¹⁹.



¹¹¹ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5057>
¹¹² <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5056>
¹¹³ https://cateee.net/lkddb/web-lkddb/SLUB_TINY.html
¹¹⁴ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L2822>
¹¹⁵ <https://lwn.net/Articles/229096/>
¹¹⁶ https://en.wikipedia.org/wiki/Slab_allocation
¹¹⁷ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L3>
¹¹⁸ <https://hammertux.github.io/slab-allocator>
¹¹⁹ <https://hammertux.github.io/img/SLUB-DS.png>

pgdatinit

“pgdatinit” is a kernel which is started by executing the “kthread_run()” function¹²⁰. The kernel thread executes the “deferred_init_memmap()” function¹²¹.

Thus, “pgdatinit” is responsible for initializing memory on every node of the system. For each node a dedicated kernel thread is created with the name pattern “pgdatinit[NodeNumber]”¹²².

Overall, the kernel thread is created in case CONFIG_DEFERRED_STRUCT_PAGE_INIT is enabled when compiling the kernel. Which states that initialization of struct pages is deferred to kernel threads¹²³.

Lastly, after the initialization flow is finished an information message is sent to the kernel ring buffer¹²⁴ - as you can see in the image below¹²⁵.

```
[ 0.212320] .... node #0, CPUs:          #1 #2 #3 #4 #5 #6 #7 #8 #9
#10 #11 #12 #13 #14 #15 #16 #17 #18 #19 #20 #21 #22 #23
[ 0.260348] smp: Brought up 1 node, 24 CPUs
[ 0.260348] smpboot: Max logical packages: 2
[ 0.260348] smpboot: Total of 24 processors activated (182404.32 BogoMIPS)
[ 0.357570] node 0 deferred pages initialised in 96ms
```

¹²⁰ https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2284

¹²¹ https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2108

¹²² https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2283

¹²³ https://cateee.net/lkddb/web-lkddb/DEFERRED_STRUCT_PAGE_INIT.html

¹²⁴ https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2177

¹²⁵ <https://www.mail-archive.com/debian-bugs-dist@lists.debian.org/msg1822096.html>

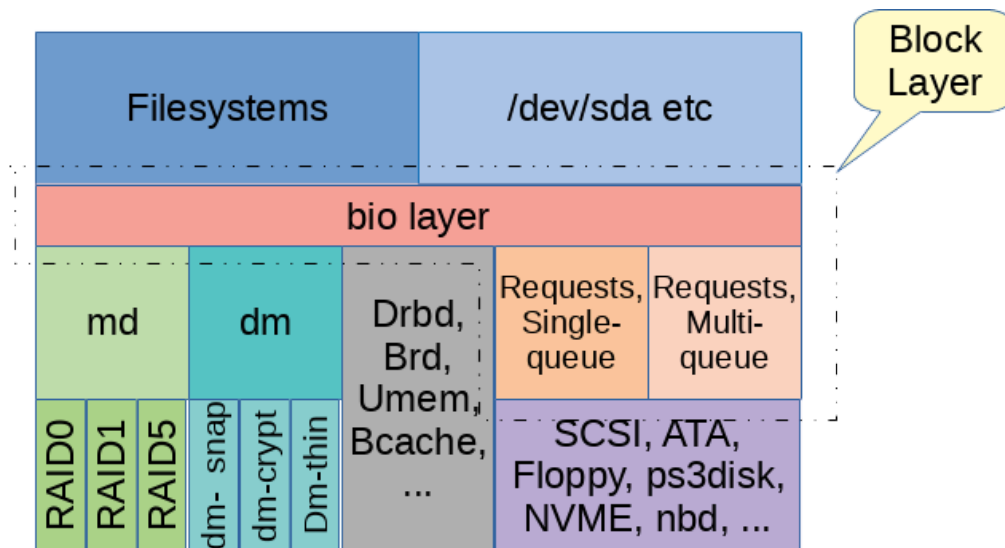
kblockd (Kernel Block Daemon)

“kblockd” is a kernel thread based on a workqueue¹²⁶ which is marked with high priority and that it can be used for memory reclaim. It is used for performing I/O disk operations.

Moreover, we can deduct based on the location of the file in the Linux source tree (/block) that “kblockd” is part of the “Block Layer” (which is responsible for managing block devices) - as shown in the diagram below¹²⁷.

Overall, one might think that we can use keventd¹²⁸ for performing I/O operations. However, because they can get blocked on disk I/O. Due to that, “kblockd” was created to run low-level disk operations like calling relevant block device drivers¹²⁹.

Thus, “kblockd” must never block on disk I/O so all the memory allocations should be GFP_NOIO. We can sum up that it is used to handle all read/writes requests to block devices¹³⁰.



¹²⁶ <https://elixir.bootlin.com/linux/v6.2.9/source/block/blk-core.c#L1191>

¹²⁷ <https://lwn.net/Articles/736534/>

¹²⁸ <https://lwn.net/Articles/11351/>

¹²⁹ <https://mirrors.edge.kernel.org/pub/linux/kernel/people/akpm/patches/2.5/2.5.70/2.5.70-mm8/broken-out/kblockd.patch>

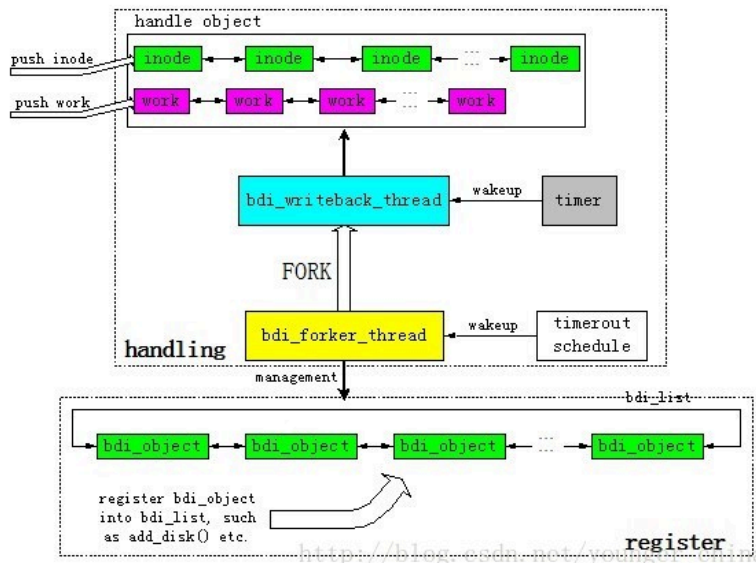
¹³⁰ <https://elixir.bootlin.com/linux/v6.3-rc4/source/block/blk-core.c#L13>

writeback

The kernel thread “writeback” is based on a workqueue¹³¹. The goal of the kernel thread is to serve all async writeback tasks¹³². Thus, “writeback” is flushing dirty information from the page cache (aka disk cache) to disks. The page cache is the main disk cache used by the kernel. The kernel references the page cache when reading from/writing to disk¹³³.

Overall, they are two ways of flushing dirty pages using writeback. The first is in case of an explicit writeback request - like syncing inode pages of a superblock. Thus, the “wb_start_writeback()” is called with the superblock information and the number of pages to flush. The second one is when there is no specific writeback request, in this case there is a timer that wakes up the thread periodically to flush dirty data¹³⁴.

Moreover, from kernel 3.2 the original mechanism of “pdflush” was changed to “bdi_writeback”. By doing so it solves one of the biggest limitations of “pdflush” in a multi-disk environment. In that case “pdflush” manages the buffer/page cache of all the disks which creates an IO bottleneck. On the other hand, “bdi_writeback” creates a thread for each disk¹³⁵. By the way, “bdi” stands for “Backing Device Information”¹³⁶. Lastly, to get an overview of the “writeback” mechanism you can checkout the diagram below¹³⁷.



¹³¹ <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L363>

¹³² <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L35>

¹³³ <https://www.oreilly.com/library/view/understanding-the-linux/0596005652/ch15s01.html>

¹³⁴ <https://lwn.net/Articles/326552/>

¹³⁵ https://blog.csdn.net/younger_china/article/details/55187057

¹³⁶ <https://lwn.net/Articles/326552/>

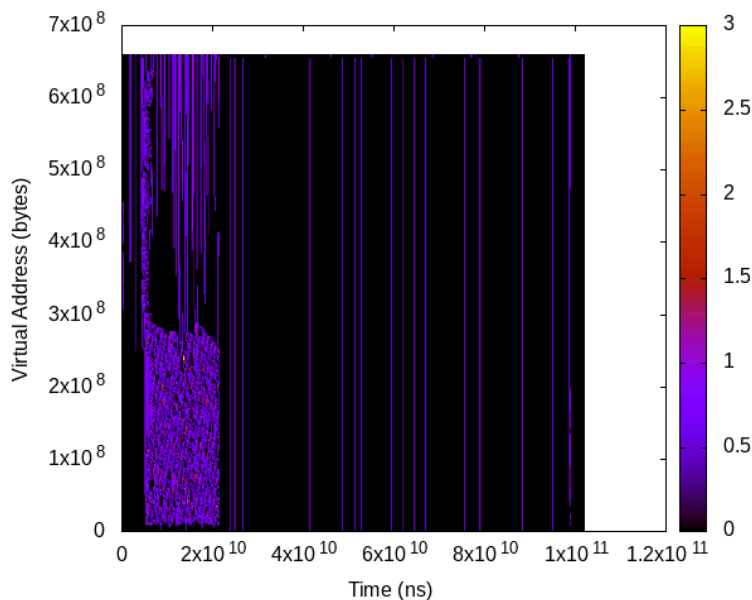
¹³⁷ https://blog.csdn.net/younger_china/article/details/55187057

kdamond (Data Access MONitor)

“kdamond” is a kernel thread which is created using the “kthread_run()” function¹³⁸ which is part of the DAMON (Data Access MONitor) subsystem. The kernel thread executes the “kdamon_fn()” function¹³⁹. Overall, DAMON provides a lightweight data access monitoring facility that can help users in analyzing the memory access patterns of their systems¹⁴⁰. Based on the documentation DAMON increases the memory usage by 0.12% and slows the workloads down by 1.39%¹⁴¹.

Also, DAMON has an API for kernel programs¹⁴². Moreover, there is also DAMOS (DAMon-Based Operations Schemas). Using that, users can develop and run access-aware memory management with no code and just using configurations¹⁴³.

Probably the best way to go over DAMON data is by using visualization. A great demonstration for that has been done by SeongJae Park using the PARSEC3/SPLASH-2X benchmarks¹⁴⁴. The output was heatmaps of the dynamic access patterns for heap area, mmap()ed area and the stack area. One example is shown in the image below, it visualizes the data access pattern of the stack area when running the parsec3-blackscholes¹⁴⁵. Lastly, there are also other mechanisms in Linux that can help with data access monitoring such as “Perf Mem” and “Idle Page Tracking”



¹³⁸ <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/damon/core.c#L632>

¹³⁹ <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/damon/core.c#L1304>

¹⁴⁰ <https://www.kernel.org/doc/html/latest/admin-guide/mm/damon/index.html>

¹⁴¹ <https://damonitor.github.io/doc/html/v20/vm/damon/eval.html>

¹⁴² <https://www.kernel.org/doc/html/v5.17/vm/damon/api.html#functions>

¹⁴³ <https://sjp38.github.io/post/damon/>

¹⁴⁴ <https://parsec.cs.princeton.edu/parsec3-doc.htm>

¹⁴⁵ <https://lwn.net/Articles/813108/>

kintegrityd (Kernel Integrity Daemon)

“kintegrityd” is a kernel thread based on a workqueue¹⁴⁶ which is responsible for verifying the integrity of block devices by reading/writing data from/to them. The function which is executed by the workqueue is “bio_integrity_verify_fn”¹⁴⁷. The function is called to complete a read request by verifying the transferred integrity metadata and then calls the original bio end_io function¹⁴⁸.

This procedure is done to ensure that the data was not changed by mistake (like in a case of a bug or an hardware failure¹⁴⁹). This mechanism is also called “bio data integrity extensions“. And it allows the user to get protection for the entire flow: from the application to storage device. The implementation is transparent to the application itself and it is part of the block layer¹⁵⁰.

Moreover, in order for it to work we should enable CONFIG_BLK_DEV_INTEGRITY, which is defined as “Block layer data integrity support”¹⁵¹. The filesystem does not have to be aware that the block device can include integrity metadata. The metadata is generated as part of the block layer when calling the submit_bio() function¹⁵². We can toggle the writing of metadata using “/sys/block/<BlockDevice>/integrity/write_generate“ and the verification of the metadata using “/sys/block/<BlockDevice>/integrity/read_verify” - as shown in the screenshot below.

Lastly, there are also file systems which are integrity aware (and they will generate/verify the metadata). There are also options for sending the metadata information from userspace, for more information I suggest reading the following Linux’s kernel documentation <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>.

```
Troller $ ls
device_is_integrity_capable  format  protection_interval_bytes  read_verify  tag_size  write_generate
Troller $ █
```

¹⁴⁶ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L455>

¹⁴⁷ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L317>

¹⁴⁸ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L313>

¹⁴⁹ <https://www.quora.com/What-is-the-purpose-of-kintegrityd-Linux-Kernel-Daemon/answer/Liran-Ben-Haim>

¹⁵⁰ <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

¹⁵¹ <https://elixir.bootlin.com/linux/v6.1/source/block/Kconfig#L60>

¹⁵² <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

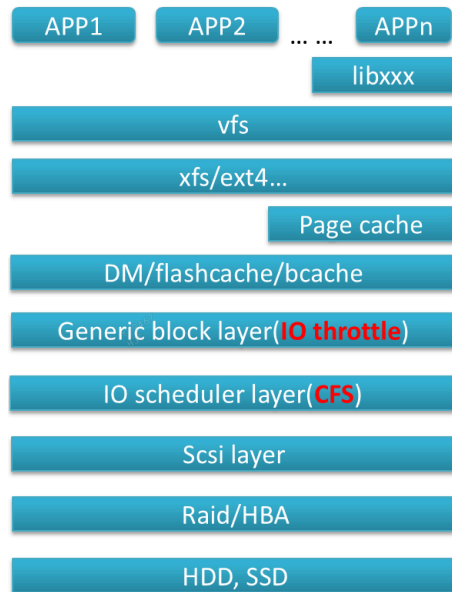
kthrotld (Kernel Throttling Daemon)

“kthrotld” is a kernel thread which was created using an workqueue¹⁵³ which acts as an interface for controlling IO bandwidth on request queues (throttling requests). Overall, read and write requests to block devices are placed on request queues¹⁵⁴.

In order to understand how request queues are used the best way is to check the source code of the kernel. The first step is going over the definition of “struct request_queue”¹⁵⁵ and then where is it referenced¹⁵⁶. By the way, in kernel version 6.1.1 it is referenced in 199 files. We can summarize that a request queue holds I/O requests in a linked list. Also, it is a best practice to create a separate request queue for every device¹⁵⁷.

Thus, we can say that “kthrotld” acts as a block throttle, which provides block QoS (Quality of Service). It is used to limit IOPS (I/O per second)/BPS (Bits per second) per cgroup (control group)¹⁵⁸.

Overall, IO throttling is done as part of the generic block layer and before the IO scheduler as seen in the diagram below¹⁵⁹. For more information on “Block Throttling” I suggest reading <https://developer.aliyun.com/article/789736>.



¹⁵³ <https://elixir.bootlin.com/linux/v6.1.1/source/block/blk-throttle.c#L2470>

¹⁵⁴ <https://www.halolinux.us/kernel-architecture/request-queues.html>

¹⁵⁵ <https://elixir.bootlin.com/linux/v6.1.1/source/include/linux/blkdev.h#L395>

¹⁵⁶ https://elixir.bootlin.com/linux/v6.1.1/C/ident/request_queue

¹⁵⁷ <https://www.oreilly.com/library/view/linux-device-drivers/0596000081/ch12s04.html>

¹⁵⁸ <https://developer.aliyun.com/article/789736>

¹⁵⁹ <https://blog.csdn.net/viveguzhou100/article/details/104044419>

scsi_eh (Small Computer System Interface Error Handling)

The kernel thread “scsi_eh” is executed using the “kthread_run” function. The name pattern of the kernel thread is “scsi_eh_<SCSI_HOST_NUMBER>”¹⁶⁰. It is the “SCSI error handler” which is responsible for all of the error handling targeting every SCSI host¹⁶¹. The kernel thread is executing the “scsi_error_handler” function¹⁶².

Moreover, a SCSI controller which coordinates between other devices on the SCSI bus is called a “host adapter”. It can be a card connected to a slot or part of the motherboard. You can see an example of a SCSI connector in the image below¹⁶³.

Lastly, SCSI stands for “Small Computer System Interface”. It is a set of standards (from ANSI) for electronic interfaces in order to communicate with peripheral hardware like CD-ROM drives, tape drives, printers, disk drives and more¹⁶⁴.. For more information about SCSI I suggest going over <https://hackaday.com/2023/03/02/scsi-the-disk-bus-for-everything/>.



¹⁶⁰ <https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/hosts.c#L504>

¹⁶¹ https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2230

¹⁶² https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2233

¹⁶³ <https://computer.howstuffworks.com/scsi.htm>

¹⁶⁴ <https://www.techtarget.com/searchstorage/definition/SCSI>

blkcg_punt_bio

“blkcg_punt_bio” is a kernel thread based on a workqueue. The workqueue itself is created in the “blkcg_init” function¹⁶⁵. It is part of the common block controller cgroup interface¹⁶⁶.

Overall, when a shared kernel thread tries to issue a synchronized block I/O (bio) request for a specific cgroup it can lead to a priority inversion. It can happen if the kernel thread is blocked waiting for that cgroup¹⁶⁷. An example of priority inversion is shown in the diagram below¹⁶⁸.

Thus, to avoid the problem mentioned above the function “submit_bio”¹⁶⁹ punts the issuing of the bio request to a dedicated work item (per-block cgroup).

It calls “blkcg_punt_bio_submit”¹⁷⁰, which will call “__blkcg_punt_bio_submit”¹⁷¹.

Priority inversion.

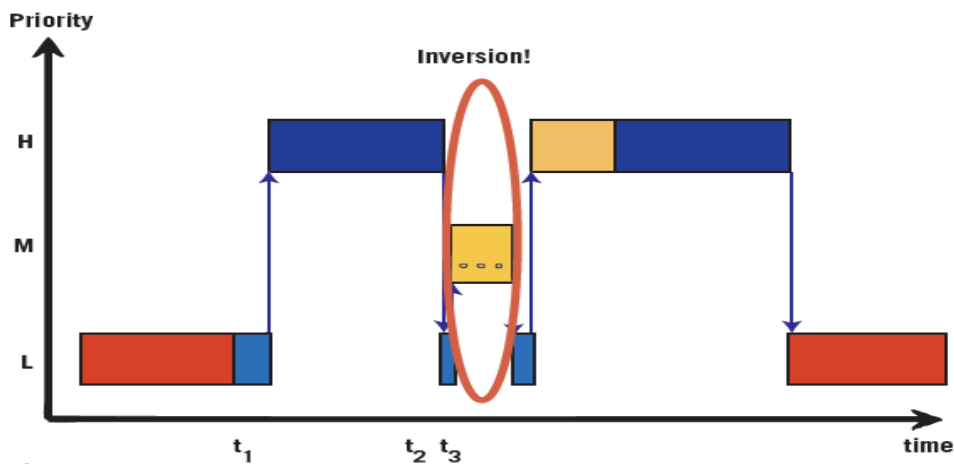


Figure 2

¹⁶⁵ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L2058>

¹⁶⁶ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L3>

¹⁶⁷ <https://patchwork.kernel.org/project/linux-block/patch/20190627203952.386785-6-tj@kernel.org/>

¹⁶⁸ <https://embeddedgurus.com/barr-code/2010/11/firmware-specific-bug-8-priority-inversion/>

¹⁶⁹ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-core.c#L829>

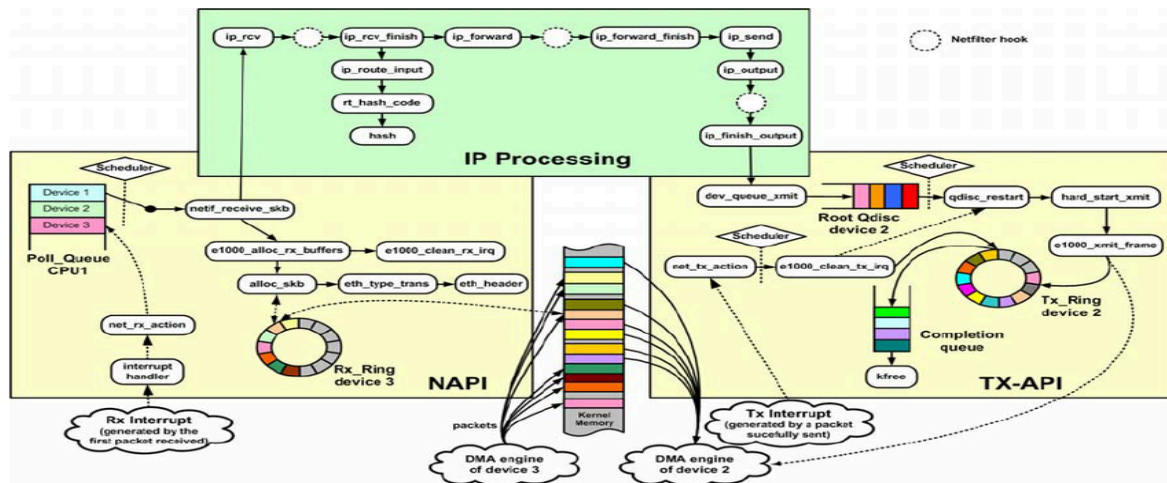
¹⁷⁰ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.h#L380>

¹⁷¹ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L1657>

napi (New API)

NAPI stands for “New API” which is used to reduce the number of received interrupts. Think about cases in which the network driver receives a large number of packets at a fast pace¹⁷². If we think about it in the case of a Gigabit network card and an MTU of 1500 the CPU will get about 90K of interrupt per second. Thus, we can say that NAPI is an extension to the Linux packet processing framework, which is done for improving performance for high speed networking. This is performed using interrupt mitigation and packet throttling. It is important to say that the addition of NAPI does not break backward compatibility¹⁷³. “napi” is a kernel thread which is created using the “kthread_run()”¹⁷⁴ function which is part of the NAPI (New API) subsystem. The name of the kernel thread is based on the pattern “napi[DeviceName]-[NAPI-ID]”. It executes the “napi_threaded_poll”¹⁷⁵ function.

Due to that, drivers that support NAPI can disable hardware interrupts as a mechanism for packet reception. In that case the network stack relies on polling for new packets at a specific interval. It might seem that polling is less efficient but in case the network device is busy any time the kernel will poll for a packet it will get something¹⁷⁶. Lastly, the way NAPI does that is by combining hardware interrupts and polling. When a hardware interrupt is received, the driver disables it and notifies the kernel to read the packets. Then a kernel software interrupt polls the network device for a specific time. When the time runs out/there is no more data the kernel will enable the hardware interrupt again¹⁷⁷. A detailed diagram of the NAPI flow is shown in the diagram below¹⁷⁸.



¹⁷² <https://www.hitchhikersguidetolearning.com/2023/04/09/handling-receive-packets-via-napi/>

¹⁷³ <https://wiki.linuxfoundation.org/networking/napi>

¹⁷⁴ <https://elixir.bootlin.com/linux/v6.4-rc4/source/net/core/dev.c#L1371>

¹⁷⁵ <https://elixir.bootlin.com/linux/v6.4-rc4/source/net/core/dev.c#L662>

¹⁷⁶ <https://lwn.net/Articles/833840/>

¹⁷⁷ <https://www.jianshu.com/p/7d4e36c0abe8>

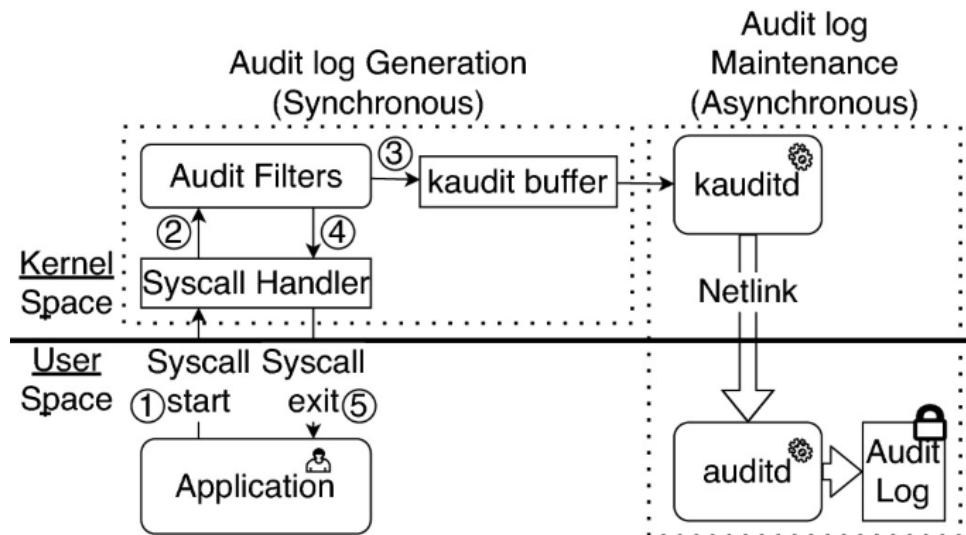
¹⁷⁸ <https://www.researchgate.net/profile/Roberto-Bruschi-2/publication/228624515/figure/fig4/AS:301797211164675@1448965470134/Detailed-scheme-of-the-forwarding-operations-in-26-kernel-NAPI.png>

kauditd (Kernel Audit Daemon)

“kauditd” is a kernel thread which is started using the “kthread_run” function¹⁷⁹. The kernel thread is calling the “kauditd_thread” function, this function is responsible for sending audit logs to userspace¹⁸⁰. Overall, the kernel mechanism in the Linux kernel has a couple of goals: integrate fully with LSMs¹⁸¹, minimal run-time overhead when performing auditing, ability to disable system call auditing at boot time, allow to be used by other parts of the kernel for auditing, netlink interface to userspace and support for filtering to minimize the information sent to user-mode¹⁸².

Thus, we can say “kauditd” is the kernel component of the “Linux Auditing System” which handles the audit events - as shown in the diagram below¹⁸³. In order to configure which set of rules are going to be loaded in the kernel audit system we can use the “/etc/audit/audit.rules” file. This file can hold configuration in one of three categories: control (configuring the audit system), file system rules monitoring rules and system call monitoring rules¹⁸⁴.

Lastly, by using the “Linux Auditing System” the system administrator can investigate what happens in the system for the purpose of debugging or in case of a security incident. We can also use the “auditctl” utility get/add/delete rules as part of Linux's kernel audit system¹⁸⁵. Also, there are great examples for “audit.rules” in GitHub (one example is <https://github.com/Neo23x0/auditd/blob/master/audit.rules>).



¹⁷⁹ <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L1700>

¹⁸⁰ <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L828>

¹⁸¹ <https://medium.com/@boutnaru/linux-security-lsm-linux-security-modules-907bbc8e8b4>

¹⁸² <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L11>

¹⁸³ https://link.springer.com/chapter/10.1007/978-3-031-17143-7_30

¹⁸⁴ <https://manpages.debian.org/unstable/auditd/audit.rules.7.en.html>

¹⁸⁵ <https://linux.die.net/man/8/auditctl>

tpm_dev_wq (Trusted Platform Module Device Work Queue)

“tpm_dev_wq” is a kernel thread base on a workqueue¹⁸⁶. It belongs a device file system interface for “Trusted Platform Module” aka TPM¹⁸⁷.

Overall, TPM is an international standard for secure cryptoprocessors. Those are microprocessors which are used for a variety of security applications such as secure boot, random number generating and crypto key storage¹⁸⁸.

Moreover, a work is queued for “tpm_dev_wq” as part of the function “tpm_common_write”¹⁸⁹. In case we are working in non-blocking mode an async job for sending the command is scheduled¹⁹⁰.

Lastly, “tpm-dev-common.c” is compiled as part of the kernel TPM device drivers as shown in the Makefile¹⁹¹. The information about the TPM module is shown in the screenshot below. I am using Ubuntu “22.04.2”, in which the TPM module is compiled directly into the kernel itself.

```
Troller $ cat /etc/issue
Ubuntu 22.04.2 LTS \n \l

Troller $ modinfo tpm
name:          tpm
filename:      (builtin)
license:       GPL
file:          drivers/char/tpm/tpm
version:       2.0
description:   TPM Driver
author:        Leendert van Doorn (leendert@watson.ibm.com)
parm:          suspend_pcr:PCR to use for dummy writes to facilitate flush on suspend. (uint)
```

¹⁸⁶ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L273>

¹⁸⁷ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L13>

¹⁸⁸ https://wiki.archlinux.org/title/Trusted_Platform_Module

¹⁸⁹ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L209>

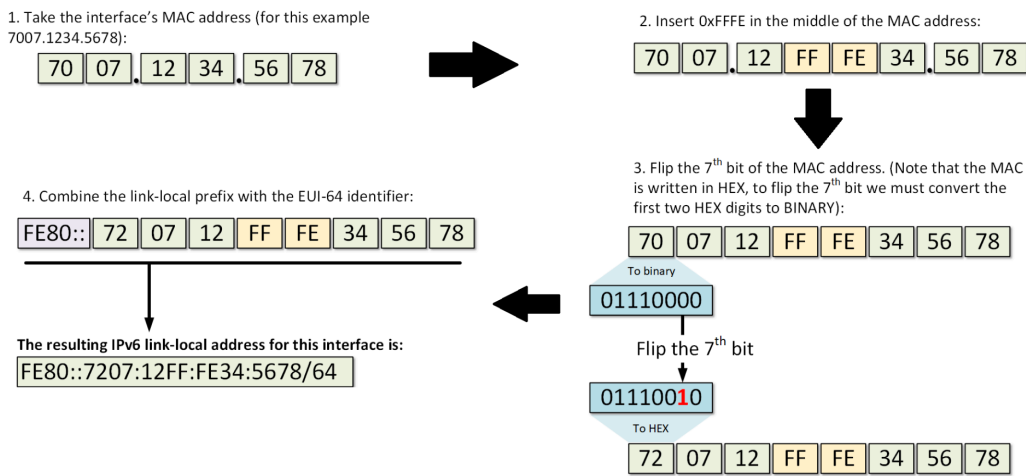
¹⁹⁰ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L202>

¹⁹¹ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/Makefile>

ipv6_addrconf (IPv6 Address Auto Configuration)

“ipv6_addrconf” is a kernel thread which is based on a workqueue¹⁹². This code is part of the Linux INET6 implementation and is responsible for the IPv6 Address auto configuration¹⁹³. Overall, each IPv6 entity in the network needs a globally unique address for communicating outside of the local segment. In order to get such an address there are a few options: manual assignment of an address, DHCPv6 (Dynamic Host Configuration Protocol version 6) and SLAAC (Stateless Address Autoconfiguration). When talking about stateless and stateful it means if there is a server/device that keep tracks of a state for each address assignment¹⁹⁴.

Moreover, the stateless address autoconfiguration has the following phases. The node configures itself with a link-local address. The most known way for doing that is using the link-local prefix “FE80::/64” and combining that with the EUI-64 identifier generated from the MAC address - as shown in the diagram below.



The flow above It is done by the function “addrconf_addr_gen”¹⁹⁵. We can see there the link-local prefix¹⁹⁶ and the call for generating the EUI-64 identifier by the function “ipv6_generate_eui64”¹⁹⁷. After that, the node performs DAD (Duplicate Address Detection) in order to ensure that the address is unique in the local segment. It is done using NDP (Neighbor Discovery Protocol), which defines 5 new packets types to ICMPv6 that allows to provide different functionality like DAD and others like parameter discovery, next hop determination and

¹⁹² <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L7292>

¹⁹³ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3>

¹⁹⁴ <https://www.networkacademy.io/ccna/ipv6/stateless-address-autoconfiguration-slaac>

¹⁹⁵ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3314>

¹⁹⁶ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3326>

¹⁹⁷ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3345>

more¹⁹⁸. If there are no issues with the link-local address it is assigned to the specific device. The DAD operation is performed by the function “`addrconf_dad_work`”¹⁹⁹.

Lastly, there is also a similar flow for configuring a global unicast address. The difference is that there is also a need for sending a “Router Solicitation” message for getting the global prefix of the segment, I will leave the details of that for a future writeup.

¹⁹⁸ <https://datatracker.ietf.org/doc/html/rfc4862>

¹⁹⁹ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L4058>

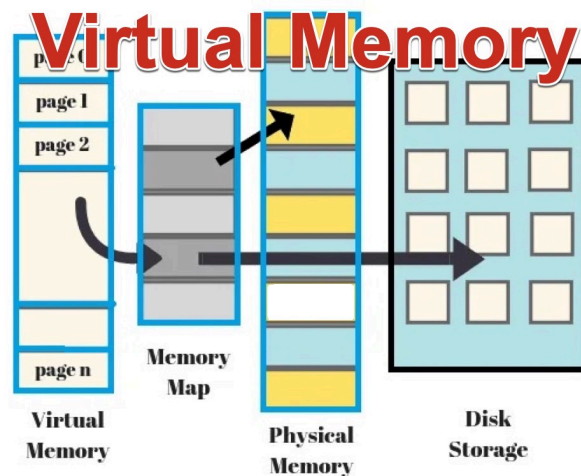
mm_percpu_wq (Per-CPU Memory Work Queue)

“mm_percpu_wq” is a kernel thread based on a workqueue which is created in the “init_mm_internals” function²⁰⁰. It is part of the the statistics management regarding virtual memory²⁰¹. An overview diagram of virtual memory is shown below²⁰².

Overall, “mm_percpu_wq” is the worker thread which updates different counters about the virtual memory of a Linux system. It is also called the “vmstat worker”²⁰³. “vmstat” stands for “Virtual Memory Statistics” which includes information such as: number of free pages, number of mapped pages, number or dirty pages, amount of memory allocated to kernel stacks and more (there are more than 150 different counters).

The statistics can be read from the file “/proc/vmstat”²⁰⁴. This proc entry is created with others (“buddyinfo”, “pagetypeinfo” and “zoneinfo”) in the same file in which “mm_percpu_mm” is allocated²⁰⁵. We can see the list of the metric counters in the source code²⁰⁶.

As it names suggested the kernel thread is responsible for accumulating the vm events among all CPUs²⁰⁷. It is done by going over all the “online” CPUs²⁰⁸. Lastly, we can use different cli tools to review the different statistic counters. One of those tools is “vmstat”²⁰⁹.



²⁰⁰ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2100>

²⁰¹ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L5>

²⁰² <https://iboysoft.com/wiki/virtual-memory.html>

²⁰³ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2021>

²⁰⁴ <https://man7.org/linux/man-pages/man5/proc.5.html>

²⁰⁵ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2123>

²⁰⁶ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L1168>

²⁰⁷ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L126>

²⁰⁸ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L117>

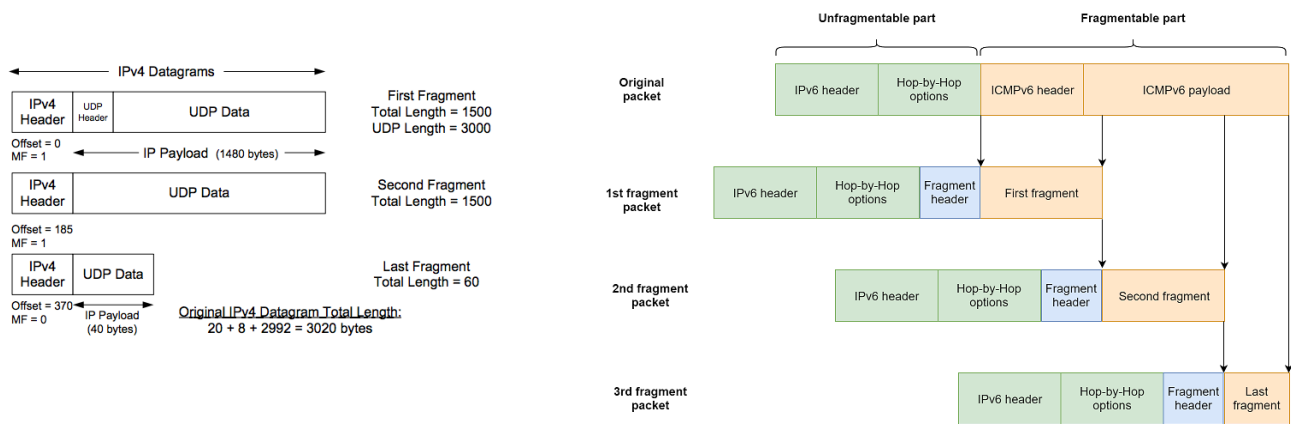
²⁰⁹ <https://linux.die.net/man/8/vmstat>

inet_frag_wq (IP Fragmentation Work Queue)

The kernel thread “inet_frag_wq” is created using a workqueue²¹⁰, we could have guessed it based on a workqueue do to the “wq” suffix. It is used for fragment management of IP packets. Thus the goal of “inet_frag_wq” is to reassemble fragmented IPv4/IPv6 packets²¹¹.

Overall, the goal of IP fragmentation is to split packets into smaller chunks in order to allow them to meet the MTU (Maximum Transmission Unit) requirement of a specific network. There is an implementation difference between IP fragmentation in IPv4 and IPv6. On IPv4 the information needed for fragmentation is part of the IPv4 header which in IPv6 there is a specific “Fragmentation Header”²¹². An illustration of the flow is shown in the diagram below both for IPv4²¹³ and IPv6²¹⁴.

Thus, “inet_frag_wq” is relevant when a fragmented IP packet arrives at a specific system. The OS stores the fragmented packets in a queue and reassembles them before they are passing the data to the upper layers of the network stack. The fragment queue is represented by "struct inet_frag_queue"²¹⁵. Moreover, we can see in the source code the function “ip_frag_reasm” which is responsible for building a new IP datagram from all of its fragments²¹⁶.



²¹⁰ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/inet_fragment.c#L211

²¹¹ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/inet_fragment.c#L6

²¹² <https://www.geeksforgeeks.org/ipv6-fragmentation-header/>

²¹³ <https://notes.shichao.io/tcpv1/ch10/>

²¹⁴ <https://blog.quarkslab.com/analysis-of-a-windows-ipv6-fragmentation-vulnerability-cve-2021-24086.html>

²¹⁵ https://elixir.bootlin.com/linux/v6.2-rc1/source/include/net/inet_frag.h#L66

²¹⁶ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/ip_fragment.c#L411

kstrp (Stream Parser)

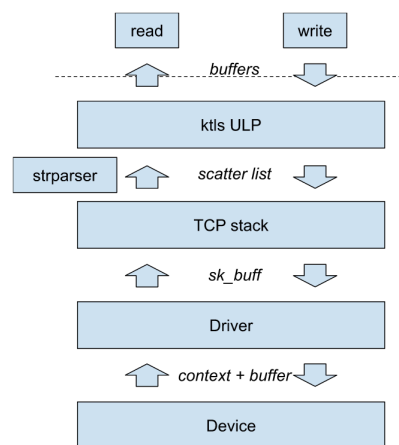
“kstrp” is based on a single threaded workqueue²¹⁷. From the source code documentation we can see that “strparser” means “Stream Parser”²¹⁸. A stream parser is a utility that gets data streams and parses the application layer protocol over those streams. A stream parser can work in one of two modes: general mode or receive callback mode.

In general mode, a sequence of socket buffers (skbs) are given to the stream parser from an outside source. Messages are parsed and delivered as the sequence is processed. This mode allows a stream parser to be applied to any arbitrary stream of data. In receive callback mode, the stream parser is called from the data_ready callback of the TCP socket. Messages are parsed and delivered as they are received on the socket²¹⁹.

Thus, we can say that we can parse application layer protocol messages in TCP. It is basically a generalization of KCM (Kernel Connection Multiplexor)²²⁰.

KMC provides a message based interface over TCP for generic application protocols. With the use of KMC applications can send/receive application messages efficiently over TCP²²¹.

Lastly, “strparser” allows intercepting packets on TCP connections. This is done at the kernel level which provides the ability to perform custom processing. The processing can be done using the BPF/Kernel module²²². One example for that is the implementation of KTLS²²³ (a Linux TLS/DTLS kernel module). An illustration of the flow is shown below²²⁴.



²¹⁷ <https://elixir.bootlin.com/linux/v6.1.1/source/net/strparser/strparser.c#L539>

²¹⁸ <https://elixir.bootlin.com/linux/v6.1.1/source/net/strparser/strparser.c#L3>

²¹⁹ <https://www.kernel.org/doc/html/v5.10/networking/strparser.html>

²²⁰ <https://lwn.net/Articles/695982/>

²²¹ <https://www.kernel.org/doc/html/latest/networking/kcm.html>

²²² <https://zhuanlan.zhihu.com/p/543663512>

²²³ https://github.com/ktls/af_ktls

²²⁴ <https://docs.kernel.org/networking/tls-offload.html>

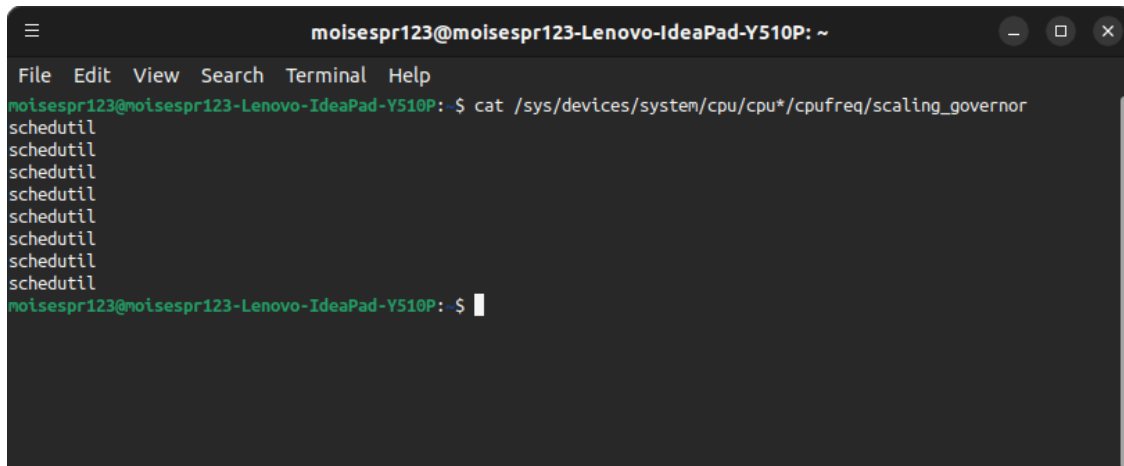
devfreq_wq

“devfreq_wq” is a kernel thread which is based on a freezable workqueue²²⁵. It is part of the Generic Dynamic Voltage and Frequency Scaling (DVFS) Framework for Non-CPU devices²²⁶.

Overall, DVFS enables Linux to scale the CPU frequency in order to minimize the power usage. It is mostly done when the full performance of the CPU is not needed. By using DVFS the system can set min/max CPU frequency. There is also the ability to set a “scaling governor” which monitors the performance requirements and decides what CPU frequency to use each time²²⁷.

Moreover, based on the Linux documentation there are 6 governors: “Performance”, “Powersave”, “Userspace”, “Ondemand”, “Conservative” and “Schedutil”²²⁸. We can also develop our own governor as a kernel module, we just need to register it using the function “cpufreq_register_governor”²²⁹.

Lastly, we can use the sysfs filesystem to configure/read information regarding “cpufreq”. An example of a file path for the first cpu is “/sys/devices/system/cpu/cpu0/cpufreq/” (if sysfs is mounted at “/sys”). It might contain the information like (but not limited to): current frequency of the CPU, the time it takes the CPU to switch frequencies (in nanosecs) and more²³⁰. An example of reading the current configure governor is shown below²³¹.



```
moisespr123@moisespr123-Lenovo-IdeaPad-Y510P: ~  
File Edit View Search Terminal Help  
moisespr123@moisespr123-Lenovo-IdeaPad-Y510P: $ cat /sys/devices/system/cpu/cpu*/cpufreq/scaling_governor  
schedutil  
schedutil  
schedutil  
schedutil  
schedutil  
schedutil  
schedutil  
schedutil  
schedutil  
moisespr123@moisespr123-Lenovo-IdeaPad-Y510P: ~ $
```

²²⁵ <https://elixir.bootlin.com/linux/v6.2.5/source/drivers/devfreq/devfreq.c#L1997>

²²⁶ <https://elixir.bootlin.com/linux/v6.2.5/source/drivers/devfreq/devfreq.c#L3>

²²⁷ https://wiki.somlabs.com/index.php/How_to_scale_CPU_frequency_with_DVFS_framework

²²⁸ <https://www.kernel.org/doc/Documentation/cpu-freq/governors.txt>

²²⁹ <https://elixir.bootlin.com/linux/v6.5-rc2/source/drivers/cpufreq/cpufreq.c#L2443>

²³⁰ <https://www.kernel.org/doc/Documentation/cpu-freq/user-guide.txt>

²³¹ <https://moisescardona.me/changing-the-cpu-governor-to-performance-in-linux/>

dmccrypt_write (Device Mapper for Transparent Encryption/Decryption)

“dmccrypt_write” is a kernel thread which is created using the “kthread_run” function²³². The name of the kernel thread is in the pattern of “dmccrypt_write/%s”, where the added string represents the device name. Overall, “dm-crypt” is a device-mapper target²³³ supported from kernel version 2.6.4²³⁴. It is responsible for transparent (aka real-time/on-the-fly encryption) block device encryption while using the kernel crypt API²³⁵.

This means the data is encrypted/decrypted while it is read/written. To enable the “dm-crypt” support we need to enable “CONFIG_DM_CRYPT” in the compilation config of the kernel²³⁶. Moreover, the function that is executed as part of the kernel thread is “dmccrypt_write” function²³⁷. This function is part of the kernel module “dm_crypt” - as shown in the screenshot below. We can use “modinfo dm_crypt” for more information, also shown in the screenshot below.

```
Troller $ modinfo dm_crypt | head -20
filename:      /lib/modules/5.15.0-78-generic/kernel/drivers/md/dm-crypt.ko
license:      GPL
description:  device-mapper target for transparent encryption / decryption
author:      Jana Saout <jana@saout.de>
srcversion:   FEC327FF4AB4CE3D2F1A54D
depends:
retpoline:    Y
intree:      Y
name:         dm_crypt
vermagic:     5.15.0-78-generic SMP mod_unload modversions
sig_id:       PKCS#7
signer:       Build time autogenerated kernel key
sig_key:      75:7A:05:56:12:13:0C:E4:F2:F6:B1:90:9C:50:42:33:83:2E:68:ED
sig_hashalgo: sha512
signature:    11:8A:EC:F9:98:EA:1E:5C:A0:81:E8:58:7F:0B:45:46:CB:FE:0F:CB:
48:90:65:7A:5C:45:11:84:C0:72:77:20:79:64:F5:EC:2F:CB:2C:69:
6D:C0:32:9D:42:32:00:DA:9F:4F:D6:F6:8C:E6:F2:DD:3B:A6:77:F0:
72:F9:2A:C6:92:33:15:33:7A:38:D4:E2:BF:FB:5D:78:11:50:7F:B5:
03:32:AF:FD:34:3B:D5:C5:24:12:DA:FC:6D:9A:49:90:F9:C6:5E:18:
32:55:E4:DD:3E:CB:14:9C:81:D7:44:96:05:F8:D6:CD:29:4D:23:4D:
Troller $ cat /proc/kallsyms | grep dmccrypt_write
0000000000000000 t dmccrypt_write [dm_crypt]
```

²³² <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L3388>

²³³ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L3689>

²³⁴ <https://elixir.bootlin.com/linux/v2.6.4/source/drivers/md/dm-crypt.c>

²³⁵ <https://gitlab.com/cryptsetup/cryptsetup/-/wikis/DMCrypt>

²³⁶ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/Makefile#L59>

²³⁷ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L1922>

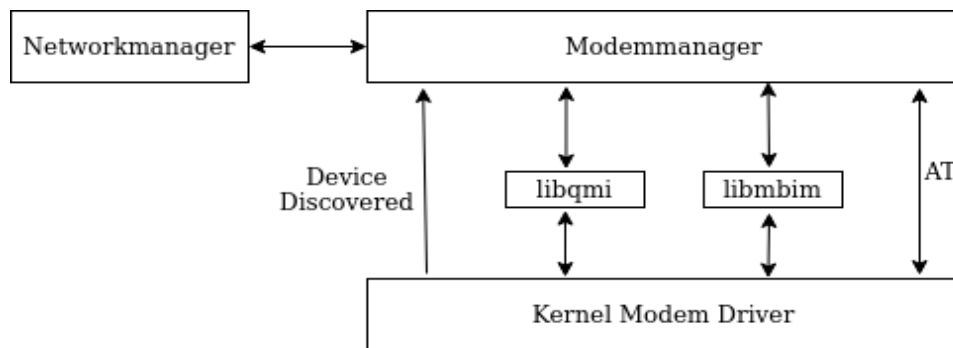
ModemManager (Modem Management Daemon)

“ModemManager” is an ELF binary located by default at “/usr/sbin/ModemManager” which is used to provide a unified high level API for communication with mobile broadband modems²³⁸. Alos, it is started by PID 1 (init/systemd) with the permission of the root user.

Overall, it is a Dbus-powered²³⁹ Linux daemon which acts as a standard RIL (Radio Interface Layer). “ModemManager” can be used by different connection managers (think about “NetworkManager” for example). Moreover, if we want to control and manage “ModemManager” we can use the CLI tool “mmcli”. By using it we can list all available modems, connect to a modem, get/set properties of the modem and more²⁴⁰.

Thus, we can summarize “ModemManager” as a system daemon that controls WWAN (2G/3G/4G/5G) devices and connections. It is the default mobile broadband management system in most Linux distributions (like Ubuntu, Debian, Fedora and Arch). By the way, it is also used by routers running OpenWRT²⁴¹.

It is important to understand that “ModemManager” leverages “libqmi”²⁴² and “libmbim”²⁴³ to communicate over QMI (Qualcomm MSM Interface) and MBIM (Mobile Interface Broadband Model) for setting connection to to the cellular network²⁴⁴. It does not matter if the modem is builtin, USB connected or bluetooth-paired. A diagram of the architecture is shown below. Lastly, if we want to go over the source code on “ModemManager” or contribute we can use its repo²⁴⁵. I also suggest going over the documentation site of “ModemManager” and the relevant libraries: libmbim, libqmi and libqrtr-glib²⁴⁶.



²³⁸ <https://manpages.ubuntu.com/manpages/trusty/man8/ModemManager.8.html>

²³⁹ <https://www.freedesktop.org/software/ModemManager/api/latest/>

²⁴⁰ <https://manpages.ubuntu.com/manpages/trusty/man8/mmcli.8.html>

²⁴¹ <https://modemmanager.org/>

²⁴² <https://github.com/linux-mobile-broadband/libqmi>

²⁴³ <https://github.com/linux-mobile-broadband/libmbim>

²⁴⁴ <https://developer.toradex.com/software/connectivity/modem-support/>

²⁴⁵ <https://gitlab.freedesktop.org/mobile-broadband/ModemManager>

²⁴⁶ <https://modemmanager.org/docs/>

kerneloops

“kerneloops” is an ELF binary located “/usr/sbin/kerneloops” (Ubuntu for example). It is used to collect kernel crash information (as part of a kernel oops) and submit them to kerneloops.org²⁴⁷. An example of such oops is shown in the screenshot below²⁴⁸. By the way, they are also known as “soft panic”²⁴⁹.

Overall, a kernel oops is a serious but non-fatal error in the Linux kernel. It is a way for the kernel to signal that it has found a problem that could potentially cause the system to crash. However, the kernel will continue to run after an oops, although it may be unstable and can lead to a kernel panic. This helps in debugging the error in order to find a solution for the problem²⁵⁰.

Moreover, if we want to debug the kernel with gdb it is suggested to compile it with “CONFIG_DEBUG_INFO” enabled, which causes the kernel to be built with full debugging information²⁵¹. Also, I recommend also enabling “CONFIG_FRAME_POINTER”, which gives very useful debugging information in case of kernel bugs - precise oopses/stacktraces/warnings²⁵².

Lastly, there is also a setting called “oops_limit” which states after what number of oops should cause a panic. The default value by the way is 10000²⁵³.

```
[12710.153112] oops init (level = 1)
[12710.153115] triggering oops via BUG()
[12710.153127] -----[ cut here ]-----
[12710.153128] kernel BUG at /home/duck/Articles/linuxoops/oops.c:17!
[12710.153132] invalid opcode: 0000 [#1] PREEMPT SMP PTI
[12710.153748] CPU: 0 PID: 5531 Comm: insmod
[12710.156191] RSP: 0018:ffffb41340e6fdd8 EFLAGS: 00010246
[12710.156849] RAX: 0000000000000019 RBX: ffffffff831015040 RCX: 0000000000000000
[12710.157513] RDX: 0000000000000000 RSI: ffffffff83bc9d39 RDI: 00000000ffffffff
[12710.158171] RBP: ffff8d6101bd1d50 R08: 0000000000000000 R09: fffffb41340e6fc90
[12710.158826] R10: 0000000000000003 R11: ffffffff83f3d1e8 R12: fffffb41340e6fde0
[12710.159483] R13: 0000000000000000 R14: 0000000000000000 R15: 0000000000000000
[12710.160143] FS: 00007f6c290b31c0(0000) GS:ffff8d6411a00000(0000) knlGS:000000
[12710.160820] CS: 0010 DS: 0000 ES: 0000 CR0: 0000000080050033
[12710.161478] CR2: 00000000004134f0 CR3: 000000018be34005 CR4: 0000000003706f0
[12710.162156] DR0: 0000000000000000 DR1: 0000000000000000 DR2: 0000000000000000
[12710.162824] DR3: 0000000000000000 DR6: 00000000fffe0ff0 DR7: 0000000000000400
[12710.163474] Call Trace:
[12710.164129] <TASK>
[12710.164779] do_one_initcall+0x56/0x230
[12710.165424] do_init_module+0x4a/0x210
[12710.166050] __do_sys_finit_module+0x9e/0xf0
[12710.166711] do_syscall_64+0x37/0x90
```

²⁴⁷ <https://linux.die.net/man/8/kerneloops>

²⁴⁸ <https://nakedsecurity.sophos.com/2023/03/13/linux-gets-double-quick-double-update-to-fix-kernel-oops/>

²⁴⁹ <https://www.opensourceforu.com/2011/01/understanding-a-kernel-oops/>

²⁵⁰ https://en.wikipedia.org/wiki/Linux_kernel_oops

²⁵¹ <https://www.oreilly.com/library/view/linux-device-drivers/0596005903/ch04.html>

²⁵² https://cateee.net/lkddb/web-lkddb/FRAME_POINTER.html

²⁵³ <https://docs.kernel.org/admin-guide/sysctl/kernel.html#oops-limit>

xargs (Extended Arguments)

“xargs” is an ELF file which is located at “/bin/xargs”, by the way “/bin” is usually a symbolic link to “/usr/bin” so we can find the binary at “/usr/bin/xargs”. It is used mainly for building commands based on the given standard input (the can be the standard output of a different command) and executing them²⁵⁴ - as shown in the screenshot below.

Overall, the default delimiter of xargs is blanks. We can override this behavior using “--delimiter” or “-d” switches. Also, “xargs” can read the input from a file using the “-a” or “--arg-file” switches. There is also the ability to limit the number of arguments (“-n” or “--max-args”), number of lines (“-l” or “--max-lines”) or the number of processes at a time (“-P” or “--max-procs”). If the value is more than the default 1 the command executed should handle race conditions, it is not done by xargs²⁵⁵.

Moreover, we can control the number of arguments at a to pass as input for “xargs” - as shown in the screenshot below. I suggest going over different xargs examples for better understanding its usage²⁵⁶. Lastly, there are different implementations for xargs two examples are Apple’s one²⁵⁷ and the busybox one²⁵⁸.

```
root@localhost:~# echo {1337..1444} | xargs -n 10
1337 1338 1339 1340 1341 1342 1343 1344 1345 1346
1347 1348 1349 1350 1351 1352 1353 1354 1355 1356
1357 1358 1359 1360 1361 1362 1363 1364 1365 1366
1367 1368 1369 1370 1371 1372 1373 1374 1375 1376
1377 1378 1379 1380 1381 1382 1383 1384 1385 1386
1387 1388 1389 1390 1391 1392 1393 1394 1395 1396
1397 1398 1399 1400 1401 1402 1403 1404 1405 1406
1407 1408 1409 1410 1411 1412 1413 1414 1415 1416
1417 1418 1419 1420 1421 1422 1423 1424 1425 1426
1427 1428 1429 1430 1431 1432 1433 1434 1435 1436
1437 1438 1439 1440 1441 1442 1443 1444
root@localhost:~# echo {1337..1444} | xargs -n 20
1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349 1350 1351 1352 1353 1354 1355 1356
1357 1358 1359 1360 1361 1362 1363 1364 1365 1366 1367 1368 1369 1370 1371 1372 1373 1374 1375 1376
1377 1378 1379 1380 1381 1382 1383 1384 1385 1386 1387 1388 1389 1390 1391 1392 1393 1394 1395 1396
1397 1398 1399 1400 1401 1402 1403 1404 1405 1406 1407 1408 1409 1410 1411 1412 1413 1414 1415 1416
1417 1418 1419 1420 1421 1422 1423 1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436
1437 1438 1439 1440 1441 1442 1443 1444
root@localhost:~# ls /bin/ba* | xargs -i echo "{}"
/bin/badblocks
/bin/base32
/bin/base64
/bin/basename
/bin/basenc
/bin/bash
/bin/bashbug
root@localhost:~#
```

²⁵⁴ <https://www.wikiwand.com/en/Xargs>

²⁵⁵ <https://man7.org/linux/man-pages/man1/xargs.1.html>

²⁵⁶ <https://www.tecmint.com/xargs-command-examples/>

²⁵⁷ https://opensource.apple.com/source/shell_cmds/shell_cmds-149/xargs/

²⁵⁸ <https://github.com/josefbacik/busybox/blob/master/findutils/xargs.c>

cpp (The C Preprocessor)

“cpp” is an ELF file which is the C preprocessor. It is mostly located at “/bin/cpp” or “/usr/bin/cpp” (where “/bin” is usually a symbolic link to “/usr/bin”). In some distributions like Ubuntu “/bin” is a symbolic link to “/usr/bin”. “cpp” is a macro processor that is used by the compiler in order to transform our code before compilation. In general, macros are abbreviations for longer constructs. We use it when writing code in C/C++/Objective-C²⁵⁹.

Moreover, the macros’ transformations can be the inclusion of header files, macro expansions and more. Header files are included using “#include” while defining a macro is done using “#define”. Also, we can use the preprocessor to instruct whether to include/or not a block of code as part of the compilation phase²⁶⁰. This is done using “ifdef”, “ifndef”, “define”, “else” and “endif”²⁶¹.

Lastly, we can say that the preprocessor prepares the source code for the compilation phase by performing different transformations - as shown in the screenshot below (we can see the replacement of the macro name with a string and the removal of a function).

```
root@localhost:/troller# cat troller.c
#define TROLLER_STRING "Tr0LL1Er"
#ifdef BLA
void troller_func()
{
    return 2222;
}
#endif
void main()
{
    char[]=TROLLER_STRING;
    return;
}
root@localhost:/troller# cpp troller.c
# 0 "troller.c"
# 0 "<built-in>"
# 0 "<command-line>"
# 1 "/usr/include/stdc-predef.h" 1 3 4
# 0 "<command-line>" 2
# 1 "troller.c"
# 12 "troller.c"
void main()
{
    char[]="Tr0LL1Er";
    return;
}
```

²⁵⁹ <https://linux.die.net/man/1/cpp>

²⁶⁰ <https://www.programiz.com/c-programming/c-preprocessor-macros>

²⁶¹ <https://www.cprogramming.com/reference/preprocessor/ifndef.html>

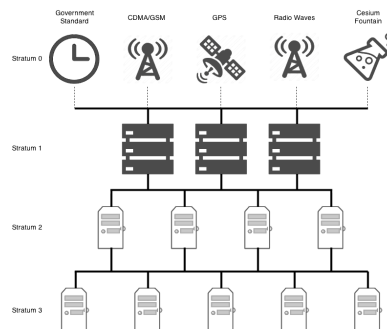
ntpd (Network Time Protocol Daemon)

“ntpd” is an ELF located by default at “/usr/bin/ntpd”. It is used to ensure that the system time is synchronized with standard time servers. It fully implements the “Network Time Protocol” (NTP) version 4, while also maintaining compatibility with older versions. “ntpd” performs most calculations using 64-bit floating-point arithmetic, only resorting to 64-bit fixed-point operations when necessary to maintain the highest possible precision, which is around 232 picoseconds²⁶².

Also, while such precision is not attainable with ordinary computers and networks today, it may be necessary in the future with gigahertz CPU clocks and gigabit LANs. “ntpd” reads its configuration from “/etc/ntp.conf” at startup. It uses the configuration in order to define the synchronization sources²⁶³. Overall, the “Network Time Protocol” was developed in 1981 by David Mills (professor at the University of Delaware). It was designed to be highly fault-tolerant/scalable, while supporting time synchronization. It is based on a Client/Server architecture using UDP with a default port number of 123²⁶⁴.

Moreover, there is an open source implementation of NTP from the University of Delaware²⁶⁵. You can also check out a port of it for the Windows operating system²⁶⁶. There are two utilities that can be used when “ntpd” is running for tasks like troubleshooting/configuration/monitoring: “ntpq”²⁶⁷ “ntpmon”²⁶⁸. Also, “Stratum” covers the accuracy of the time source.

Lastly, “Stratum 0” sources are the most accurate time sources, such as GPS, Cesium clocks, or cell networks. “Stratum 1” sources are systems that get their time from “Stratum 0” sources. “Stratum 2” sources get their time from “Stratum 1” sources, and so on. The lower the stratum number, the more accurate the time source. “Stratum 16” represents an unsynchronized clock, which is not reliable - as shown in the NTP stratum hierarchy diagram below²⁶⁹.



²⁶² <https://linux.die.net/man/8/ntpd>

²⁶³ <https://www.mankier.com/8/ntpd>

²⁶⁴ <https://www.techtarget.com/searchnetworking/definition/Network-Time-Protocol>

²⁶⁵ <http://www.ntp.org/>

²⁶⁶ <https://www.meinbergglobal.com/english/sw/ntp.htm>

²⁶⁷ <https://www.mankier.com/1/ntpq>

²⁶⁸ <https://www.mankier.com/1/ntpmon>

²⁶⁹ <https://chrisshort.net/ntp-i-need-you-to-go-ahead-and-love-it/>

gold (The GNU ELF Linker)

“gold” is a symbolic link to an ELF file which is the GNU ELF linker. As an example under Ubuntu it points to `x86_x64-linux-gnu-gold` which can be a symbolic link to `x86_x64-linux-ld.gold`. “gold” is mostly located at `/bin/gold` or `/usr/bin/gold`. In some distributions like Ubuntu `/bin` is a symbolic link to `/usr/bin`. We can also find “gold” located at `/usr/bin/ld.gold`. It was developed by Ian Lance Taylor and other members of Google, which measured it 5 times faster than the old GNU linker when linking C++ applications²⁷⁰.

Moreover, “gold” today is part of GNU binutils²⁷¹. As opposed to the GNU linker it does not use Binary File Descriptor library (BFD), which limits “gold” to process only ELF files but results in cleaner and faster implementations without the need for an abstraction layer²⁷².

Thus, we can specify “gold” as the the linker to use by setting LD as part of a makefile, setting the LD environment variable²⁷³ or by using the `“-fuse-ld=gold”` option of gcc²⁷⁴.

Lastly, we can also go over the paper that introduced gold titled “A new ELF Linker” by Ian Lance as part of the “2008 GCC Developers’ Summit”. The main difference between the GNU linker is that “gold” has a second walk over the input file to read the relocations, and the omission of the linker script - as shown in the image below taken from the paper about “gold”²⁷⁵. “gold” (sometimes also “ld.gold”) can take different command line arguments which can affect the linking processes²⁷⁶.

At a very high level, the GNU linker follows these steps:

- Walk over all the input files, identifying the symbols and the input sections.
- Walk over the linker script, assigning each input section to an output section based on its name.
- Walk over the linker script again, assigning addresses to the output sections.
- Walk over the input files again, copying input sections to the output file and applying relocations.



At a high level, `gold` follows these steps:

- Walk over the input files, identifying the symbols and the input sections.
- Walk over the input files again, reading the relocations and building the PLT and GOT.
- Assign output sections to output segments, and assign addresses to the output segments.
- Walk over the input files again, copying input sections to the output file and applying relocations.

²⁷⁰ <https://opensource.googleblog.com/2008/04/gold-google-releases-new-and-improved.html>

²⁷¹ <https://sourceware.org/binutils/>

²⁷² [https://en.wikipedia.org/wiki/Gold_\(linker\)](https://en.wikipedia.org/wiki/Gold_(linker))

²⁷³ [https://en.wikipedia.org/wiki/Gold_\(linker\)](https://en.wikipedia.org/wiki/Gold_(linker))

²⁷⁴ <https://man7.org/linux/man-pages/man1/gcc.1.html>

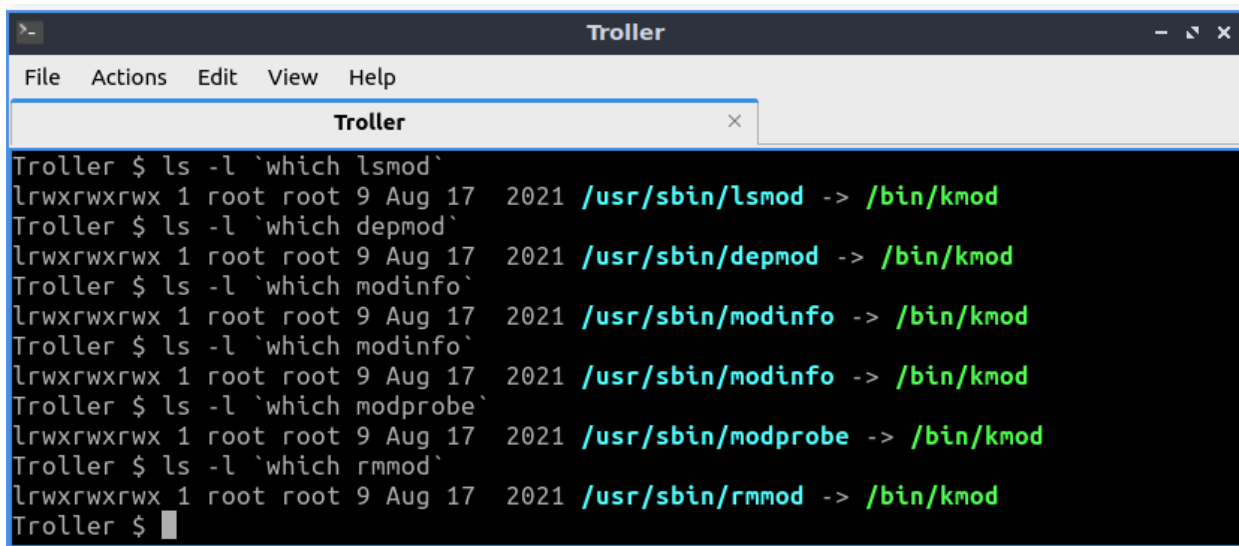
²⁷⁵ <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/34417.pdf>

²⁷⁶ <https://www.mankier.com/1/ld.gold>

kmod (Linux Kernel Module Handling)

“kmod” is an ELF file which is located at “/bin/kmod”, by the way “/bin” is usually a symbolic link to “/usr/bin” so we can find the binary at “/usr/bin/kmod”. It is used mainly for managing Linux kernel modules²⁷⁷. It is important to understand that “kmod” is a multi-call binary that implements utilities used to manage kernel modules. Hence, most users will not run it directly²⁷⁸.

Overall, by using “kmod” we can load/remove/insert/show information/resolve dependencies of kernel modules²⁷⁹. “kmod” is developed by Lucas De Marchi²⁸⁰. Lastly, the most well known utilities used to control kernel modules are just symbolic links to “kmod” examples are: lsmod, rmmmod, modinfo, modprobe and depmod - as shown in the screenshot below.



```
Troller $ ls -l `which` lsmod`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/lsmod -> /bin/kmod
Troller $ ls -l `which` depmod`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/depmod -> /bin/kmod
Troller $ ls -l `which` modinfo`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/modinfo -> /bin/kmod
Troller $ ls -l `which` modinfo`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/modinfo -> /bin/kmod
Troller $ ls -l `which` modprobe`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/modprobe -> /bin/kmod
Troller $ ls -l `which` rmmmod`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/rmmmod -> /bin/kmod
Troller $
```

²⁷⁷ https://linux-kernel-labs.github.io/refs/heads/master/labs/kernel_modules.html

²⁷⁸ <https://man7.org/linux/man-pages/man8/kmod.8.html>

²⁷⁹ <https://linuxhint.com/linux-kmod-command/>

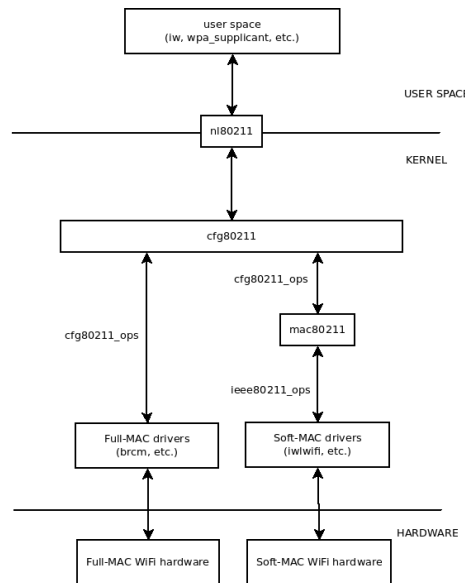
²⁸⁰ <https://github.com/kmod-project/kmod>

cfg80211 (Wireless Configuration)

“cfg80211” is a Linux kernel thread which is based on an ordered workqueue²⁸¹. It is defined as part of “/net/wireless/core.c” in the “cfg80211_init()” function, it is the part of the Linux wireless configuration interface²⁸². It interfaces with “nl80211”, together they are basically a replacement of wext aka “Wireless Extensions”²⁸³.

Overall, the “cfg80211” (subsystem for the Linux kernel regarding wireless configuration) was created by “Johannes Berg”, he also works on “mac80211” which is the Linux wireless stack²⁸⁴. “cfg80211” uses “struct cfg80211_ops” as a backend description for wireless configuration. This structure is registered by card drivers/wireless stacks for handling configuration requests on their interfaces²⁸⁵.

Moreover, “cfg80211” is a thin layer between userspace and drivers/mac80211 which includes mainly sanity checks and protocol translations. We can summarize the main flows of “cfg80211” as: device registration, regulatory enforcement, station management (AP), key management (AP only), mesh management, virtual interface management and scanning²⁸⁶. Lastly, “cfg80211” interfaces with “nl80211” (as shown in the below) using netlink and not ioctl²⁸⁷. By the way, the name “nl80211” represents 802.11 netlink interface²⁸⁸.



²⁸¹ <https://elixir.bootlin.com/linux/v6.7.5/source/net/wireless/core.c#L1719>

²⁸² <https://elixir.bootlin.com/linux/v6.7.5/source/net/wireless/core.c#L4>

²⁸³ https://elinux.org/images/7/75/DeRosier_WirelessInterfacing.pdf

²⁸⁴ <https://johannes.sipsolutions.net/Projects/>

²⁸⁵ <https://elixir.bootlin.com/linux/v6.7.5/source/include/net/cfg80211.h#L4501>

²⁸⁶ https://wireless.wiki.kernel.org/_media/en/developers/documentation/control.pdf

²⁸⁷ <https://stackoverflow.com/questions/21456235/how-do-the-nl80211-library-cfg80211-work>

²⁸⁸ <https://elixir.bootlin.com/linux/v6.7.5/source/include/uapi/linux/nl80211.h#L4>

kdmflush (Kernel Device Mapper Flush)

“kdmflush” is a kernel thread which is based on a workqueue²⁸⁹. Its name is created in the pattern of “kdmflush/%s” (where %s is the mapped device name). It is used by the “Device Mapper” (dm) in order to queue up deferred work to other context if doing them immediately so would be problematic²⁹⁰.

It is part of the “Device Mapper” framework and used in order to process deferred work that it has queued up from other contexts where doing immediately so would be problematic²⁹¹. Also, the kernel parameter “vm.dirty_background_ratio” controls the percentage of system memory that can be filled with “dirty” pages (memory pages not written to disk) before “kdmflush” kicks in to write them to disk²⁹².

Lastly, there are also other kernel parameters which affect the handling of “dirty” pages, we can check them out using the “sysctl” utility²⁹³ - as shown in the screenshot below. By the way, “kdmflush” is created as part of the “alloc_dev” function that is responsible for allocating and initializing a blank device with a given minor²⁹⁴.

```
root@localhost:~# sysctl -a | grep dirty_
vm.dirty_background_bytes = 0
vm.dirty_background_ratio = 10
vm.dirty_bytes = 0
vm.dirty_expire_centisecs = 3000
vm.dirty_ratio = 20
vm.dirty_writeback_centisecs = 500
root@localhost:~#
```

²⁸⁹ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm.c#L2131>

²⁹⁰ <https://www.compuhoy.com/how-do-you-check-which-process-is-using-more-disk-in-linux/>

²⁹¹ <https://askubuntu.com/questions/986211/what-is-kdmflush>

²⁹² <https://www.cnblogs.com/cobbliu/articles/11792193.html>

²⁹³ <https://linux.die.net/man/8/sysctl>

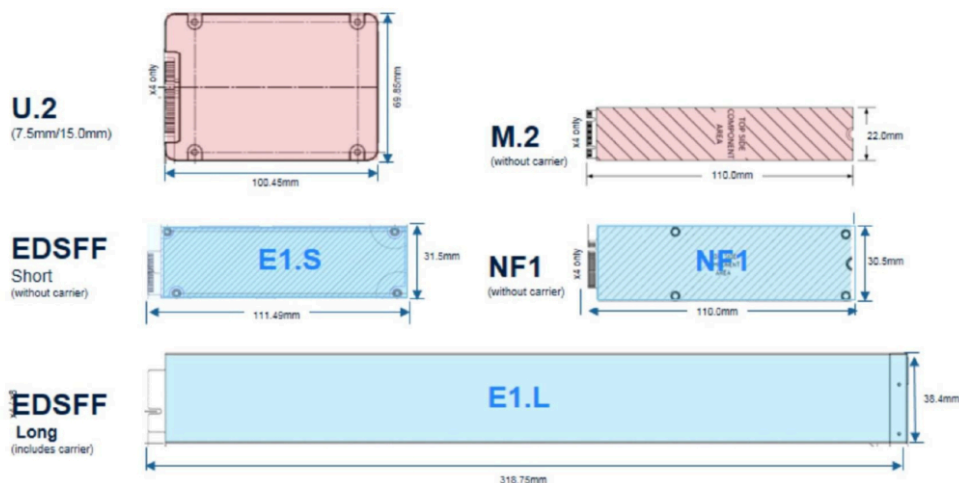
²⁹⁴ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm.c#L2044>

nvme-wq (Non-Volatile Memory Express Work Queue)

“nvme-wq” is a Linux kernel based on a workqueue, which is created as part of the “nvme_core_init” function in “/drivers/nvme/host/core.c”²⁹⁵. It is used for hosting NVMe related work which is not reset or delete²⁹⁶. Reset/delete are handled by other kernel threads which we are going to elaborate on in future writeups.

Overall, NVMe is a family of specifications that define how host software communicates with non-volatile memory across multiple transports (like PCI Express/RDMA/TCP). It is used as an industry standard for SSDs (solid state drives) in all form factors (U.2/M.2/AIC/EDSFF) - as shown in the diagram below²⁹⁷. By the way, NVM Express is a non-profit consortium²⁹⁸.

Lastly, “nvme-wq” is used to host work such as: scanning, AEN (Asynchronous Event Notifications) handling, FW activation, periodic reconnects and keep-alive²⁹⁹. Examples of those operations being queued from the Linux source are the following functions: “nvme_queue_scan”³⁰⁰, “nvme_enable_aen”³⁰¹, and “nvme_fw_act_work”³⁰².



²⁹⁵ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L4831>

²⁹⁶ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L93>

²⁹⁷ <https://venturebeat.com/business/meet-edsff-1pb-of-flash-storage-in-a-single-rack/>

²⁹⁸ <https://nvmexpress.org/>

²⁹⁹ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L197>

³⁰⁰ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L136>

³⁰¹ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L1682>

³⁰² <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L4257>

] (Checking File Types and Comparing Values)

“[” is an ELF binary located at “/usr/bin/[“ (or “/bin/[“) - as shown in the screenshot below. It is an equivalent to the “test” utility³⁰³. Thus, it is used for checking file types and comparing values³⁰⁴. The binary is part of the “coreutils” package³⁰⁵ - as shown in the screenshot below (taken from “copy.sh”, using the “Arch Linux”).

Overall, using “[” we can check expressions, compare integers, check the type of a file³⁰⁶, check if a file has been modified since it was last read, if the file exists and the user has read/execute access and more³⁰⁷. We can retrieve the result of the check using “\$?” - as shown in the screenshot below.

Lastly, in some shells like bash there is also an implementation of “[” as a builtin command³⁰⁸ - as shown in the screenshot below. Also, “[” is commonly used in scripts. For example, searching for it in scripts (using “grep.app”) yields more than 43K results³⁰⁹.

```
root@localhost:~# builtin [
-bash: [: missing `]'
root@localhost:~# builtin /usr/bin/[
-bash: builtin: /usr/bin/: not a shell builtin
root@localhost:~# file /usr/bin/[
/usr/bin/: ELF 32-bit LSB pie executable, Intel 80386, version 1 (SYSV), dynamically linked, interpreter /lib/ld-linux.so.2, BuildID[sha1]=, for GNU/Linux , stripped
root@localhost:~# pacman -Qo /usr/bin/[
/usr/bin/[ is owned by coreutils .0
root@localhost:~# [ -d /root ]
root@localhost:~# echo $?
0
root@localhost:~# [ -d /defugbhmfd ]
root@localhost:~# echo $?
1
```

³⁰³ <https://superuser.com/questions/334549/what-is-usr-bin-and-how-do-i-use-it>

³⁰⁴ <https://www.man7.org/linux/man-pages/man1/test.1.html>

³⁰⁵ <https://github.com/coreutils/coreutils/blob/master/src/test.c#L36>

³⁰⁶ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³⁰⁷ <https://linux.die.net/man/1/test>

³⁰⁸ <https://linux.die.net/man/1/builtins>

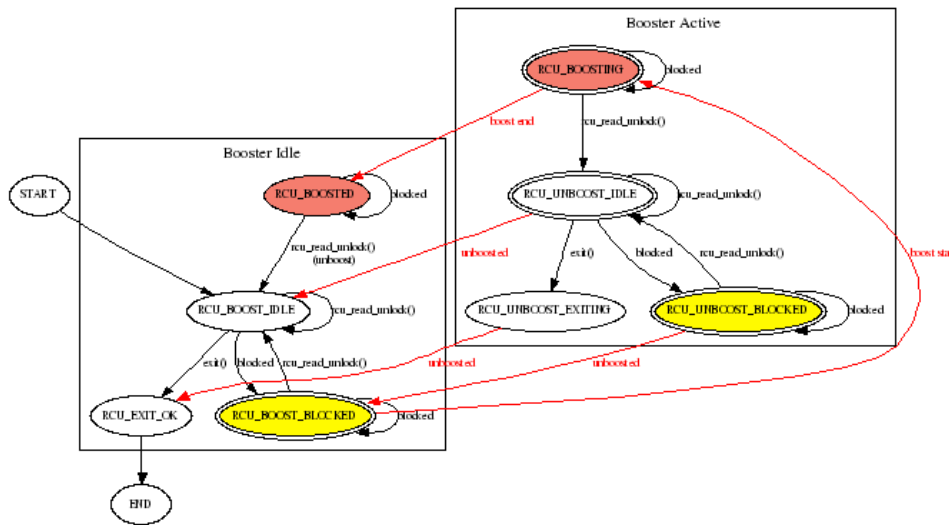
³⁰⁹ [https://grep.app/search?q=%5B%20&filter\[path\]\[0\]=scripts/](https://grep.app/search?q=%5B%20&filter[path][0]=scripts/)

rcub (Read-Copy Update Boost)

“rcub” is a Linux kernel thread created as part of the “rcu_spawn_one_boost_kthread” function³¹⁰. This function is used to create an RCU-boost kernel thread for a specific node. Due to that, the name of the kernel thread is in the format of “rcub/%d” (where %d is replaced with the node id). The kernel thread is created only for preemptible RCU³¹¹.

Overall, RCU (Read-Copy Update) is a mutual exclusion mechanism³¹². RCU is sometimes in place of read/write locks, due to the fact that based on properties of RCU updaters can block readers. Each task has its own RCU-booster state in order to avoid failure scenarios -as shown in the diagram below³¹³. We can enable RCU priority boosting using the “CONFIG_RCU_BOOST” configuration item when compiling the kernel³¹⁴.

Lastly, The “rcub” kernel thread executes the code of the “rcu_boost_kthread” function³¹⁵. We can summarize that “rcub” is used for boosting the priority of preempted RCU readers that block the current preemptible RCU grace period for too long. It is used for preventing heavy loads from blocking RCU callback invocation for all flavors of RCU. Thus, it is relevant for working with real-time apps\heavy loads³¹⁶.



³¹⁰ https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1202
³¹¹ https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1188
³¹² <https://medium.com/@boutnaru/linux-rcu-read-copy-update-c2793ee7a4cd>
³¹³ <https://lwn.net/Articles/220677/>
³¹⁴ <https://lwn.net/Articles/777214/>
³¹⁵ https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1110
³¹⁶ https://cateee.net/lkddb/web-lkddb/RCU_BOOST.html

dmesg (Print/Control the Kernel Ring Buffer)

“dmesg” is an ELF binary located at “/usr/bin/dmesg“ (or “/bin/dmesg“). The Linux kernel writes different messages to a kernel ring buffer (during the boot time and while the system is running). Thus, the kernel ring buffer holds various log messages, the buffer is fixed in size so in case it is full older entries are overwritten. There are different log facilities that can write to the kernel ring buffer like (but not limited to): “kern” (kernel messages), “user” (user level message), “mail” (mail system), “daemon” (system daemons), “auth” (security/authorization messages), “syslog” - as shown in the screenshot below³¹⁷.

Overall, whether we can execute “dmesg” without root permissions depends on the “kernel.dmesg_restrict” sysctl key. The “0” value indicates there are no restrictions, while “1” restricts access to those users that have the CAP_SYSLOG capability. Also, we can use the “CONFIG_SECURITY_DMESG_RESTRICT” kernel config option for setting the default value while compiling the kernel³¹⁸.

Moreover, “dmesg” is based on reading information for the character device “/dev/kmsg”³¹⁹. The utility can get different arguments for clearing the ring buffer after printing it (“-c”/“--read-clear”), setting the level of logging messages, enabling human readable format (“-H”/“--human”), using JSON format (“-J”/“--json”), printing human-readable timestamps (“-T”/“--ctime”) and more³²⁰.

Lastly, “dmesg” is part of the “util-linux” package, we can go over the source code of the utility in repo of the package³²¹. We can also force “dmesg” to use the “syslog” syscall³²² to read kernel messages instead of using “/dev/kmsg”³²³. By the way, it does not mean we will get the same number of entries - as shown in the screenshot below.

```
Troller $ dmesg
dmesg: read kernel buffer failed: Operation not permitted
Troller $ sudo dmesg -H -S | wc -l #using syslog syscall
1636
Troller $ sudo dmesg -H | wc -l #reading /dev/kmsg
1745
Troller $ sudo dmesg -H -f kern | head -3
[ 14:15] Dynamic Preempt: voluntary
[ +0.000034] rcu: Preemptible hierarchical RCU implementation.
[ +0.000002] rcu: RCU restricting CPUs from NR_CPUS=8192 to nr_cpu_ids=
```

³¹⁷ <https://linuxize.com/post/dmesg-command-in-linux/>

³¹⁸ https://linuxsecurity.expert/kb/sysctl/kernel_dmesg_restrict/

³¹⁹ <https://www.kernel.org/doc/Documentation/ABI/testing/dev-kmsg>

³²⁰ <https://man7.org/linux/man-pages/man1/dmesg.1.html>

³²¹ <https://github.com/util-linux/util-linux/blob/master/sys-utils/dmesg.c>

³²² <https://man7.org/linux/man-pages/man2/syslog.2.html>

³²³ <https://github.com/util-linux/util-linux/blob/master/sys-utils/dmesg.c#L343>

login (Begin Session on The System)

Overall, “login” is an ELF binary located by default at “/usr/bin/login” which is used to begin a session on a Linux system. It is used when signing onto a Linux system, if no arguments are passed a prompt of a user is shown³²⁴ - as shown in the screenshot below³²⁵.

Overall, login is part of the “util-linux” package which is a standard package that is distributed by the “Linux Kernel Organization”. It is important to understand that there are also other executables in that package like: kill, more, renice, su and more³²⁶. We can check the source code of “login” as part of the “util-linux” github repository³²⁷.

Moreover, login is based on PAM (Package Authentication Modules) which provides a framework for system-wide user authentication. You can see that also using “ldd”³²⁸ which will show libpam.so and probably also libpam_misc.so³²⁹.

Lastly, there are different configuration files that affect the behavior of “login” (beside PAM conf files). Among those configuration files are: “/etc/login.def”, /etc/motd, /etc/passwd and /etc/nologin - more information about them in future writeups. Also, there are logging based files which are handled by “login” such as: /var/run/utmp, /var/log/wtmp and /var/log/lastlog - more on them in future writeups³³⁰.

```
CentOS Linux 7 (Core)
Kernel 3.10.0-123.el7.x86_64 on an x86_64

tecmint login: root
Password:
Last login: Wed Apr 12 10:28:43 from 192.168.56.1
[root@tecmint ~]#
[root@tecmint ~]# _
```

³²⁴ <https://man7.org/linux/man-pages/man1/login.1.html>

³²⁵ <https://www.tecmint.com/understanding-shell-initialization-files-and-user-profiles-linux/>

³²⁶ <https://en.wikipedia.org/wiki/Util-linux>

³²⁷ <https://github.com/util-linux/util-linux/blob/master/login-utils/login.c>

³²⁸ <https://medium.com/@boutnaru/linux-instrumentation-part-4-ldd-888502965a9b>

³²⁹ <https://medium.com/@boutnaru/the-linux-security-journey-pam-pluggable-authentication-module-388496a8785c>

³³⁰ <https://linux.die.net/man/1/login>

su (Substitute\Switch User)

“su” is an ELF binary located at “/usr/bin/su” (or “/bin/su”) and used for running a command with a substitute user\group identifier. “su” stands for “Substitute User” or “Switch User”³³¹. “su” is part of the “util-linux” package, we can go over the implementation of the command as part of the package’s Github repository³³².

Overall, when executing “su [USERNAME]” we are prompted for the password of [USERNAME]. In case we don’t provide any [USERNAME] by default we are prompted for the password of root³³³ - as shown in the screenshot below. By the way, if we are running as root we can change to a different user without providing its password - as shown in the screenshot below.

Lastly, there is a difference between “su” to “su -”. The first one retains the user’s environment variables, working directory, current user’s shell setting and the target user’s PATH variable is not updated. The second option “su -” resets the environment variables, changes the working directory to the target home folder, resets all shell’s settings and updates the PATH variable³³⁴. By the way, it is recommended to use “su --login” instead of “su -” as a shortcut to avoid side effects caused by mixing environments³³⁵.

```
root@localhost:~# id
uid=0(root) gid=0(root) groups=0(root)
root@localhost:~# su troller
[troller@localhost root]$ id
uid=1000(troller) gid=1000(troller) groups=1000(troller)
[troller@localhost root]$ su
Password:
[troller@localhost root]$ su root
Password:
```

³³¹ <https://linuxize.com/post/su-command-in-linux/>

³³² <https://github.com/util-linux/util-linux/blob/master/login-utils/su-common.c>

³³³ <https://www.tecmint.com/difference-between-su-and-su-commands-in-linux/>

³³⁴ <https://www.geeksforgeeks.org/difference-between-su-and-su-command-in-linux/>

³³⁵ <https://man7.org/linux/man-pages/man1/su.1.html>

kacpid (Kernel Advanced Configuration and Power Interface Daemon)

“kacpid” is a Linux kernel thread which is based on an workqueue³³⁶. It is defined as part of “drivers/acpi/osl.c” in the “acpi_os_initialize1()” function, as part of the ACPI support of Linux³³⁷.

Overall, ACPI (Advanced Configuration and Power Interface) is a standard developed by Toshiba, Microsoft and Intel. It is used for power and system management. Think about controlling the amount of power each device is given (by putting a device on standby/powering it off) - we can see the ACPI states in the table below³³⁸. Also, it can be used for checking battery levels, PCI IRQ routing, CPUs, NUMA domains, fan speeds, temperature sensors and more³³⁹.

Lastly, work is queued for the “kacpid” kernel as part of the “acpi_os_execute()” function³⁴⁰. This kernel thread is used in case of a GPE (General Purpose Event) handler is needed³⁴¹. We can think about a GPE as an interrupt. In which the hardware is informing the OS (using ACPI) that “something” has happened. Examples are plugging/unplugging your AC adapter, closing/opening the lid of your laptop³⁴².

Table 4.3 ACPI operating states

Global System State	Sleep State	Description
G0 Working	(S0)	The computer is fully functional. Software, such as the AutoSave function used with Microsoft products, can be optimized for performance or lower battery usage.
		Requires less power than the G0 state and has multiple sleeping states.
G1 Sleeping	(S1)	CPU is still powered, and unused devices are powered down. RAM is still being refreshed. Hard disks are not running.
	(S2)	CPU is not powered. RAM is still being refreshed. System is restored instantly upon user intervention.
	(S3)	Power supply output is reduced. RAM is still being refreshed. Some info in RAM is restored to CPU and cache.
	(S4)	Lowest-power sleep mode and takes the longest to come up. Info in RAM is saved to hard disk. Some manufacturers call this the hibernate state.
G2	(S5)	Also called soft off. Power consumption is almost zero. Requires the operating system to reboot. No information is saved anywhere.
G3		Also called off, or mechanical off. This is the only state where the computer can be disassembled. You must power on the computer to use it again.

³³⁶ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1665>

³³⁷ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi>

³³⁸ <https://www.slideserve.com/reece/chapter-4-disassembly-and-power>

³³⁹ <https://wiki.osdev.org/ACPI>

³⁴⁰ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1113>

³⁴¹ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1105>

³⁴² <https://askubuntu.com/questions/148726/what-is-an-acpi-gpe-storm>

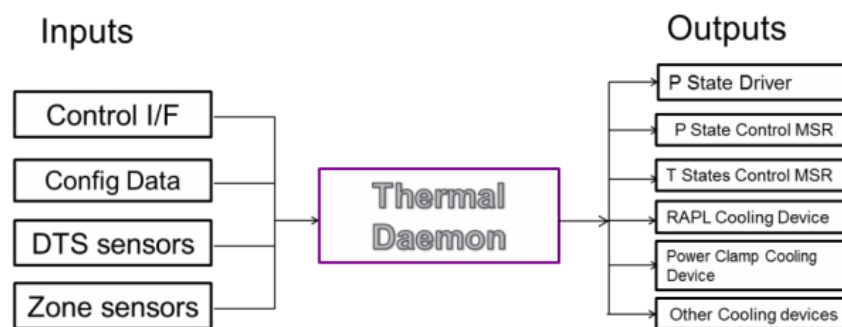
thermald (Thermal Daemon)

“thermald” is a Linux daemon which is responsible for controlling/monitoring the temperature of laptops/tablets/PCs containing the latest Intel CPU (at least Sandy Bridge). It is used for activating cooling methods when the system temperature reaches a specific temperature value³⁴³. We can go over the source code of the thermal daemon in GitHub³⁴⁴.

Overall, the goal is to prevent the BIOS from thermal throttling (by using proactive control) and to avoid the need of relying on the ACPI configuration. Thus, the “thermal daemon” provides a proactively user-mode controlled temperature solution that leverages existing kernel infrastructure and can be easily enhanced if needed. This is done by using sensors to calculate set a points and based on predefined configuration use the best cooling method³⁴⁵ - the general architecture is shown in the diagram below.

Lastly, “thermald” has two modes of operations: “Zero Mode Configuration” and “User Defined Configuration Mode”. By using the latest kernel drivers/modules like DTS temperature sensor (the output can be read at “/sys/devices/platform/coretemp.x interface) , Intel’s P state driver (performance state), RAPL (running average power limit) and the power clamp driver the “thermal daemon” allows developers to more precisely control the temperature set point for a given use case. It is important to know that not all are accessible to all hardware configurations. For example the “Intel’s P state” drivers are specific to “Sandy Bridge” and “Ivy Bridge” CPUs³⁴⁶.

Block Diagram



³⁴³ <https://wiki.debian.org/thermald>

³⁴⁴ https://github.com/intel/thermal_daemon

³⁴⁵ <https://web.archive.org/web/20211123224127/https://01.org/linux-thermal-daemon/documentation/introduction-thermal-daemon>

³⁴⁶ <https://www.linux.com/news/linux-thermal-daemon-monitors-and-controls-temperature-tablets-laptops/>

dhcpcd (Dynamic Host Configuration Protocol Daemon)

“dhcpcd” is an ELF binary located by default at “/usr/bin/dhcpcd” (or “/bin/dhcpcd”). It is an implementation of a DHCP (Dynamic Host Configuration Protocol) server for Linux. DHCP allows automatic configuration of network elements. Among that configuration can be included attributes such as: IP address, subnet mask and dns servers. The protocol operates on top UDP port 67 (server)/68 (client)³⁴⁷.

Overall, we can find the “dhcpcd” configuration in the “dhcpcd.conf” file located at “/etc”. Among the different configurations we can include: DHCP scope, subnets (including net masks) and subnet specific parameters like default gateway, lease timeout and more³⁴⁸. “dhcpcd” is the older ISC (Internet Systems Consortium) implementation of a DHCP server, which is not maintained since January 2023 and thus marked as “end of life”³⁴⁹. By the way, the ISC replacement is called “Kea” (more on that in a future writeup).

Lastly, it is important to understand that “dhcpcd” and “dhcpcd” are not the same. The first is a server implementation while the second is a client implementation of DHCP. Also, “dhcpcd” includes both an implementation for DHCPv4 and DHCPv6³⁵⁰. More information about “dhcpcd” can be found as part of the ISC documentation³⁵¹. On some Linux distributions we can use “systemctl”³⁵² for performing different operations of the “dhcpcd” daemon like checking its status - as shown in the screenshot below³⁵³.

```
[shovon@linuxhint-dhcp-server-6c024 ~]# sudo systemctl status dhcpcd
● dhcpcd.service - DHCPv4 Server Daemon
   Loaded: loaded (/usr/lib/systemd/system/dhcpcd.service; disabled; vendor preset: disabled)
   Active: active (running) since Thu 2020-03-12 21:28:06 +06; 5s ago
     Docs: man:dhcpcd(8)
           man:dhcpcd.conf(5)
  Main PID: 2878 (dhcpcd)
    Status: "Dispatching packets..."
     Tasks: 1 (Limit: 5935)
    Memory: 8.6M
   CGroup: /system.slice/dhcpcd.service
           └─2878 /usr/sbin/dhcpcd -f -cf /etc/dhcp/dhcpcd.conf -user dhcpcd -group dhcpcd --no-pid

Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]:
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: No subnet declaration for ens160 (192.168.20.175).
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: ** Ignoring requests on ens160. If this is not what
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: you want, please write a subnet declaration
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: in your dhcpcd.conf file for the network segment
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: to which interface ens160 is attached. **
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]:
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: Sending on Socket/fallback/fallback-net
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 dhcpcd[2878]: Server starting service.
Mar 12 21:28:06 linuxhint-dhcp-server-6c024 systemd[1]: Started DHCPv4 Server Daemon.
[shovon@linuxhint-dhcp-server-6c024 ~]#
```

³⁴⁷ <https://www.spiceworks.com/tech/networking/articles/what-is-dhcp/>

³⁴⁸ <https://linux.die.net/man/5/dhcpcd.conf>

³⁴⁹ <https://www.isc.org/blogs/isc-dhcp-eol/>

³⁵⁰ <https://wiki.archlinux.org/title/dhcpcd>

³⁵¹ <https://kb.isc.org/docs/aa-00333>

³⁵² <https://www.man7.org/linux/man-pages/man1/systemctl.1.html>

³⁵³ https://linuxhint.com/dhcp_server_centos8/